



- (51) International Patent Classification:
G06F 19/12 (2011.01)
- (21) International Application Number:
PCT/US2016/059233
- (22) International Filing Date:
27 October 2016 (27.10.2016)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data:
62/247,729 28 October 2015 (28.10.2015) US
62/394,721 14 September 2016 (14.09.2016) US
62/395,273 15 September 2016 (15.09.2016) US

(71) Applicants: **THE BOARD INSTITUTE INC.** [US/US]; 415 Main Street, Cambridge, MA 02142 (US). **MAS-SACHUSETTS INSTITUTE OF TECHNOLOGY** [US/US]; 77 Massachusetts Ave., Cambridge, MA 02142 (US). **THE REGENTS OF THE UNIVERSITY OF CALIFORNIA** [US/US]; 111 Franklin Street, 12th Floor, Oakland, CA 94607 (US).

(72) Inventors; and

(71) Applicants: **REGEV, Aviv** [IL/US]; 15a Ellsworth Ave, Cambridge, MA 02139 (US). **ADAMSON, Brittany** [US/US]; 1111 Franklin Street, 12th Floor, Oakland, CA 94607 (US). **CLEARY, Brian** [US/US]; 42 Curtis Ave., Somerville, MA 02144 (US). **CONG, Le** [CN/US]; 100 Memorial Drive, Apt. 8-21b, Cambridge, MA 02142 (US). **DIXIT, Atray** [US/US]; 203 Harvard Street, #2, Cambridge, MA 02139 (US). **GAUBLomme, Jellert** [BE/US]; 36 Oxford Street, Gsas Mailcenter Conant 104, Cambridge, MA 02138 (US). **LANDER, Eric, S.** [US/US]; 74r Fayerweather Street, Cambridge, MA 02138

(US). **NORMAN, Thomas** [CA/US]; 1111 Franklin Street, 12th Floor, Oakland, CA 94607 (US). **PAMAS, Oren** [IL/US]; 50 Green Street, Brookline, MA 02448 (US). **ROZENBLATT-ROSEN, Orit** [US/US]; 6 Juniper Road, Sharon, MA 02067 (US). **SHALEK, Alexander, K.** [US/US]; 2408 Massachusetts Ave., Lexington, MA 02421 (US). **WEISSMAN, Jonathan** [US/US]; 111 Franklin Street, 12th Nfloor, Oakland, CA 94607 (US).

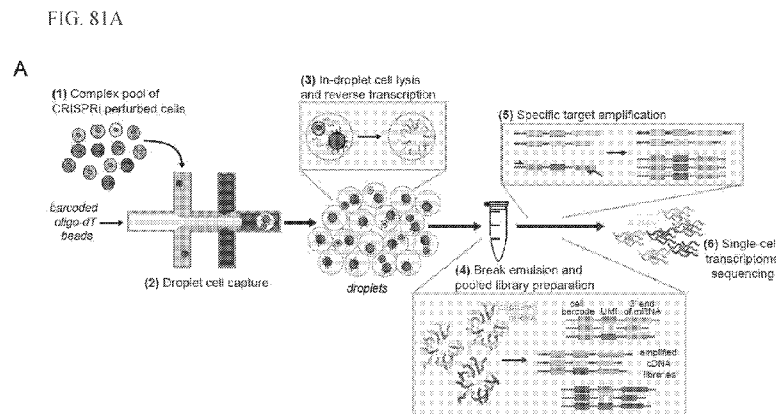
(74) Agents: **LU, Deborah, L.** et al.; Vedder Price P.C., 1633 Broadway, New York, NY 10019 (US).

(81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

(84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

[Continued on next page]

(54) Title: ASSAYS FOR MASSIVELY COMBINATORIAL PERTURBATION PROFILING AND CELLULAR CIRCUIT RE-CONSTRUCTION



(57) Abstract: The present invention provides tools and methods for the systematic analysis of genetic interactions, including higher order interactions. The present invention provides tools and methods for combinatorial probing of cellular circuits, for dissecting cellular circuitry, for delineating molecular pathways, and/or for identifying relevant targets for therapeutics development.

WO 2017/075294 A1

Published:

— *with international search report (Art. 21(3))*

**ASSAYS FOR MASSIVELY COMBINATORIAL PERTURBATION PROFILING AND
CELLULAR CIRCUIT RECONSTRUCTION****INCORPORATION BY REFERENCE**

[0001] This application claims priority and benefit of U.S. provisional application Serial numbers 62/247,729 filed October 28, 2015, 62/394,721 filed September 14, 2016, and 62/395,273 filed September 15, 2016.

[0002] Reference is also made to U.S. provisional application serial numbers 62/247,630 filed October 28, 2015, 62/247,656 filed October 28, 2015, and 62/372,393 filed August 9, 2016.

[0003] The foregoing applications, and all documents cited therein or during their prosecution (“appln cited documents”) and all documents cited or referenced in the appln cited documents, and all documents cited or referenced herein (“herein cited documents”), and all documents cited or referenced in herein cited documents, together with any manufacturer’s instructions, descriptions, product specifications, and product sheets for any products mentioned herein or in any document incorporated by reference herein, are hereby incorporated herein by reference, and may be employed in the practice of the invention. More specifically, all referenced documents are incorporated by reference to the same extent as if each individual document was specifically and individually indicated to be incorporated by reference.

FEDERAL FUNDING LEGEND

[0004] This invention was made with government support under grant numbers HG006193 MH105960, GM102706, CA168370, and DA036858 awarded by the National Institutes of Health. The government has certain rights in the invention.

FIELD OF THE INVENTION

[0005] The present invention provides methods and tools for analyzing genetic interactions. The present invention relies on a systematic approach of causing single and combinatorial genome-wide perturbations in cells, with subsequent molecular profiling at the single cell level. Applications include dissection of cell circuitry and delineation of functional or molecular pathways. The present invention is also relevant for therapeutics target discovery.

BACKGROUND OF THE INVENTION

[0006] Regulatory circuits in cells process signals, trigger decisions, and orchestrate physiological responses under diverse conditions. Diseases, in turn, arise from circuit malfunctions: one or more components are missing or defective; a key component is over- or under-active. To understand mechanisms underlying disease and develop more effective treatments, it would be highly advantageous to be able to provide a comprehensive picture of all cellular components, to identify the circuits in which they function, and to delineate how these components and circuits are integrated to form cellular responses.

[0007] Genomic research on dissecting cellular circuitry has generally had distinct phases: genomic observations and perturbation of single components.

[0008] Early advances in functional genomics made it possible to observe molecular profiles in different cells. Such global analysis has been very powerful in drawing hypotheses that relate regulators to their targets from statistical correlations. However, it is also very limited: the hypotheses were mostly not tested, and because correlation is not causation, many hypotheses may be found partially or fully incorrect.

[0009] In recent years, efforts were implemented in order to determine causation. Genomic profiles were used to infer a molecular model, on an increasingly large scale, based on genetic manipulations. However, the approach of testing genes individually has limitations: because genes involved in biological circuits have non-linear interactions, one cannot predict how a cellular circuit functions simply by summing up the individual effects. Indeed, biological systems are not linear: the combined effect of multiple factors is not simply the sum of their individual effects. This is a direct outcome of the biochemistry underlying molecular biology, from allosteric protein changes to cooperative binding, and is essential for cells to process complex signals.

[0010] It has remained an insurmountable stumbling block to achieving a quantitative and predictive understanding of circuits on a genomic scale, with far-reaching implications for basic and translational science. For example, despite decades of work, one still cannot predict how the enhancer controlling the transcription of the interferon beta gene (IFN β) behaves in response to viral and other stimuli. In another example, p38 α , a serine/threonine kinase with key roles in inflammation, has been studied for two decades, and yet it remains unclear how it balances control of inflammatory and anti-inflammatory cytokines, and many therapeutics programs

launched to target this protein were hampered by unwarranted and unexpected effects. Finally, in genomic studies ranging from yeast to mammals, many molecular events (e.g., transcription factor binding) appear 'functionally silent' upon factor perturbation, and only some expression variation is explained with available mechanistic data.

[0011] Citation or identification of any document in this application is not an admission that such document is available as prior art to the present invention.

SUMMARY OF THE INVENTION

[0012] A current phase of genomic research on dissecting cellular circuitry involves combinatorial probing of circuits. It would be desirable to provide a combinatorial approach: perturbing multiple components, at a large enough scale that will allow one to reliably reconstruct cellular circuits; for example, simultaneously or at or about the same time or in parallel. Such a combinatorial genomics approach has generally been considered intractable, because it required: (1) the ability to perturb many genes simultaneously (or at or about the same time or in parallel) in the same cell; (2) the ability to readout genomic profiles in individual cells, so that the effect of many perturbations can be assessed in parallel in a pool of cells; and (3) the development of mathematics and computational tools, because even millions of experiments are very few compared to the staggering size of the possible combinatorial space.

[0013] The invention involves Massively Parallel Combinatorial Perturbation Profiling (MCP) to address or identify genetic interactions. Biological systems are not linear: the combined effect of multiple factors is not simply the sum of their individual effects. This is a direct outcome of the biochemistry underlying molecular biology, from allosteric protein changes to cooperative binding, and is essential for cells to process complex signals. However, heretofore, it has remained an insurmountable stumbling block to quantitative and predictive biology on a genomic scale, with far-reaching implications e.g., from basic research to clinical translation. The invention provides a combinatorial approach: perturbing multiple components simultaneously. Yet, even recent experimental advances in genetics and genomics heretofore have not been enough; for instance, studying all 2-way, 3-way and 4-way interactions would, prior to the instant invention, require performing $\sim 10^8$, $\sim 10^{12}$, and $\sim 10^{16}$ combinations of perturbations, respectively; and, beyond the technological and cost issues, the sheer number of

mammalian cells required becomes infeasible because, for example, even if one focused on 100 genes, exhaustive analysis of 5-way interactions would require 70 million perturbations.

[0014] It would be desirable to provide tools and methods for the systematic analysis of genetic interactions, including higher order interactions.

[0015] It would be desirable to provide tools and/or methods for combinatorial probing of cellular circuits, for instance, for dissecting cellular circuitry, for delineating molecular pathways and/or interactions (e.g., intercellular and/or intracellular pathways or interactions), for identifying relevant targets and/or for identifying impact or effect of perturbations or stimuli or mutation; for instance for therapeutics development and/or cellular engineering and/or any cellular manipulation and/or ascertaining internal cell function and/or for bioproduction (e.g., production of antibodies from new sources, expression of products from organisms or cells that previously did not naturally express such products, increasing or decreasing expression of endogenous products, and the like), new plants or animal models.

[0016] The present invention involves cellular circuits (both intracellular and extracellular circuits). For instance, a cellular, e.g., regulatory circuit combines trans inputs (such as the levels and activities of factors, e.g., transcription factors, non-coding RNAs, e.g., regulatory RNAs and signalling molecules) and cis inputs (such as sequences, e.g., regulatory sequences in the promoter and enhancer of a gene); for instance, to determine the level of mRNA produced from a gene.

[0017] Reconstruction of a cellular, e.g., regulatory circuit is to identify inputs, e.g., all identifiable inputs (for example, proteins, non-coding RNAs and cis-regulatory elements), their physical 'wirings' (or connections) and the transcriptional functions that they implement; for instance, as to regulation of the level of mRNA.

[0018] A model should address (advantageously simultaneously or in parallel) providing a functional description of the input–output relationships (for example, if regulator A is induced, then target gene B is repressed to a particular extent), and providing a physical description of the circuit (for example, regulator A binds to the promoter of gene B in sequence Y, modifies its chromatin and leads to repression). Networks, e.g., regulatory networks, control complex downstream cellular phenotypes (such as cell death, proliferation and migration).

[0019] Reconstructing the connectivity of a network can be through the monitoring of hundreds to thousands of cellular parameters (massively parallel monitoring or hundreds to

thousands of cellular parameters), such as the levels of mRNAs. Hence “massively parallel” can mean undertaking a particular activity hundreds to thousands to millions, e.g., from 100 to 1000 or to 10,000 or to 100,000 or to 1,000,000 or up to 1,000,000,000 times (or as otherwise indicated herein or in figures herewith), in parallel, e.g., simultaneously or at or about the same time. *See, e.g.*, Amit et al., “Strategies to discover regulatory circuits of the mammalian immune system,” NATURE REVIEWS (IMMUNOLOGY) 11: 873-880 (DECEMBER 2011).

[0020] The present invention relates to methods of measuring or determining or inferring RNA levels, e.g., massively parallel measuring or determining or inferring of RNA levels in a single cell or a cellular network or circuit in response to at least one perturbation parameter or advantageously a plurality of perturbation parameters or massively parallel perturbation parameters involving sequencing DNA of a perturbed cell, whereby RNA level and optionally protein level may be determined in the single cell in response to the at least one perturbation parameter or advantageously a plurality of perturbation parameters or massively parallel perturbation parameters. The invention thus may involve a method of inferring or determining or measuring RNA in a single cell or a cellular network or circuit, e.g., massively parallel inferring or determining or measuring of RNA level in a single cell or a cellular network or circuit in response to at least 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99 or 100 or massively parallel perturbation parameter(s) comprising optionally so perturbing the cell or the cells or each cell of the cellular network or circuit with the perturbation parameter(s) and sequencing of the perturbed cell(s), whereby RNA level(s) and optionally protein level(s) is / are determined in the cell(s) in response to the perturbation parameter(s).

[0021] Genetic screens are used to infer gene function in mammalian cells, but it has remained difficult to assay complex phenotypes – such as genome-wide transcriptional profiles – in large-scale screens. Moreover, it has been traditionally difficult to assay the transcriptional phenotype of genetic perturbations at scale. Preferably, a genomewide scale transcriptome phenotype associated with a perturbation would be possible.

[0022] Here, Applicants developed Perturb-seq, which combines single cell RNA-seq and CRISPR/Cas9 based perturbations identified by unique polyadenylated barcodes to perform

many, tens of thousands in certain embodiments, of such assays in a single pooled experiment. By randomly integrating more than one sgRNA in each cell, Perturb-Seq is extended to test transcriptional phenotypes caused by genetic interactions. Applicants develop a computational framework, MIMOSCA (Multi-Input Multi-Output Single Cell Analysis) to identify the regulatory effects of individual perturbations and their combinations at different levels of resolution: from effects on each individual gene to functional signatures to proportional changes in cell types. Applicants demonstrate Perturb-seq by analyzing 200,000 cells across three screens: transcription factors controlling the immune response of dendritic cells to LPS, transcription factors bound in the K562 cell line, and cell cycle regulators in the same cell line. Perturb-Seq accurately identified known regulatory relations, and its individual gene target predictions were validated by ChIP-Seq binding profiles. Applicants posit new functions for regulatory factors affecting cell differentiation, the anti-viral response and mitochondrial function during immune activation, and uncovered an underlying circuit that balances these different programs through positive and negative feedback loops. Using Perturb-Seq Applicants identified genetic interactions including synergistic, buffering and dominant genetic interactions that could not be predicted from individual perturbations alone. Perturb-Seq can be flexibly applied to diverse cell metadata, to customize design and scope of pooled genomic assays.

[0023] Applicants also applied perturb-seq to dissect the mammalian unfolded protein response (UPR). First, Applicants used perturb-seq to build an epistatic map of UPR-mediated transcription. Then, Applicants conducted a genome-scale CRISPRi screen to identify genes whose depletion perturbs ER homeostasis and subjected a subset of our hits (using a ~100 element sublibrary) to perturb-seq, revealing high precision functional gene clusters. Single-cell analyses revealed both bifurcations in behavior in cells subject to the same perturbation as well as differential activation of the three UPR branches across hits, including an isolated feedback loop between the translocon and the IRE1 α branch of the UPR. These studies provide insight into how the three separate sensors of ER homeostasis allow the cell to monitor distinct types of ER stress and, more generally, highlight the ability of perturb-seq to dissect complex cellular responses.

[0024] By combining droplet based single cell transcriptomics with CRISPR-Cas based perturbations, Applicants demonstrate an approach that allow researchers to perform thousands of these assays in a single pooled experiment. This pooled approach represents a 10-fold

improvement in cost over current methods for obtaining perturbation transcription profiles. Leveraging the discrete nature of the single cell measurements, the screening approach also has the ability to resolve novel phenotypes such as the effect of a perturbation on cell type composition or cell cycle phase, and filter unperturbed cells whose presence otherwise dilutes the measured effect in population measurements. Finally, by randomly integrating more than one sgRNA in each cell, applicants show the ability to extend the method to test combinatorial effects. Together, these statistical and experimental methods enable researchers to perform large-scale screening of perturbations, including systematic dissection of epistatic effects, using RNA transcription as a phenotype.

[0025] In one aspect, the present invention provides for a method of reconstructing a cellular network or circuit, comprising introducing at least 1, 2, 3, 4 or more single-order or combinatorial perturbations to a plurality of cells in a population of cells, wherein each cell in the plurality of the cells receives at least 1 perturbation; measuring comprising: detecting genomic, genetic, proteomic, epigenetic and/or phenotypic differences in single cells compared to one or more cells that did not receive any perturbation, and detecting the perturbation(s) in single cells; and determining measured differences relevant to the perturbations by applying a model accounting for co-variates to the measured differences, whereby intercellular and/or intracellular networks or circuits are inferred. The measuring in single cells may comprise single cell sequencing. The single cell sequencing may comprise cell barcodes, whereby the cell-of-origin of each RNA is recorded. The single cell sequencing may comprise unique molecular identifiers (UMI), whereby the capture rate of the measured signals, such as transcript copy number or probe binding events, in a single cell is determined. The model may comprise accounting for the capture rate of measured signals, whether the perturbation actually perturbed the cell (phenotypic impact), the presence of subpopulations of either different cells or cell states, and/or analysis of matched cells without any perturbation.

[0026] The single-order or combinatorial perturbations may comprise 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99 or 100 perturbations. The perturbation(s) may target genes in a pathway or intracellular network.

[0027] The measuring may comprise detecting the transcriptome of each of the single cells. The perturbation(s) may comprise one or more genetic perturbation(s). The perturbation(s) may comprise one or more epigenetic or epigenomic perturbation(s). At least one perturbation may be introduced with RNAi- or a CRISPR-Cas system. At least one perturbation may be introduced via a chemical agent, biological agent, an intracellular spatial relationship between two or more cells, an increase or decrease of temperature, addition or subtraction of energy, electromagnetic energy, or ultrasound.

[0028] The cell(s) may comprise a cell in a model non-human organism, a model non-human mammal that expresses a Cas protein, a mouse that expresses a Cas protein, a mouse that expresses Cpf1, a cell *in vivo* or a cell *ex vivo* or a cell *in vitro*. The cell(s) may also comprise a human cell.

[0029] The measuring or measured differences may comprise measuring or measured differences of DNA, RNA, protein or post translational modification; or measuring or measured differences of protein or post translational modification correlated to RNA and/or DNA level(s).

[0030] The perturbing or perturbation(s) may comprise(s) genetic perturbing. The perturbing or perturbation(s) may comprise(s) single-order perturbations. The perturbing or perturbation(s) may comprise(s) combinatorial perturbations. The perturbing or perturbation(s) may comprise gene knock-down, gene knock-out, gene activation, gene insertion, or regulatory element deletion. The perturbing or perturbation(s) may comprise genome-wide perturbation. The perturbing or perturbation(s) may comprise performing CRISPR-Cas-based perturbation. The perturbing or perturbation(s) may comprise performing pooled single or combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs. The perturbations may be of a selected group of targets based on similar pathways or network of targets.

[0031] The perturbing or perturbation(s) may comprises performing pooled combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs. Each sgRNA may be associated with a unique perturbation barcode. Each sgRNA may be co-delivered with a reporter mRNA comprising the unique perturbation barcode (or sgRNA perturbation barcode).

[0032] The perturbing or perturbation(s) may comprise subjecting the cell to an increase or decrease in temperature. The perturbing or perturbation(s) may comprise subjecting the cell to a chemical agent. The perturbing or perturbation(s) may comprise subjecting the cell to a biological agent. The biological agent may be a toll like receptor agonist or cytokine. The

perturbing or perturbation(s) may comprise subjecting the cell to a chemical agent, biological agent and/or temperature increase or decrease across a gradient.

[0033] The cell may be in a microfluidic system. The cell may be in a droplet. The population of cells may be sequenced by using microfluidics to partition each individual cell into a droplet containing a unique barcode, thus allowing a cell barcode to be introduced.

[0034] The perturbing or perturbation(s) may comprise transforming or transducing the cell or a population that includes and from which the cell is isolated with one or more genomic sequence-perturbation constructs that perturbs a genomic sequence in the cell. The sequence-perturbation construct may be a viral vector, preferably a lentivirus vector. The perturbing or perturbation(s) may comprise multiplex transformation or transduction with a plurality of genomic sequence-perturbation constructs.

[0035] In another aspect, or in alternative embodiments of aspects described herein, the present invention provides for a method wherein proteins or transcripts expressed in single cells are determined in response to a perturbation, wherein the proteins or transcripts are detected in the single cells by binding of more than one labeling ligand comprising an oligonucleotide tag, and wherein the oligonucleotide tag comprises a unique constituent identifier (UCI) specific for a target protein or transcript. The single cells may be fixed in discrete particles. The discrete particles may be washed and sorted, such that cell barcodes may be added, e.g. sgRNA perturbation barcodes as described above. The oligonucleotide tag and sgRNA perturbation barcode may comprise a universal ligation handle sequence, whereby a unique cell barcode may be generated by split-pool ligation. The labeling ligand may comprise an oligonucleotide label comprising a regulatory sequence configured for amplification by T7 polymerase. The labeling ligands may comprise oligonucleotide sequences configured to hybridize to a transcript specific region. Not being bound by a theory, both proteins and RNAs may be detected after perturbation. The oligonucleotide label may further comprise a photocleavable linker. The oligonucleotide label may further comprise a restriction enzyme site between the labeling ligand and unique constituent identifier (UCI). The ligation handle may comprise a restriction site for producing an overhang complementary with a first index sequence overhang, and wherein the method further comprises digestion with a restriction enzyme. The ligation handle may comprise a nucleotide sequence complementary with a ligation primer sequence and wherein the overhang complementary with a first index sequence overhang is produced by hybridization of the ligation

primer to the ligation handle. The method may further comprise quantitating the relative amount of UCI sequence associated with a first cell to the amount of the same UCI sequence associated with a second cell, whereby the relative differences of a cellular constituent between cell(s) are determined. The labeling ligand may comprise an antibody or an antibody fragment. The antibody fragment may be a nanobody, Fab, Fab', (Fab')₂, Fv, ScFv, diabody, triabody, tetrabody, Bis-scFv, minibody, Fab₂, or Fab₃ fragment. The labeling ligand may comprise an aptamer. The labeling ligand may be a nucleotide sequence complementary to a target sequence.

[0036] Single cell sequencing may comprise whole transcriptome amplification.

[0037] The method in aspects of the invention may comprise comparing an RNA profile of the perturbed cell with any mutations in the cell to also correlate phenotypic or transcriptome profile and genotypic profile.

[0038] In another aspect, or in alternative embodiments of aspects described herein, the present invention provides for a method comprising determining genetic interactions by causing a set of P genetic perturbations in single cells of the population of cells, wherein the method comprises: determining, based upon random sampling, a subset of π genetic perturbations from the set of P genetic perturbations; performing said subset of π genetic perturbations in a population of cells; performing single-cell molecular profiling of the population of genetically perturbed cells; inferring, from the results and based upon the random sampling, single-cell molecular profiles for the set of P genetic perturbations in cells. The method may further comprise: from the results, determining genetic interactions. The method may further comprise: confirming genetic interactions determined with additional genetic manipulations.

[0039] The set of P genetic perturbations or said subset of π genetic perturbations may comprise single-order genetic perturbations. The set of P genetic perturbations or said subset of π genetic perturbations may comprise combinatorial genetic perturbations. The genetic perturbation may comprise gene knock-down, gene knock-out, gene activation, gene insertion, or regulatory element deletion. The set of P genetic perturbations or said subset of π genetic perturbations may comprise genome-wide perturbations. The set of P genetic perturbations or said subset of π genetic perturbations may comprise k-order combinations of single genetic perturbations, wherein k is an integer ranging from 2 to 15, and wherein the method comprises determining k-order genetic interactions. The set of P genetic perturbations may comprise combinatorial genetic perturbations, such as k-order combinations of single-order genetic

perturbations, wherein k is an integer ranging from 2 to 15, and wherein the method comprises determining j -order genetic interactions, with $j < k$.

[0040] The method in aspects of this invention may comprise performing RNAi- or CRISPR-Cas-based perturbation. The method may comprise an array-format or pool-format perturbation. The method may comprise pooled single or combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs. The method may comprise pooled combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

[0041] The random sampling may comprise matrix completion, tensor completion, compressed sensing, or kernel learning. The random sampling may comprise matrix completion, tensor completion, or compressed sensing, and wherein π is of the order of $\log P$.

[0042] The cell may comprise a eukaryotic cell. The eukaryotic cell may comprise a mammalian cell. The mammalian cell may comprise a human cell. The cell may be from a population comprising 10^2 to 10^8 cells and DNA or RNA or protein or post translational modification measurements or variables per cell comprise 50 or more.

[0043] The perturbation of the population of cells may be performed *in vivo*. The perturbation of the population of cells may be performed *ex vivo* and the population of cells may be adoptively transferred to a subject. The population of cells may comprise tumor cells. The method may comprise a lineage barcode associated with single cells, whereby the lineage or clonality of single cells may be determined.

[0044] The perturbing may be across a library of cells to thereby obtain RNA level and/or optionally protein level, whereby cell-to-cell circuit data at genomic or transcript or expression level is determined. The library of cells may comprise or is from a tissue sample. The tissue sample may comprise or is from a biopsy from a mammalian subject. The mammalian subject may comprise a human subject. The biopsy may be from a tumor. The method may further comprise reconstructing cell-to-cell circuits.

[0045] The method may comprise measuring open chromatin and may comprise fragmenting chromatin inside isolated intact nuclei from a cell, adding universal primers at cutting sites, and uniquely tagging DNA that originated from the cell.

[0046] The method may comprise measuring protein and RNA levels and may comprise CyTOF.

[0047] In another aspect, the present invention provides for a method of determining any combination of protein detection, RNA detection, open chromatin detection, protein-protein interactions, protein-RNA interactions, or protein-DNA interactions comprising any of the preceding methods.

[0048] In another aspect, the present invention provides for a method for screening compounds or agents capable of modifying a cellular network or circuit comprising performing any method as described herein, wherein perturbing further comprises exposing the cell to each compound or agent.

[0049] In another aspect, the present invention provides for a method of identifying a therapeutic, and to a therapeutic identified by the method described herein.

[0050] In another aspect, the present invention provides a method of reconstructing a cellular network or circuit, comprising introducing at least 1, 2, 3, 4 or more single-order or combinatorial perturbations to each cell in a population of cells; measuring genomic, genetic and/or phenotypic differences of each cell and coupling combinatorial perturbations with measured differences to infer intercellular and/or intracellular networks or circuits. The single-order or combinatorial perturbations can comprise 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98 or 99 or massively parallel perturbations. The perturbation(s) can comprise one or more genetic perturbation. The perturbation(s) can comprise one or more epigenetic or epigenomic perturbation. The perturbation can be introduced with RNAi- or a CRISPR-Cas system. For example, reference is also made to Dahlman et al., Nature Biotechnology (2015) doi:10.1038/nbt.3390 Published online 05 October 2015 to allow efficient orthogonal genetic and epigenetic manipulation. Dahlman et al., Nature Biotechnology (2015) doi:10.1038/nbt.3390 have developed a CRISPR-based method that uses catalytically active Cas9 and distinct single guide (sgRNA) constructs to knock out and activate different genes in the same cell. These sgRNAs, with 14- to 15-bp target sequences and MS2 binding loops, can activate gene expression using an active *Streptococcus pyogenes* Cas9 nuclease, without inducing double-stranded breaks. Dahlman et al., Nature Biotechnology (2015) doi:10.1038/nbt.3390 use these 'dead RNAs' to perform orthogonal gene knockout and transcriptional activation in human cells.

[0051] The at least one perturbation can be introduced via a chemical agent, an intracellular spatial relationship between two or more cells, an increase or decrease of temperature, addition or subtraction of energy, electromagnetic energy, or ultrasound. The cell can comprise a cell in a model non-human organism, a model non-human mammal that expresses a Cas protein, a mouse that expresses a Cas protein, a cell *in vivo* or a cell *ex vivo* or a cell *in vitro*. The measuring or measured differences can comprise measuring or measured differences of DNA, RNA, protein or post translational modification; or measuring or measured differences of protein or post translational modification correlated to RNA and/or DNA level(s). The method can include sequencing, and prior to sequencing: perturbing and isolating a single cell with at least one labeling ligand specific for binding at one or more target RNA transcripts, or isolating a single cell with at least one labeling ligand specific for binding at one or more target RNA transcripts and perturbing the cell; and/or lysing the cell under conditions wherein the labeling ligand binds to the target RNA transcript(s).

[0052] The method in aspects of this invention may also include, prior to sequencing perturbing and isolating a single cell with at least one labeling ligand specific for binding at one or more target RNA transcripts, or isolating a single cell with at least one labeling ligand specific for binding at one or more target RNA transcripts and perturbing the cell; and lysing the cell under conditions wherein the labeling ligand binds to the target RNA transcript(s). The perturbing and isolating a single cell may be with at least one labeling ligand specific for binding at one or more target RNA transcripts. The isolating a single cell may be with at least one labeling ligand specific for binding at one or more target RNA transcripts and perturbing the cell.

[0053] The perturbing of the present invention may involve genetic perturbing, single-order genetic perturbations or combinatorial genetic perturbations. The perturbing may also involve gene knock-down, gene knock-out, gene activation, gene insertion or regulatory element deletion. The perturbation may be genome-wide perturbation. The perturbation may be performed by RNAi- or CRISPR-Cas-based perturbation, performed by pooled single or combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs or performing pooled combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

[0054] In another aspect, the methods described herein may be used for a diagnostic assay. In one embodiment, T cells are obtained from a subject and perturb-seq is performed on the cells. In

another embodiment, T cells are obtained from a subject and gene expression of single cells is determined. Upon determining gene expression, perturb-seq is performed on a subset of genes differentially expressed. Perturb-seq can inform proper therapies to administer to a subject and can test many targets in a single experiment. In another embodiment, tumor cells are obtained from a subject. The tumor cells may also include cells of the tumor microenvironment, such as immune cells. The cells may be assayed for gene expression and differentially expressed genes can be assayed using the perturb-seq methods described herein. Not being bound by a theory, perturb-seq may allow assaying many targets and perturbations in a single experiment.

[0055] Accordingly, it is an object of the invention not to encompass within the invention any previously known product, process of making the product, or method of using the product such that Applicants reserve the right and hereby disclose a disclaimer of any previously known product, process, or method. It is further noted that the invention does not intend to encompass within the scope of the invention any product, process, or making of the product or method of using the product, which does not meet the written description and enablement requirements of the USPTO (35 U.S.C. §112, first paragraph) or the EPO (Article 83 of the EPC), such that Applicants reserve the right and hereby disclose a disclaimer of any previously described product, process of making the product, or method of using the product. It may be advantageous in the practice of the invention to be in compliance with Art. 53(c) EPC and Rule 28(b) and (c) EPC. All rights to explicitly disclaim any embodiments that are the subject of any granted patent(s) of applicant in the lineage of this application or in any other lineage or in any prior filed application of any third party is explicitly reserved. Nothing herein is to be construed as a promise.

[0056] It is noted that in this disclosure and particularly in the claims and/or paragraphs, terms such as "comprises", "comprised", "comprising" and the like can have the meaning attributed to it in U.S. Patent law; e.g., they can mean "includes", "included", "including", and the like; and that terms such as "consisting essentially of" and "consists essentially of" have the meaning ascribed to them in U.S. Patent law, e.g., they allow for elements not explicitly recited, but exclude elements that are found in the prior art or that affect a basic or novel characteristic of the invention.

[0057] These and other embodiments are disclosed or are obvious from and encompassed by, the following Detailed Description. For the purpose of clarity and a concise description features

are described herein as part of the same or separate embodiments, however, it will be appreciated that the scope of the invention may include embodiments having combinations of all or some of the features described. The various embodiments as described herein for specific aspects of the invention, may be part of any aspect of the invention as described herein.

BRIEF DESCRIPTION OF THE DRAWINGS

[0058] The patent or application file contains at least one drawing executed in color. Copies of this patent or patent application publication with color drawing(s) will be provided by the Office upon request and payment of the necessary fee.

[0059] The following detailed description, given by way of example, but not intended to limit the invention solely to the specific embodiments described, may best be understood in conjunction with the accompanying drawings.

[0060] **Figure 1** illustrates a schematic overview of embodiments using hydrogel embedding of single cells, followed by lipid clearing and DNA-tagged antibody labeling. Also shown, are low-throughput and high-throughput readouts.

[0061] **Figure 2A-C** illustrates a proof of principle. **A)** Overlay of transmission light microscopy image epifluorescent measurement of CD51/Alexa488 (green). Black arrow indicates an example of a hydrogel embedded cell, the white arrow indicates an empty hydrogel droplet, staining negative for CD51. **B)** same cells, stained for genomic DNA DAPI (blue) and intracellular PCNA/Alexa647 (red). **C)** Strong Pearson correlation (0.98) measured between fluorescence levels of 'endogenous' GFP levels from 293/GFP cells, and a detecting anti-GFP antibody conjugated to Alexa647. Similar staining with the BD cytofix/perm kit yielded a Pearson correlation of 0.36.

[0062] **Figure 3** illustrates measuring protein levels by staining of aggregations of cellular constituents with high affinity reagents (antibodies) linked to an oligonucleotide with the structure [5' Amino Modifier]-[~6bp spacer]-[PhotoCleavable linker]-[~4bp spacer]-[Illumina PCR primer]-[~8-16bp UMI]-[~21bp UCI]-[~20bp universal ligation handle]. Note: UMI may be omitted in case of incorporation of a UMI in a split and pool index.

[0063] **Figure 4** illustrates hybridization of a ligation primer that binds to the universal ligation handle on oligonucleotide label with a sticky end needed for ligation of index A is produced.

[0064] **Figure 5** illustrates split-pool ligation using single-cell hydrogel drops as the basic unit and ligation of Index A, B and [C + PCR primer].

[0065] **Figure 6** illustrates staining in bulk with adjacent oligo's that hybridize to an RNA transcript or single guide RNA (sgRNA) at sites adjacent to each other.

[0066] **Figure 7** illustrates single probe detection of an RNA transcript or sgRNA using a single DNA probe that specifically binds to the target transcript.

[0067] **Figure 8** illustrates dual probe detection of an RNA transcript or sgRNA using adjacently binding probes that are ligated, such that only dually detection events are amplified.

[0068] **Figure 9** illustrates staining with the ligation primer and performing split-pool ligation with an Index A containing a UMI in index C such that sequencing starts with a random region (improves cluster detection) and ligation primer is no longer separately added, but pre hybridized before staining.

[0069] **Figure 10** illustrates an example of the generation of an Index A + UMI.

[0070] **Figure 11** illustrates measuring protein-protein complexes by performing a restriction enzyme digestion to generate an oligonucleotide containing two UCI and a compatible end for ligation to an index A for split-pool ligation.

[0071] **Figure 12** illustrates oligonucleotide structures for measuring protein-protein complexes. Oligo 1: [5' Amino Modifier]-[~6bp spacer]-[PhotoCleavable linker]-[~4bp spacer]-[Illumina PCR primer]-[~21bp UCI]-[~11bp Hybridization sequence 1]. Oligo 2: [5' Amino Modifier]-[~6bp spacer] -[RE site for sticky overhang] - [~21bp UCI]-[~11bp Hybridization sequence 1 complement].

[0072] **Figure 13** illustrates measuring protein-RNA complexes using proximity hybridization. The Final oligonucleotide to sequence contains the UCI protein, UCI RNA and UMI + USI via split-pool ligation protocol.

[0073] **Figure 14** illustrates high throughput single-cell ATAC-seq.

[0074] **Figure 15** illustrates high throughput single-cell measuring protein-DNA complexes.

[0075] **Figure 16** illustrates staining with an antibody bound to an oligonucleotide label and performing split-pool ligation with an Index C containing a UMI in index C such that sequencing starts with a random region (improves cluster detection) and ligation primer is no longer separately added, but pre hybridized before staining.

[0076] **Figure 17** illustrates alternative embodiments of measuring RNA levels.

[0077] **Figure 18** illustrates the generation of an Index C + UMI.

[0078] **Figure 19** illustrates a brightfield microscopy image showing hydrogel droplet encapsulated cells with magnetic particles embedded into the droplets to enable magnetic separation, aiding in clean up and washing steps in multiple reactions. Greatly enhances automation and therefore throughput.

[0079] **Figure 20** illustrates a novel probe for detection of complexes consisting of more than 2 cellular constituents at the same time. The probe includes a Unique Location Identifier (ULI). It can be constructed by rolling circle amplification.

[0080] **Figure 21** illustrates the overall scheme to measure the proximity of 3 or more proteins, RNA or DNA molecules. The hybridization sequence of the ligand bound oligo binds to the complementary hybridization site on the ULI probe. By extension, each ligand bound oligo incorporates the same ULI. Restriction enzyme digestion generates a 4bp overhang. Sticky end ligation is used to attach a USI + UMI.

[0081] **Figure 22** is an illustrative example of a Design Matrix according to aspects of the present invention.

[0082] **Figure 23** is an illustrative example of an Observed Measurement Matrix according to aspects of the present invention.

[0083] **Figure 24** is a plot of expected counts versus observed counts for a predetermined set of probes under a certain dilution criteria.

[0084] **Figure 25** is a plot of gene aggregations (observed versus expected) across random rows of a measurement vector for a 50 trial experiment.

[0085] **Figure 26** is plot of expected counts versus observed counts for a predetermined set of probes under a certain dilution criteria for Purified cDNA High.

[0086] **Figure 27** is plot of expected counts versus observed counts for a predetermined set of probes under a certain dilution criteria for unPurified cDNA High.

[0087] **Figure 28** is plot of expected counts versus observed counts for a predetermined set of probes under a certain dilution criteria for Purified cDNA Low.

[0088] **Figure 29** is plot of expected counts versus observed counts for a predetermined set of probes under a certain dilution criteria for unPurified cDNA Low.

[0089] **Figure 30-33** illustrate results of computational simulations where counts are aggregated across all genes from randomly chosen rows.

[0090] **Figure 34** illustrates an estimate number of cells (Y axis) needed to capture a given minimal number of cells per subtype (X axis) for BCs (rarest subtype ~5%), AC, and RGCs (rarest subtype ~1%).

[0091] **Figure 35** illustrates aDrop-Seq. From left: Device; clear separation of a mixture of mouse and human cells; and t-SNE of 44,808 single cells from a mouse retina, distinguishing 39 subpopulations (colors).

[0092] **Figure 36** illustrates a Drop-ATAC. Left: schematic. Right: Number of human (x-axis) vs. mouse (y-axis) reads in each barcode in a species cell mixing experiment.

[0093] **Figure 37** illustrates aSingle cell ATAC-Seq in DCs.

[0094] **Figure 38** illustrates aSingle cell protein + RNA by PEA. Top: PEA: Bottom: PEA distinguishes single THP-1 and K562 cells.

[0095] **Figure 39** illustrates a NGS proteomics in hydrogel droplets.

[0096] **Figure 40** illustrates a Rapid sequential FISH of 5 genes using five sequential rounds.

[0097] **Figure 41** illustrates a Linkage PCR in gel droplets of two CRISPR barcodes.

[0098] **Figure 42** illustrates a Proportion of sites (x-axis) where perturbation of bound TF (y-axis) affects target's expression.

[0099] **Figure 43** illustrates a Perturbation model of Th17 cell differentiation.

[00100] **Figure 44** illustrates a Genetic (left) vs. molecular (right) models.

[00101] **Figure 45** illustrates a Hog network. Each set of incoming edges is quantified for each 1–5-way contribution (not shown).

[00102] **Figure 46** illustrates MCPP assays in Cas9-expressing DCs. Cells are infected with a library of sgRNA-expressing lentiviruses that target regulators; DCs are stimulated with LPS; finally, millions of cells are profiled by scRNA-seq or NGS proteomics to monitor changes in gene or protein expression by one or multiple sgRNAs per cell.

[00103] **Figure 47** illustrates a CRISPR screen. TLR pathway with screen hits marked (blue). Illustrative flow cytometry staining of Tnf for (shaded) vs. sgRNA controls (empty).

[00104] **Figure 48A-G** illustrates Perturb-seq: pooled approach to obtain screening of transcriptional profiles of perturbations. (A) Overview of approach. From left: A complex pool of lentiviruses, each carrying a guide targeting a specific genes (A-E) in a vector including matched PAPIs is transduced at a given MOI (γ_{MOI}) into cells that are subsequently selected (γ_{MOI}) through growth and/or fluorescence. Individual cells are profiled with droplet scRNA-seq, to

simultaneously recover their RNA profiles and PAPIs. **(B)** Perturb-Seq vector. A transcribed PAPI is pre-associated with a specific sgRNA, and enriched after WTA using dial-out PCR. **(C)** Perturb-Seq screens in this study. **(D)** The observed distribution of number of guides detected per cell (X axis) in BMDCs at 3h post-LPS stimulation. **(E)** Effect on target. Hif1a transcript expression (Y axis) in cells carrying a an sgRNA targeting Hif1a (sgHif1a-1, right) compared to all other cell (left). Box plots denote [describe all aspects: mean/median? Percentiles, outliers etc]. **(F)** Distribution of on-target effects in BMDC at 3h post-LPS. Shown is the distribution of on-target expression (X axis) in cells carrying the corresponding targeting guides (blue) and permuted results for a single permutation (green). Rectangle is the 99% confidence interval for the permuted mean. Mean on-target effect of individual guides are in tickmarks, including one outlier exceeding even the permuted data. **(G)** Linear modeling framework. Applicants fit the coefficients of a (regulatory) matrix (β) to best explain the observed expression profiles of each cell (matrix Y) given the covariates in the design matrix (X), including the observed sgRNAs in each cell and additional cell covariates. See also **Figure 55**.

[00105] **Figure 49A-J** illustrates MIMOSCA: A linear model stratifies analysis framework to stratify the observed expression variation.. **(A)** A linear modeling framework for continuous phenotypes. A linear model relates a continuous phenotype (gene expression, PC scores, etc; arrow) to a covariate (here guide identity). **(B-D)** Accounting for differences in cell complexity and state. Scatter plots show the relation for every cell (dot) between the expression levels of the highly expressed Ccl17 gene (B, Y axis) or its residual expression after a model is fit (C, D; Y axis) and the number of transcripts in the same cell (X axis; $\text{sum}(\log(\text{transcripts detected}))$) in the original data (B), after including the quality measure as a covariate residuals (C) and after also including in addition cell state proportion (D). **(E)** Effect on cell states. Cartoon illustration of a hypothetical scenario where cells belong to either of two states (red, blue) and perturbation by sgRNA_1 increases the proportion of cells in one state over the other. **(F,G)** Accounting for cell states. Shown is the effect on Ccl17 expression (Y axis) in cells with (right box) and without (left box) sg-Rela-3 in the original model (left) and when conditioned on cell state proportions (right). **(H)** Expectation maximization-like framework to separate effectively wild type cells within the set of cells containing a particular sgRNA. Left: The distribution of number of cells with the sgStat1-3 guide with a given fit (X axis) to the model of the effect of the Stat1-3 perturbation. Right: the distribution of inferred probabilities of an actual perturbation in the same cells. **(I)**

Contribution of each component in the model (Y axis) to the % variance explained (X axis) based on the r^2 values from cross-validation. (J) Correlation matrix between genes in the residuals of the model. The addition of our covariates reduces the residual structure in the data. See also **Figure 56**.

[00106] **Figure 50A-F** illustrates the analysis of the role of 24 TFs in the response of BMDCs to LPS. (A) Key TF modules in BMDCs at 3h post-LPS. Heat map shows the Pearson correlation (red-blue colorbar) between the coefficients of the regulatory matrix β for every pair of guides (rows, columns), based on a model that does not include the cell state covariates. Yellow rectangles: four key TF modules high (members marked on right). Leftmost column: effect of the guide on-target. (B) Agreement between guides targeting the same gene. Shown is the distribution of correlations between guides targeting the same gene (grey) or different genes (blue). (C) Cell states. Shown are enrichments ($-\log_{10}(\text{P-value})$) for induced (red) and repressed (blue) genes) for GO gene sets (rows) in each of seven cell states (columns) defined for genetically unperturbed BMDC at 3h post-LPS. Key terms are marked on right. (D) TFs group by their effects on cell state proportion. Shown are the Z-score for enrichment (red) or depletion (blue) of guides in cells with each of seven states (columns; as in C). (E, F) Accounting for cell states emphasizes TF-specific effects. (E) Heat map shows the Pearson correlation (red-blue colorbar) between the coefficients of the regulatory matrix β for every pair of guides (rows, columns), based on a model that accounts for cell states as covariates. Yellow rectangles: specific modules all with guides to a single gene (gene marked on right). Leftmost column: effect of the guide on-target. (F) The distribution of correlations between guides targeting the same gene (grey) or different genes (blue), with a model including cell state covariates. See also **Figure 57**.

[00107] **Figure 51A-G** illustrates Perturb-Seq recovers a gene regulatory circuit for BMDCs balancing cell states and immune responses. (A,B) Four TF modules control five transcriptional programs in BMDCs at 3h post-LPS. (A) Shown is the regulatory coefficient of each guide (labeled and color coded columns) on each gene (rows) based on a model that does not account for cell states as co-variables. Guides and genes are clustered by the similarity of their profiles. Four TF modules (M1-4) and five target programs (P1-5) are marked. Green-white heatmap shows the enrichment in bound targets of each TF based on ChIP-Seq (Garber et al., 2012) in each program (rows). (B) Bi-partite graph, based on (A) of association of TF modules (top) to

target programs (bottom). Blue (red) arrows: TF in module activates (inhibits) program (perturbation has the opposite effect). Bottom: TFs from modules that are members of the regulated program (blue: activator of the program; red: repressor of the program). (C,D) TF-TF circuit in BMDCs at 3h post-LPS. (C) Heatmap as in (A) but only showing genes (rows) that encode the TFs targeted by the guides. (D) TF-TF circuit based on the associations in (C). TFs are shown as nodes. TF modules are shadowed in grey. Blue (red) arrows: TF in module activates (inhibits) program (perturbation has the opposite effect). (E-F) ChIP-Seq support for genetic network. (E) Expected effects of perturbation of a TF (green) on a target gene when the factor does not regulate (or bind) the gene (top); directly binds and represses it (middle) or directly binds and activates it (bottom). (F) Regulatory effect of each guide (rows) on its bound target (colorbar) as defined by ChIP-Seq in BMDCs at four time points (columns, 0, 0.5, 1, and 2h post-LPS). (G) The proportion of bound targets at 120 min post-LPS for the TF targeted by each guide (rows) that are repressed (blue), activated (red) or unaffected (grey) by the perturbation. Asterisks denote significant enrichments (as in F). See also **Figure 58**.

[00108] **Figure 52A-F** illustrates dissecting genetic interactions between perturbations using Perturb-seq. (A) Extending the linear model framework to include higher order polynomial functions as the product of single covariates. (B) Interaction analysis on cell states for BMDCs at 3h post-LPS. Heatmap shows the enrichment (red) or depletion (blue) of single, pair and triplets of guides (rows) in cells in each of seven states (as in **Figure 50C**). (C) A three-way genetic interaction enhances cell state 2. Box plot shows the distribution of probabilities of assignment to cell state 2 marginally contributed by the individual perturbation of each of NFKB1, Stat2, and Rel, their pair wise interaction and three-way interaction. (D) Genetic interaction categories. A schematic illustrating the 27 possible schemes of genetic interaction between two genes (A,B), when considering all permutations of significantly positive (red), negative (blue) or no (white) regulatory coefficients marginally associated with each of the two individual guides or their combination. The four key categories of relationships (buffering, synergistic, dominant, and additive) are color coded on left. (E) The distribution genes in each of the 27 categories (rows) for every pair of perturbations assayed for interaction (columns). Key examples are marked on bottom. (F) Genetic interaction between Rela and Nfkb1 associated with co-binding. Shown are the marginal regulatory coefficients for Rela, Nfkb1 and their interaction term for each gene (rows) with at least one non-zero coefficient, sorted by the four

key categories (color code, left). ChIP-Seq enrichment of individually bound and co-bound targets is shown for each of the groups (right). See also **Figure 59**.

[00109] **Figure 53A-I** illustrates the analysis of the role of 10 non-essential TFs and 14 cell cycle regulators in K562 cells. **(A,B)** Grouping of TFs by their regulatory coefficients. Shown are the Pearson correlation matrix of the coefficients of the regulatory matrix β of the TFs without including the cell state covariates (A) and after including them (B). **(C)** Cell states in genetically unperturbed K562 cells. Shown are enrichments ($-\log_{10}(\text{P-value})$) for induced (red) and repressed (blue) genes) for GO gene sets (rows) in each of seven cell states (columns) defined for genetically unperturbed BMDC at 3h post-LPS. Key terms are marked on right. **(D)** TFs effects on cell state proportions. Shown are the Z-score for enrichment (red) or depletion (blue) of guides in cells in each of seven states (columns; as in C). **(E,F)** Agreement of guide effects across time points. Violin plots show the distribution of correlations between guides targeting the same gene (grey), different genes (blue) and a gene and an intergenic region (red) with each of two time points (T1=7d, T2=14d) and across the two time points, in either a model that does not (E) or does (F) include cell state covariates. **(G)** Pearson correlation matrix of the coefficients of the regulatory matrix β of the cell cycle regulators without including the cell state covariates. **(H)** The regulatory effect (color bar, defined by the average regulatory coefficients in our model) or the guides targeting each cell cycle gene (rows) on the genes associated with each of six signatures of phases of the cell cycle and apoptosis (columns). **(I)** Distribution of fitness effects (X axis) across the guides targeting each gene (Y axis) by comparing the number of cells with the guide to its abundance in the initial pool. See also **Figure 60**.

[00110] **Figure 54A-D** illustrates a power analysis and prospects for Perturb-seq. **(A,B)** Saturation analysis. Shown is the effect of the number of cells (Y axis) and reads (X axis) on our ability to recover a given level of correlation (color bar) with either the per gene transcriptional defects (A) or gene signature effects (*e.g.*, cell type proportions) (B) Applicants observed with our full data. **(C)** The tradespace of number of cells (X axis) and measurements per cell (Y axis) required for screening scale transcriptome measurements. **(D)** Extensions of Perturb-Seq. The perturb seq framework can be extended by scaling the number of cells (arrow, left) or by incorporating additional types of cell covariates, such as lineage, marker expression, or temporal tracers.

[00111] **Figure 55A-L** illustrates the performance of Perturb-seq (A-E) Log-likelihood functions as a function of detection probability (X axis) and MOI (Y axis) for our zero-truncated zero inflated Poisson distribution for the indicated experiment (label on top). (A-C) Line plots are cumulative distributions for the observed distribution of guides per cell (blue) and the expectation from the maximum likelihood estimate (green). (F) Specificity of PAPI detection. Scatter plot shows the percentage of reads for a PAPI within a given cell (CBC – cell barcode) (Y axis) as a function of the $\log_2(\text{reads})$ it received in that cell (X axis). (G) Effect on target. Cebpb transcript expression (Y axis) in cells carrying an sgRNA targeting Cebpb (sgCebpb_1, right) compared to all other cell (left). Box plots. (I) Relationship between overall mean expression of the on-target gene (X axis) and the observed effect on its expression (Y axis) by the guides that target it, in BMDCs at 3h post-LPS. (J,K) Relationship between population expression measurements and 10 cell average (top) and 100 cell average for BMDCs (J) and K562 cells (H). (L) Relationship between transcript length (X axis) and the difference between population expression and single cell average expression (Y axis).

[00112] **Figure 56A-I** illustrates the performance of MIMOSCA framework. (A-C) Contribution of each component in the model (Y axis) to the % variance explained (X axis) based on the r^2 values from cross-validation in each of 3 screens (labeled on top). (D-F) Significance of regulatory coefficients. Shown are the distributions of signed $\log_{10}(\text{FDR})$ for each of three sgRNAs (marked on top). Values capped at 3; Zero coefficients (due to shrinkage by regularization) have no assigned FDR. (G) Relationship between number of cells/sgRNA (X axis) and number of significant genes for each sgRNA (Y axis). One outlier guide (sgStat1_2) is marked and – while included in all subsequent plots – was not considered in subsequent biological analysis. (H) Computational flowchart for MIMOSCA. Raw input data (orange) is processed into intermediate data types that require little computation (purple) or a more intensive processing (green). Finally, output data (blue) is generated and interpreted biologically. (I) Significant guide effects. Scatter plots show conceptual (left) and real (middle and right) examples of the sum of difference in the sum of the squared residuals between a model fit including a guide covariate and a model not including the guide covariate for positive and negative values. Each dot is a cell. Conceptually (left column), a pure population of perturbed cells that is well separated from the other cells should have few genes that are better fit without the guide covariate (left column, top and middle), whereas a poor performing guide (or a well-

performing control) should have few such genes (left column, bottom panel). A guide targeting *Rela* in BMDC (3h post-LPS) (top middle and top right) has a distinguishable effect compared to a non-targeting guide (below) in a model without cell state covariates (middle) and even when considering cell states (right). Bottom middle and right plots show the mean deviation of each guide (X axis). Red line: the non-targeting control.

[00113] **Figure 57A-F** illustrates analysis analyses of the role of 24 TFs in BMDCs. **(A)** Cell states in BMDCs pre-stimulation. Shown are enrichments ($-\log_{10}(\text{P-value})$) for induced (red) and repressed (blue) genes) for GO gene sets (rows) in each of four cell states (columns) defined for genetically unperturbed BMDC at 0h (pre-stimulation). Key terms are marked on right. **(B)** TFs group by their effects on cell state proportions pre-stimulation. Shown are the Z-score for enrichment (red) or depletion (blue) of guides in cells with each of four states (columns; as in A). **(C)** Distinct effects pre- and post-stimulation. Boxplots comparing the correlations of regulatory coefficients between guides targeting the same genes within the 3hr LPS stimulated cells (light blue) or within unstimulated cells (white), and across the two conditions either without modeling cell state co-variates (dark grey) or with modeling cell states (blue). **(D)** Distribution of number of detected transcripts per cell (X axis) for cells harboring guides targeting distinct genes (Y axis). **(E)** Relationship between the relative abundance of an sgRNA in the input plasmid pool (X axis) and the relative amount of cells with that sgRNA (Y axis). **(F)** Assessing potential fitness effects. Shown is the distribution of fold changes of sgRNA abundance compared to the input abundance (X axis) for the guides (dots) targeting each gene (Y axis).

[00114] **Figure 58A-D** illustrates additional aspects of the regulatory circuitry of BMDCs. **(A)** Relation between cell states and regulatory programs. Shown is the significance ($-\log_{10}(\text{P-value})$ of the hypergeometric test) of the overlap between the genes in each of the programs P1-P5 (as in **Figure 51A**) and the genes induced in each of the seven transcriptional states of genetically unperturbed DCs at 3h post-LPS (as in **Figure 50C**). **(B)** TF control of transcriptional programs in BMDCs at pre-stimulation (0h). Shown is the regulatory coefficient of each guide (labeled columns) on each gene (rows) based on a model that does not account for cell states as co-variates. Guides and genes are clustered by the similarity of their profiles. Four target programs (P1-4) are marked. Green-white heatmap shows the enrichment in bound targets of each TF based on ChIP-Seq (Garber et al., 2012) in each program (rows) based on either a lenient, genome-wide background (top) or a strict background (bottom) restricted only to the genes in the

four programs. (C) Enrichment in bound genes based on a lenient background. Green-white heatmap shows the enrichment in bound targets of each TF based on ChIP-Seq (Garber et al., 2012) in each program (rows) based on either a lenient, genome-wide background for the same model as in **Figure 51**. (D) Relation between TFs and transcriptional programs in BMDCs at 3h post-LPS after accounting for cell states. Shown is the regulatory coefficient of each guide (labeled columns) on each gene (rows) based on a model that accounts for cell states as co-variables. Guides and genes are clustered by the similarity of their profiles. Five target programs (P1-5) are marked (distinct from the programs of **Figure 51A**).

[00115] **Figure 59** illustrates the genetic interactions for K562 TFs.

[00116] **Figure 60A-D** illustrates additional analysis of the role TFs and cell cycle regulators in K562 cells. (A) Assessing potential fitness effects of TF perturbations in K562 cells. Shown is the distribution of fold changes of sgRNA abundance compared to the input abundance (X axis) for the guides (dots) targeting each gene (Y axis). (B) TF control of transcriptional programs in K562 cells. Shown is the regulatory coefficient of each guide (labeled columns) on each gene (rows) based on a model that either does not (B) or does (C) account for cell states as co-variables. Guides and genes are clustered by the similarity of their profiles. Target programs are marked with key enriched annotations. Green-white heatmap shows the enrichment in bound targets of each TF based on ChIP-Seq in each program (rows) based on either a lenient, genome-wide background (top) or a strict background (bottom) restricted only to the genes in the programs. (C) Cell cycle genes in our screen. Shown is the cell biological classification of our genes (Figure reproduced from Neumann et al). (D) Effects of perturbing cell cycle genes on transcriptional programs in K562 cells. Plot as in (B) but for the cell cycle regulators.

[00117] **Figure 61A-L** illustrates a saturation analysis for differential expression, related to **Figure 54**. (A) Theoretical probability of having a successful perturbation in every target as a function of the number of perturbations ($1-p^n$), assuming independence. (B-L) Saturation analysis. Shown is the effect of the fraction of cells (Y axis) and reads (X axis) on our ability to recover a given level of correlation (color bar) with either the PCA scores (B-D), gene signature effects (E-H) or per gene transcriptional defects (I-L) based on the BMDC 3h stimulated screen data. The total number of cells (1.0) is, on average, 300 cells/perturbation and the total number of transcripts per cell is, on average 5,000. PCA scores (as the expression matrix Y) were obtained by projecting the data onto PCs defined by the wildtype cells for the top 11 PCs conditioned on

different effect sizes. Effect size units are arbitrary. For gene signatures, effects sizes are in units of average $\log_2(\text{UMI})$ across the gene signature. For an individual gene level, effect sizes are in units of $\log_2(\text{UMI})$.

[00118] **Figure 62** illustrates Perturb-Seq analysis of the role of TFs in the response of BMDCs to LPS.

[00119] **Figure 63A-D** illustrates Perturb-Seq analysis of the role of cell cycle regulators in K562 cells. (A) Regulation of apoptosis [“virtual FACS”]: control, all guides, selected example. (B) Individual and synthetic effects on cell fitness (Livnat model). (C) Individual and synthetic effects on G1/S and G2/M cell states (“virtual FACS” plus Livant model). (D) Transcriptional patterns underlying effect on cell states.

[00120] **Figure 64** illustrates Perturb-Seq analysis of the role of non-essential TFs in K562 cells.

[00121] **Figure 65** illustrates the prospects for Perturb-Seq.

[00122] **Figure 66** illustrates GO terms associated with the following sgRNAs.

[00123] **Figure 67** illustrates Gene set (GO/MSigDB) analysis.

[00124] **Figure 68** illustrates a graphical embodiment of a perturb-Seq experiment.

[00125] **Figure 69** illustrates RT primers that are used to tag the 3' ends of transcripts.

[00126] **Figure 70** illustrates a vector for introducing a guide RNA and the actual transcript insert after capture.

[00127] **Figure 71** illustrates a graphical embodiment of a perturb-Seq experiment.

[00128] **Figure 72** illustrates that capture primers contain two barcodes specifying cell and molecule identity. Recover guide information is measured by co-expressing a polyadenylated guide barcode from a guide RNA vector.

[00129] **Figure 73** illustrates read statistics per cell.

[00130] **Figure 74** illustrates guide identification.

[00131] **Figure 75** illustrates capture of the guide barcode transcripts during RNAseq and specific amplification.

[00132] **Figure 76** illustrates singles, doubles and triplets during droplet sequencing.

[00133] **Figure 77** illustrates a graphical embodiment of a perturb-Seq experiment.

[00134] **Figure 78** illustrates a triple guide Drop-Seq vector.

[00135] **Figure 79** illustrates a triple guide Drop-Seq vector.

[00136] **Figure 80** illustrates a comparison of methods for single cell RNA-seq.

[00137] **Figure 81A-F** A robust strategy for systematic genetic modifier screens using single cell expression profiling. (A) Schematic of the perturb-seq platform. (1) Cells transduced with a complex pool of CRISPRi guide RNAs (2) are encapsulated with cell lysis buffer and gel beads, which deliver DNA oligos, in droplets using the Chromium™ instrument (10X Genomics). DNA oligos encode cell barcodes (unique to each bead), unique molecular identifiers or UMIs (unique to each bead oligo), and an oligo-dT homopolymer region. (3) After capture, beads dissolve and released oligos prime cellular mRNAs for cDNA synthesis within droplets. (4-6) The cDNA libraries are prepared for deep sequencing in pooled format, and transcripts of particular interest are separately enriched by specific amplification. (B) Method for introducing and identifying perturbations using CRISPR guide RNAs. A lentiviral vector designed for single copy integration into cells expressing Cas9 or dCas9 fused to effector domains (for CRISPRi/a) carries two expression units: (1) a polymerase III-driven sgRNA and (2) a guide barcode selection cassette (top) in reverse orientation. The later transcript can be selectively amplified during library preparation (bottom). (C) Performance of guide barcode capture by specific amplification. The plot shows the distribution of guide barcode UMIs counted per cell in a pilot experiment of ~5,700 cells containing different perturbations. (D) Performance of perturbation identification. The high coverage of the guide barcode enables efficient identification of perturbation identity, and confident rejection of cell doublets (apparent as multiple guide barcodes attached to a single cell barcode). (E) Characterization of perturb-seq vector by GFP knockdown. GFP+ K562 cells with dCas9-KRAB were transduced with either the perturb-seq vector or the original CRISPRi vector expressing a GFP-targeting sgRNA or a perturb-seq vector expressing a negative control sgRNA. GFP levels were measured by flow cytometry 11 d after transduction. Plotted are kernel density estimates of normalized flow cytometry counts for infected (BFP+) cells. (F) Characterization of perturb-seq vector in pooled format via single-cell RNA-seq. Using identities inferred from the capture of guide barcodes, Applicants computationally separate our pilot perturb-seq experiment into subpopulations of cells containing each sgRNA. The single-cell RNA-seq profiles allow us to assess average knockdown of the gene targeted within each subpopulation (relative to cells containing a control sgRNA). The plot shows three examples.

[00138] **Figure 82A-G** Strategy for multiplexed delivery of CRISPR guide RNAs in a single expression vector. (A) Schematic of the unfolded protein response. (B) Schematic of three-guide vector. Three guide RNA expression cassettes (a U6 promoter, a sgRNA targeting sequence, and a sgRNA constant region (cr)) are fused to express the three sgRNAs. The lentiviral backbone is the same as that of the perturb-seq vector. (C) Characterization of initial three-guide vector by GFP knockdown. GFP⁺ K562 cells with UCOE-dCas9-KRAB were transduced with either the single perturb-seq vector expressing a GFP-targeting sgRNA, an initial three-guide vector expressing a GFP-targeting sgRNA from the hU6 promoter, or a single vector expressing a negative control sgRNA. Cells were selected to purity with puromycin and GFP levels measured by flow cytometry 7 d after transduction. Plotted are kernel density estimates of normalized flow cytometry counts. (D) Design and characterization of constant region variants. Top: Schematic representation of sgRNA constant region, with location and description of changes indicated. Changes that were combined to generate cr2 and cr3 are labeled in orange and purple, respectively. Bottom: RFP levels of *E. coli* CRISPRi reporter strain expressing a mRFP-targeting sgRNA and the denoted constant region variant. RFP levels, determined by flow cytometry, are plotted as mean fluorescence of three replicates relative to a strain expressing a negative control sgRNA. Error bars denote standard deviations. (E) Characterization of new sgRNA expression cassettes by GFP knockdown. GFP⁺ K562 cells with dCas9-KRAB were transduced with different constructs expressing a GFP-targeting sgRNA from different U6 promoters and with different constant regions, as indicated. GFP levels were measured by flow cytometry 10 d after transduction. Plotted are kernel density estimates of normalized flow cytometry counts for infected (BFP⁺) cells. (F) Characterization of final three-guide vector by GFP knockdown. GFP⁺ K562 cells with UCOE-dCas9-KRAB were transduced with three-guide vectors expressing a GFP-targeting sgRNA from the position indicated in parentheses and two different negative control sgRNAs from the other two positions or a three-guide vector expressing three negative control sgRNAs. Cells were treated as described in (B). Plotted are kernel density estimates of normalized flow cytometry counts. Plot for the single perturb-seq construct is the same as in panel (B). (G) Characterization of the three-guide vector via perturb-seq. Using the three-guide vector, Applicants simultaneously targeted the three main UPR sensors and determined the average reduction in expression (relative to cells containing three negative

control guides) via perturb-seq. Each target showed >80% depletion, and depletion was stable across two chemical treatments that induce the UPR. See also Figure 88.

[00139] **Figure 83A-F** Epistatic analysis of the three transcriptional arms of the unfolded proteins response using perturb-seq (A) Schematic of perturb-seq epistasis experiment. (B) Unbiased identification and decoupling of single-cell behaviors via low rank independent component analysis. Individual single-cell RNA-seq profiles are noisy, but shared regulation and correlated gene expression mean that patterns exist within the population. To identify these patterns, Applicants computationally construct a low-dimensional approximation of gene expression within the population to remove noise, and then use ICA to identify a small number of independent programs of gene expression. The figure shows this approach applied to our combinatorial knockdowns treated with thapsigargin. Low rank ICA identifies distinct sets of components that vary either across perturbation or across the cell cycle. When the first set of components are projected via t-sne, the cells (each dot) group via perturbation. Conversely when Applicants project the second set, the cells arrange in a circular pattern by cell cycle position. The “bulge” highlighted by the dashed line is enriched for cells that have PERK active. (C) Identification of a PERK- and cell-cycle-dependent subpopulation in thapsigargin-treated cells. The plots show t-sne projections of control (+DMSO) and thapsigargin-treated cells with or without PERK. Low rank ICA identifies a component (IC) that is bimodal within each perturbation subpopulation and marks G1 cells, and that is particularly disenriched in G1 cells within the thapsigargin-treated subpopulation. (D) Cell cycle composition of +DMSO, +Tg, and +Tg cells with PERK depleted. (E) Genetic interactions between G1 and PERK activation. Applicants split each perturbation subpopulation into G1 and non-G1 cells based on the value of IC, and constructed average expression profiles in each condition. Applicants then examined how the 50 genes that most influenced IC varied, exposing both synergistic and antagonistic interactions between progression through G1 and PERK activation. (F) Epistatic interactions among the three branches of the UPR. Applicants isolated 104 genes that strongly varied in our large-scale perturb-seq experiment, and clustered them by their expression pattern within the population. The heatmap shows average expression profiles for each perturbation and chemical treatment. Patterns of induction determine the branch specificity of each gene. The bottom panel shows an unbiased decomposition of the total response into three components obtained via ICA, showing that many ATF6/IRE1 α targets have some overlapping regulation. See also Figure 89.

[00140] **Figure 84A-G** Genome-scale CRISPRi screening for genetic stresses that perturb the IRE1 branch of the unfolded protein response. (A) Schematic of UPRE (mCh) and constitutive EF1a (GFP) reporter cassettes. (B) K562 reporter cells (cBA011) were transduced with the indicated sgRNAs and treated with 2 $\mu\text{g}/\text{mL}$ tunicamycin or DMSO after 4 days. Approximately, 12 hours later these cells were evaluated by flow cytometry. Data is representative of two independent experiments. (C) Schematic of sgRNA screen. K562 cells stably expressing the mCh/GFP reporter cassettes and a dCas9-BFP-KRAB fusion protein were transduced with pooled genome-scale hCRISPRi libraries. Transduced cells were selected, sorted for high and low mCh/GFP ratio, and processed for sequencing of the sgRNA-containing DNA cassettes. (D) Volcano plots of UPRE reporter gene phenotypes and p-values from hCRISPRi-v2 reporter screen. Data generated from negative control sgRNAs are indicated in gray. Screen hits and select genes are indicated in pink and red, respectively. (E) Gene UPRE reporter phenotypes from replicates of hCRISPRi-v2 screen in K562 cells (cBA011). Phenotypes generated from negative control sgRNAs are indicated in gray. (F) UPRE reporter gene phenotypes from hCRISPRi-v2 reporter screen by functional category. Red indicates screen hits. (G) Comparison of mCherry UPRE reporter signal to EF1a driven GFP in cBA011 cells transduced with 257 sgRNAs targeting 152 hit genes from the hCRISPRi-v2 screen and 3 distinct negative controls. Data represent \log_2 averages of median fluorescence signals (normalized to untransduced cells) across four possible experiments (n=2-7 technical or experimental replicates), normalized medians from all four experiments are included in the \log_2 averages for each of the control sgRNAs. See also Figure 90.

[00141] **Figure 85A-D** A large-scale perturb-seq experiment interrogating ER homeostasis (A) Functional clustering of hits from perturb-seq analysis of ER homeostasis. Applicants picked ~100 guides from our genome-wide screen and subjected them to perturb-seq analysis, totaling ~65,000 single cells. Average expression profiles were created from all cells bearing guides targeting the same gene and hierarchically clustered. The figure shows a heatmap of correlations between expression profiles for all perturbations along with functional annotations. (B) Cell cycle analysis of perturb-seq hits. For each perturbation, the fraction of cells in each cell cycle stage was computationally identified. The figure shows the change in composition induced by each perturbation relative to control cells (containing the NegCtrl-2 guide). (C) Target knockdown efficiency. Average depletion of the sgRNA target was assessed within each

subpopulation. Genes targeted by multiple guides have multiple possibly overlapping dots. Error bars are 95% confidence intervals estimated by bootstrapping. (D) Average phenotypes of perturb-seq hits. The UPRE score from the primary genome-wide screen is compared to three computationally derived scores measuring activation of the three branches of the UPR for each perturbation. The final panel is the \log_{10} number of genes differentially expressed relative to control cells (containing the NegCtrl-2 guide) measured by the Kolmogorov-Smirnov test at $P < 0.01$. Genes targeted by multiple guides have multiple possibly overlapping dots. See also Figure 91.

[00142] Figure 86A-L Single-cell information reveals a bifurcated UPR within a population and allows unbiased discovery of UPR-controlled genes (A) Single-cell analysis of HSPA5-perturbed cells. The figure shows t-sne plots of guide identity, cell cycle position, and UMI count per cell in HSPA5-perturbed cells and control cells (containing the NegCtrl3 guide). (B) Low rank ICA analysis of HSPA5-perturbed cells identifies two subpopulation-defining independent components. The right panel shows a discretized breakdown of the cells based on applying a threshold to IC1. (C) Branch activation scores in HSPA5-perturbed cells. (D) Expression of UPR genes in HSPA5-perturbed cells. The figure shows the normalized expression matrix of all HSPA5-perturbed cells. Each row is a cell, and each column is a gene in the same order as Figure 83F. The cells have been ordered by increasing value of IC1. (E) Cell cycle composition of HSPA5-perturbed cells. (F) Mean expression of HSPA5 across subpopulations. Error bars are 95% confidence intervals. (G) Strategy for using correlated expression to identify functionally related genes. (H) Unbiased identification of induced gene expression programs in perturb-seq experiment. 200 example genes induced by perturbation of ER homeostasis were identified computationally, and clustered based on their co-expression across the 65,000 cells in the perturb-seq experiment. Top part of figure shows normalized expression across all of the perturbations in the perturb-seq experiment. Bottom part shows normalized expression in epistasis experiment, to assess UPR dependence. Full-size version of figure in Figure 92A. (I) Implicit correlation information in perturb-seq orders the three branches of the UPR. UPR-responsive genes (from Figure 83F) whose expression were altered in the perturb-seq experiment were clustered using single-cell co-expression data from the epistasis experiment, the perturb-seq experiment, and unperturbed control cells. The figure shows the cophenetic correlation coefficients between dendrogram orderings, measuring how similarly genes cluster, along with a

visual guide to the movement of major groups. Full version of figure in Figure 92B. (J) Strategy for enriching cells perturbed for a trait of interest. E.g., a group of genes of interest may be induced in only a fraction of the population across many different perturbations. (K) Enriching cells strengthens correlations and identifies related genes. Cells enriched for a set of bait cholesterol biosynthesis genes were identified. Within the enriched population, a group of genes clustered with the bait genes and showed stronger correlated expression than in unperturbed control cells. Right panel shows Reactome annotations and SREBP binding data for the group. (L) Combining correlation information with other types of data. An identical analysis to the previous panel was applied to heat shock genes. The figure shows the correlation matrix among the top hits, and compares this to the most enriched transcription factor binding sites measured by Encode. See also Figure 92.

[00143] Figure 87A-F Translocon loss-of-function preferentially activates IRE1 UPR signaling. (A) Single-cell analysis of cells depleted for *SEC61B* in perturb-seq experiment. Panels show the sgRNA identity and IRE1 activation score for each cell. (B) Identical analysis for *SEC61A1*. (C) RT-PCR probing for XBP1 splicing. (D) K562 cells (cBA011) transduced and sorted for expression of the indicated sgRNAs were collected on the indicated days post transduction for analysis of XBP1 mRNA splicing (top) and expression of *SSR2* and *CHOP* mRNA (bottom). Data represent means relative to *ACTB* mRNA and normalized to cells transfected with a negative control sgRNA (dotted lines) \pm standard error of technical replicates ($n = 3$). (E) K562 cells (cBA010) transduced with the indicated sgRNAs and selected with puromycin were treated with 0.5 μ M thapsigargin for 1.5 hours or left untreated 6 days post transduction. Cells were then collected for analysis of XBP1 mRNA splicing and expression of *CHOP* mRNA. Data represent means relative to *ACTB* mRNA and normalized to cells transfected with NegCtrl-3 \pm standard error of technical replicates ($n = 3$). (F) A model in which IRE1 α actively monitors the function and number of translocons and acts to increase them as needed. See also Figure 93.

[00144] Figure 88A-F Design and characterization of three-guide vector (related to Figure 82). (A) Characterization of initial three-guide vector by GFP knockdown. GFP⁺ K562 cells with dCas9-KRAB were transduced with either the single perturb-seq vector expressing a GFP-targeting guide RNA, initial three-guide vectors expressing a GFP-targeting guide RNA from the promoter indicated in parentheses and negative control guide RNAs from the other two

promoters, or a single vector expressing a negative control guide RNA. GFP levels were measured by flow cytometry 10 d after transduction. Plotted are kernel density estimates of normalized flow cytometry counts for infected (BFP+) cells. Traces for the single perturb-seq construct and the negative control are the same as in Figure 82E. (B) Characterization of h7SK promoter in the context of the perturb-seq vector. Experiment was conducted as described in (A). Traces for the single perturb-seq construct and the negative control are the same as in Figure 82E. (C) Characterization of GFP+ K562 cells with increased dCas9-KRAB levels. BFP levels report on expression level of the dCas9-BFP-KRAB construct. Increase in dCas9-KRAB is measured by change in BFP fluorescence relative to WT K562 cells. Plotted are kernel density estimates of normalized flow cytometry counts. (D) Crystal structure of Cas9 bound to guide RNA and target DNA (PDB ID code 4OO8 (Nishimasu et al. 2014)), highlighting location of constant region mutations. Cas9 is shown in gray, target ssDNA in yellow, and the guide RNA in orange (targeting region) and cyan (constant region). Constant region bases that were mutated are highlighted in red. (E) Characterization of RNA polymerase III promoters from different mammalian species by GFP knockdown. GFP+ K562 cells with dCas9-KRAB were transduced with vectors expressing a GFP-targeting guide RNA from the different promoters, in the context of the perturb-seq vector. GFP levels were measured by flow cytometry either 9 d (experiment 1) or 8 d after transduction (experiment 2). % knockdown was calculated after subtracting GFP levels of WT K562 and calculating GFP levels relative to GFP+ K562 cells transduced with a negative control vector. Abbreviations: bU6, bovine U6; sU6, sheep U6; buU6, buffalo U6; pU6, pig U6. (F) Cloning strategy for final three-guide vector. In step 1, protospacers are ligated into the individual backbones. In step 2, three guide RNA expression cassettes are amplified by PCR and inserted into the perturb-seq backbone in a single reaction by four-piece Gibson assembly to obtain the final barcoded three-guide vector.

[00145] **Figure 89A-B** Perturb-seq analytical pipeline (related to Figure 83). (A) Schematic of the analytical pipeline used in the paper. Each step is explained in the Methods, and each single-cell figure has a dedicated section in the Methods describing its construction. (B) Example analysis of thapsigargin-treated cells, related to Figure 83B. The left panels show t-sne projections of the whole population derived using all differentially expressed genes, as described in the Methods. The middle panels show the 16 independent components found by low rank ICA overlaid on the t-sne plot. The right panels show the average values of the four components

identified as varying by perturbation within each of the subpopulations, and the average values of the four components identified as varying through the cell cycle in each cell cycle phase. Further details are in the Methods.

[00146] **Figure 90A-D** CRISPRi screens used to select UPR-modulating sgRNAs for perturb-seq (related to Figure 84). (A) K562 cells (cBA011) were treated with indicated concentrations of tunicamycin in 0.16% DMSO or DMSO alone. Cells were evaluated by flow cytometry or collected for analysis of XBP1 mRNA splicing at indicated times. Data represent average median fluorescence \pm SD (n=3) normalized to DMSO. (B) Comparison of gene phenotypes from the hCRISPRi-v1 and hCRISPRi-v2 screens. Genes chosen for analysis on the perturb-seq platform (83) are indicated in red. (C) Comparison of UPRE reporter gene phenotypes from the hCRISPRi-v2 with gene growth phenotypes from a previously reported genome-scale hCRISPRi-v2 screen (27661255). Select hits are indicated in red. (D) Top eleven annotated functional clusters from DAVID enrichment analysis. Representative names were chosen for each cluster.

[00147] **Figure 91A-F** Perturb-seq screen performance (related to Figure 85). (A) Similarity of phenotypes between guides targeting the same gene. Average expression profiles were created for each sgRNA-containing subpopulation, and hierarchically clustered. Guides targeting a common gene are indicated by color. (B) Shift in sgRNA target expression upon depletion. The distribution of expression of each targeted gene is compared between control cells (containing the NegCtrl2 guide) and each sgRNA-containing subpopulation. sgRNAs are ordered by target expression. (C) Homogeneity of knockdown. Applicants computationally separated each sgRNA-containing subpopulation into top- and bottom-third most perturbed cells based on the deviation of their RNA-seq profiles from the distribution of expression seen in control cells (Methods). The plot shows the average difference in percentage knockdown between these two subpopulations for each sgRNA (gray dot), along with a kernel density estimate of the distribution (black). (D) Expression of UPR genes in perturb-seq experiment. The plot shows the average normalized expression within each perturbed subpopulation of all of the genes identified as UPR-responsive in Figure 83F. The thapsigargin data from that figure is repeated to the right for comparison. (E) Alternate scoring system for branch activation. An alternative method of scoring branch activation was developed using independent component analysis. For comparison with Figure 85D. See Methods for details. (F) Performance of random forest scoring system for

branch activation. To cross-validate branch activation scores, three alternative scoring systems were prototyped in the epistasis experiment of Figure 83, where branch activation can be inferred from perturbation identity: (1) group score, based on expression of defined gene lists, (2) ICA score, seen in Figure 91E, (3) random forest score, seen in Figure 85D. Cells were defined as “active” for a given branch if they were treated with tunicamycin or thapsigargin, and that branch of the UPR was not depleted. The plots show the distributions of scores in active and inactive cells for each of the three branches and each of the three scoring systems. See Methods for details.

[00148] **Figure 92A-B** Functionally clustering genes using single-cell correlation information (related to Figure 86). (A) Full-size version of Figure 86H. (B) Full-size version of Figure 86I.

[00149] **Figure 93A-C** Depletion of individual translocon components SEC61A1, SEC61B, or SEC61G upregulate expression of complex partner genes but have distinct growth phenotypes (related to Figure 87). (A) ATF6 and PERK branch activation scores for *SEC61B*- and *SEC61A1*-perturbed subpopulations. For comparison with Figure 87A. (B) K562 cells (cBA011) transduced and sorted for expression of the indicated sgRNAs were collected 6 days post transduction for analysis of *SEC61A1*, *SEC61B*, *SEC61G*, and *ALG2* expression. Data represent means relative to *ACTB* mRNA and normalized to cells transfected with NegCtrl-3 \pm standard error of technical replicates ($n = 3$). Data from the cells with NegCtrl-1, SEC61A1-2, and SEC61B-2 sgRNAs are represented here and in Figure 87D. (C) Phenotypes for individual sgRNAs targeting *SEC61A1*, *SEC61B*, and *SEC61G* from two growth screens, one reported elsewhere (27661255) and the other conducted in parallel to the first with essentially the same conditions. Data for 10 library negative control sgRNAs were randomly chosen for inclusion. Guides used separately elsewhere are numbered.

DETAILED DESCRIPTION OF THE INVENTION

[00150] Functional genomics has transformed our ability to identify the components of molecular circuits in cells. However, despite over a decade of intensive research, Applicants are still incredibly limited in our ability to predict how perturbing one or more components will affect the cell’s phenotype. For example, the recent explosion in studies on the genetic basis of human disease, from genome-wide association studies (Manolio et al., 2009; Shendure and Akey, 2015; Visscher et al., 2012) to rare Mendelian disease (Altshuler et al., 2008; Bamshad et

al., 2011; Botstein and Risch, 2003) to cancer genetics (Melnikov et al., 2012) (Lawrence et al., 2014; Martincorena and Campbell, 2015), have identified many disease-associated genes, but the mechanisms by which they contribute to pathogenesis remain largely obscure. It is increasingly clear that a lack of biological understanding of the complex nature of molecular networks severely constrains our progress towards therapeutic intervention (Sawyers, 2004).

[00151] Genomic research on dissecting molecular systems in mammalian cells has proceeded along two main paths. Observational studies use global analysis of molecular profiles in different cells to draw correlational hypotheses that relate regulators to their targets. Conversely, genetic perturbation screens, with RNAi, and more recently CRISPR/Cas9-based genome editing, have opened the way for comprehensive analysis of gene function in mammalian cells *in vitro* and *in vivo*. Such screens are designed in either: (1) an individual (“arrayed”) format, where each perturbation is delivered and its effect is measured separately, or in (2) a pooled format, where the effects of perturbations are measured *en masse*. Pooled phenotypic readouts rely on measuring cell autonomous effects, such growth, drug resistance (Bandyopadhyay et al., 2010; Bassik et al., 2013a; Kampmann et al., 2014; Shalem et al., 2014; Wang et al., 2014), or the expression of a single gene (Melnikov et al., 2012; Parnas et al., 2015; Rajagopal et al., 2016; Smith et al., 2013). In general, pooled screens are more efficient, scalable, and less prone to batch effects than arrayed screens, but have thus far been limited to lower-content readouts. As a result, they are unable to discern between multiple molecular mechanisms that result in similar downstream phenotypes. Such distinctions require a more comprehensive readout, such as the transcriptional profile of a cell. Typically, detailed profiling of ‘hits’ from screens requires individual testing, in a time and labor intensive manner (Parnas et al., 2015).

[00152] Transcriptional profiles are a rich readout of the cell’s molecular state, but have been challenging to assess in a screening scale. A few studies tested the transcriptional profiles in a hundred or more individual perturbation experiments, in either follow up studies of pre-validated hits, or in model organisms, such as yeast, where the most recent effort assessed ~1,500 KO strains individually. Indeed, even signature assays (Amit et al., Science. 2009 9;326(5950):257-63) were only measured in large screens in the context of highly resourced centralized efforts (LINCS program). This “readout barrier” is common to all genetic screening, *in vitro* or *in vivo*, from unicellular model organisms to mice or human cells.

[00153] Indeed, it has been challenging to bridge the gap between the use of transcriptional or other rich profiles and the scale of pooled screens. In mammalian cells, a few studies tested transcriptional profiles in a hundred or more individual perturbation experiments, in either follow up studies of pre-validated hits (Berger et al., 2016; Parnas et al., 2015), or in model organisms, such as yeast (Hughes et al., 2000), where a recent effort assessed ~1,500 knock out (KO) strains individually (Kemmeren et al., 2014). Indeed, even signature assays (Amit et al., 2009) were only performed in large screens in the context of highly resourced centralized efforts (Duan et al., 2014; Lamb et al., 2006). This “readout barrier” is common to all genetic screening, *in vitro* or *in vivo*, from unicellular model organisms to mice or human cells.

[00154]

[00155] Notably, a rich readout – such as a genomic profile – would be particularly important to understand the combined effect of multiple factors, in systems where this effect cannot simply be predicted by the sum of their individual effects. Comprehensive analysis of the effects of pairs of genes on cell viability (“synthetic lethality”) has been performed in yeast (Boone et al., 2007; Costanzo et al., 2010; Tong, 2004). In mammals, the effects of interactions on cell viability (Bassik et al., 2013b; Wong et al., 2016) or cell morphology characteristics (Laufer et al., 2013) were assessed for a small number of selected pairs of genes. Few studies have examined the relationship between a high dimensional phenotype, like gene expression, and pairwise interactions. One study in yeast that explored the marginal transcriptional effect associated with double knockouts of a few transcription factors compared to their respective single knockouts (Capaldi et al., 2008), highlighted the potential power of the approach to determine how a regulatory circuit is wired and functions. Very few studies have examined higher order interactions beyond a limited handful (Haber et al., 2013; Tsumura et al., 2006). To the best of our knowledge, there have been few if any studies coupling a high content readout with scalable generation of cells containing multiple perturbations.

[00156] It has been suggested, by the Applicants (Macosko et al., 2015) and others (Liberali et al., 2014; Macosko et al., 2015; Shendure and Fields, 2016) that scalable methods of transcriptionally assaying perturbations using single cell profiling, if developed, could provide an important new tool in the understanding the functions of genes and circuits.

[00157] Here, applicants address the challenge of scaling the phenotypic measurements in genetic screens by developing Perturb-Seq and provide this ability by developing Perturb-seq,

which combines pooled CRISPR screening with a massively parallel single cell RNA-seq readout. Applicants leverage the modularity of sgRNAs in the CRISPR/Cas9 system to perform pooled, multi-locus gene perturbation in mammalian cells (Chen et al., 2015; Cong et al., 2013; Gilbert et al., 2014, 2013; Konermann et al., 2014; Qi et al., 2013; Ran et al., 2015; Zalatan et al., 2015), including in primary cells (Platt et al., 2014). Applicants rely on the scale and single cell nature of massively parallel single cell RNA-seq (Fan et al., 2015; Klein et al., 2015; Macosko et al., 2015) to provide a rich, genome-scale cell autonomous readout for up to millions of cells in a day's work (with appropriate parallelization (Joensson and Andersson Svahn, 2012)). Perturb-Seq uses a CRISPR lentiviral vector, which both delivers an sgRNA to a cell, and reports on the identity of the delivered sgRNA by an expressed barcode on a polyadenylated transcript, captured efficiently by scRNA-seq. Perturb-Seq can be used to assess the effect of single gene perturbations. By increasing the multiplicity of infection, a greater fraction of cells possesses more than one guide, and the approach is more powered to test for epistatic effects. Applicants further developed an integrated computational framework to decipher the effect of individual perturbations and the marginal contributions of pairwise interactions on the level of each expressed gene, gene modules, and global cell states and types.

[00158] Applicants demonstrate the power of Perturb-seq in the context of two distinct biological systems: perturbing key transcription factors (TFs) implicated in cell differentiation and the immune response in bone marrow derived dendritic cells (BMDCs), post-mitotic primary innate immune cells responding to the pathogen component lipopolysaccharide (LPS); and perturbing TFs or(MOI) cell cycle regulators in K562 cells, a rapidly dividing erythroleukemia cell line. Applicants further develop an integrated computational framework to decipher the effect of individual perturbations and the marginal contributions of pairwise interactions on the level of each transcript, transcriptional module, and global cell states. Our framework can be extended to include other high dimensional molecular phenotypes as they become available, as well as incorporate a variety of cell metadata such as lineage information.

[00159] The present invention provides tools and methods for the systematic analysis of genetic interactions, including higher order interactions.

[00160] The present invention provides tools and methods for combinatorial probing of cellular circuits, for dissecting cellular circuitry, for delineating molecular pathways, and/or for identifying relevant targets for therapeutics development.

[00161] The present invention provides tools and methods for pooled screening of perturbations and genome scale readouts in single cells, thus allowing relevant phenotypes to be correlated to specific perturbations.

[00162] In one aspect, the invention provides a method for determining genetic interactions. This method involves causing a set of P genetic perturbations in cells, wherein the method may comprise: determining, based upon random sampling, a subset of π genetic perturbations from the set of P genetic perturbations; performing said subset of π genetic perturbations in a population of cells; performing single-cell molecular profiling of the population of genetically perturbed cells of step; inferring single-cell molecular profiles for the set of P genetic perturbations in cells.

[00163] The population of cells with a plurality of genomic sequence or perturbation conditions involves a plurality of cells and perturbations to be tested and measurements sampled to obtain meaningful data and to infer appropriate circuits. The number of genes perturbed, and how many are perturbed simultaneously (the order of the perturbation, pairs, triplets, etc.) varies. In a tissue with n cell types, the rarest present in $m\%$, how many cells X do you need to sequence so that you have at least Y of the rarest subtype.

[00164] For example, ~ 500 cells ensures $\geq 95\%$ chance of including ≥ 10 of each cell type, based on the following calculation (Fig. 34). Assume the most conservative scenario that of M cell subtypes (for example, 12), all but one having the lowest predicted proportion (for example, $p_{min}=5\%$). Assuming that the Central Limit Theorem holds (a reasonable assumption when solving to detect at least 10 cells of each type) the number of cells of each type i , termed T_i , will distribute as $E[T_i] = N * p_{min}$, $STDV[T_i] = \sqrt{N * p_{min} * (1 - p_{min})}$. The minimal N (total number of cells to profile) can be solved such that all $(m-1)$ subtypes have at least n cells (the last, majority, subtype easily clears this threshold since its proportion is much higher). Applicants confirmed with simulation that the strategy conservatively holds in practice even for $n < 10$, and take a margin of additional (conservative) error, to allow for subsequent failed RNA-Seq experiments ($< 20-30\%$, depending on protocol).

Modelling genetic interactions

[00165] The method of the invention may be used for determining genetic interactions, including modelling and/or analyzing such interactions. Such genetic interactions form part of

cellular circuitry, in that the interactions reflect connections of components within one or more cellular pathways. Such pathways may be intracellular pathways or intercellular pathways.

[00166] In some embodiments, the method of the invention may further comprise determining genetic interactions.

[00167] In some embodiments, the method of the invention may further comprise confirming genetic interactions with additional genetic manipulations.

[00168] The method may further comprise a validation step, wherein additional manipulations are performed in order to confirm previously identified genetic interactions. Such validation step may include in vivo or in vitro experiments, such as gene inactivation, gene deletion, gene activation or overexpression, and combinations thereof. Such genetic manipulations may be performed with any genetic tool available in the art, comprising but not limited to RNAi, CRISPR-Cas based gene editing, nucleic acid transfection, etc.

Genetic perturbations

[00169] In one aspect, said set of P genetic perturbations or said subset of π genetic perturbations may comprise single-order genetic perturbations. Within the meaning of the present invention, single-order genetic perturbation means that a given cell undergoes a single genetic perturbation (one perturbation per cell).

[00170] In one aspect, said set of P genetic perturbations or said subset of π genetic perturbations may comprise combinatorial genetic perturbations. Within the meaning of the present invention, combinatorial or higher-order genetic perturbation means that a given cell undergoes a combination of k single-order genetic perturbations (k perturbations per cell), with $k > 1$. In some embodiments, k is an integer ranging from 2 to 15. In some embodiments, $k = 2, 3, 4, 5, 6, 7, 8, 9$ or 10 .

[00171] Within the meaning of the present invention, said genetic perturbation may comprise gene knock-down (gene repression or gene inactivation), gene knock-out (gene deletion), gene activation, gene insertion, or regulatory element deletion.

[00172] Combinations of different types of genetic perturbations are also envisioned within the meaning of the present invention. For example, a combination of genetic perturbations may comprise a knock-down for a first gene, combined to an activation of a second gene, etc.

[00173] In one aspect, said set of P genetic perturbations or said subset of π genetic perturbations may comprise genome-wide perturbations. Genome-wide perturbations are genetic

perturbations that affect loci across the entire genome. Genome-wide perturbation may include single perturbations of >100, >200, >500, >1,000, >2,500, >5,000, >10,000, >15,000 or >20,000 single genomic loci. The present invention encompasses k -order combinations of genome-wide perturbation.

[00174] In some embodiments, the method may comprise determining k -order genetic interactions.

[00175] In some embodiments, said set of P genetic perturbations may comprise combinatorial genetic perturbations, such as k -order combinations of single-order genetic perturbations, wherein k is an integer ranging from 2 to 15, and step (e) may comprise determining j -order genetic interactions, with $j < k$. Such embodiments rely on sampling higher-order interactions in order to more efficiently infer lower order ones. Given a limited number of possible assays, one is more powered to determine lower order interactions (e.g., 2-, 3-way) from measuring higher order interactions (e.g., 5-way) than from allotting all assays to the lower order, because any higher order interaction carries some information about all interaction terms up to that order (e.g., in compressed sensing, it informs in convolved form on additional Fourier coefficients). Thus, even if most interactions are low order (2- or 3-way) these embodiments are more powered to detect them.

CRISPR-Cas systems

[00176] In some embodiments, RNAi- or CRISPR-Cas-based perturbation may be performed. Said perturbation may be performed (e.g. "delivered") in an array-format or pool-format. Some embodiments may comprise pooled single or combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs, wherein each sgRNA comprises a unique molecular identifier. In some embodiments, a step may comprise pooled combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs, wherein each sgRNA comprises a unique molecular identifier and is co-delivered with a reporter mRNA.

[00177] CRISPR-Cas systems, including CRISPR-Cas9 systems, as used herein, refer to non-naturally occurring systems derived from bacterial Clustered Regularly Interspaced Short Palindromic Repeats loci. These systems generally comprise an enzyme (Cas protein, such as Cas9 protein) and one or more RNAs. Said RNA is a CRISPR RNA and may be an sgRNA. Said RNA and/or said enzyme may be engineered, for example for optimal use in mammalian cells, for optimal delivery therein, for optimal activity therein, for specific uses in gene editing, etc.

[00178] sgRNA refers to a CRISPR single-guide RNA. This RNA is a component of a CRISPR-Cas system. The sequence of the sgRNA determines the target sequence for gene editing, knock-down, knock-out, insertion, etc. For genome-wide approaches, it is possible to design and construct suitable sgRNA libraries. Such sgRNAs may be delivered to cells using vector delivery such as viral vector delivery. Combination of CRISPR-Cas-mediated perturbations may be obtained by delivering multiple sgRNAs within a single cell. This may be achieved in pooled format. In the case of sgRNA viral vector delivery, combined perturbation may be obtained by delivering several sgRNA vectors to the same cell. This may also be achieved in pooled format, and number of combined perturbations in a cell then corresponds to the MOI (multiplicity of infection). Using CRISPR-Cas systems, one may generally implement MOI values of up to 10, 12 or 15.

[00179] The CRISPR-Cas system may be implemented in order to cause massively combinatorial molecular perturbations (MCMP), including single-order and combinatorial genome-wide genetic perturbations.

[00180] CRISPR-Cas-based gene editing allows to perform pooled genome-scale screens with expression readouts in primary cells (A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks. Parnas O., Jovanovic M., Eisenhaure TM., Herbst RH., Dixit A., Ye CJ., Przybylski D., Platt RJ., Tirosch I., Sanjana NE., Shalem O., Satija R., Raychowdhury R., Mertins P., Carr SA., Zhang F., Hacohen N., Regev A. A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks. *Cell* Jul 15. (2015) 2015 Jul 30;162(3):675-86. doi: 10.1016/j.cell.2015.06.059. Epub 2015 Jul 16).

[00181] In some embodiments, the present invention involves combinatorial perturbations by way of CRISPR-Cas (such as CRISPR-Cas9) assays. In accordance with the present invention, sampling a far-from-exhaustive number of higher order perturbations, when coupled with complex genomic readouts, may suffice to resolve most non-linear relations. Accordingly, in some aspects, the present invention relies on pooled, combinatorial perturbations with genomic readout into Massively Combinatorial Perturbation Profiling (MCP).

[00182] In some embodiments, the method of the invention may comprise one or more CRISPR-Cas-based assays. Such CRISPR-Cas assays are advantageous for implementing a precise perturbation of genes and their expression levels.

[00183] In some embodiments, CRISPR-Cas systems may be used to knockout protein-coding genes by frameshifts (indels). Embodiments include efficient and specific CRISPR-Cas9 mediated knockout (Gilbert, L. A., Horlbeck, M. A., Adamson, B., Villalta, J. E., Chen, Y., Whitehead, E. H., Guimaraes, C., Panning, B., Ploegh, H. L., Bassik, M. C., Qi, L. S., Kampmann, M. & Weissman, J. S. Genome-Scale CRISPR-Mediated Control of Gene Repression and Activation. *Cell*. 159, 647-661, doi:10.1016/j.cell.2014.09.029 (2014). PMID:4253859; Ran, F. A., Cong, L., Yan, W. X., Scott, D. A., Gootenberg, J. S., Kriz, A. J., Zetsche, B., Shalem, O., Wu, X., Makarova, K. S., Koonin, E. V., Sharp, P. A. & Zhang, F. In vivo genome editing using *Staphylococcus aureus* Cas9. *Nature*. 520, 186-191, doi:10.1038/nature14299 (2015). PMID:4393360), including a CRISPR mediated double-nicking to efficiently modify both alleles of a target gene or multiple target loci (Ran, F. A., Hsu, P. D., Lin, C. Y., Gootenberg, J. S., Konermann, S., Trevino, A. E., Scott, D. A., Inoue, A., Matoba, S., Zhang, Y. & Zhang, F. Double nicking by RNA-guided CRISPR Cas9 for enhanced genome editing specificity. *Cell*. 154, 1380-1389, doi:10.1016/j.cell.2013.08.021 (2013). PMID:3856256; Wang, H., Yang, H., Shivalila, C. S., Dawlaty, M. M., Cheng, A. W., Zhang, F. & Jaenisch, R. One-step generation of mice carrying mutations in multiple genes by CRISPR-Cas-mediated genome engineering. *Cell*. 153, 910-918, doi:10.1016/j.cell.2013.04.025 (2013). PMID:3969854) and implementation of a smaller Cas9 protein for delivery on smaller vectors (Ran, F. A., Cong, L., Yan, W. X., Scott, D. A., Gootenberg, J. S., Kriz, A. J., Zetsche, B., Shalem, O., Wu, X., Makarova, K. S., Koonin, E. V., Sharp, P. A. & Zhang, F. In vivo genome editing using *Staphylococcus aureus* Cas9. *Nature*. 520, 186-191, doi:10.1038/nature14299 (2015). PMID:4393360).

[00184] CRISPR-mediated activation or inactivation (CRISPRa/i) systems may be used to activate or inactivate gene transcription. Briefly, a nuclease-dead (deactivated) Cas9 RNA-guided DNA binding domain (dCas9) (Qi, L. S., Larson, M. H., Gilbert, L. A., Doudna, J. A., Weissman, J. S., Arkin, A. P. & Lim, W. A. Repurposing CRISPR as an RNA-guided platform for sequence-specific control of gene expression. *Cell*. 152, 1173-1183, doi:10.1016/j.cell.2013.02.022 (2013). PMID:3664290) tethered to transcriptional repressor domains that promote epigenetic silencing (e.g., KRAB) forms a "CRISPRi" (Gilbert, L. A., Larson, M. H., Morsut, L., Liu, Z., Brar, G. A., Torres, S. E., Stern-Ginossar, N., Brandman, O., Whitehead, E. H., Doudna, J. A., Lim, W. A., Weissman, J. S. & Qi, L. S. CRISPR-mediated

modular RNA-guided regulation of transcription in eukaryotes. *Cell*. 154, 442-451, doi:10.1016/j.cell.2013.06.044 (2013). PMID:3770145; Konermann, S., Brigham, M. D., Trevino, A. E., Hsu, P. D., Heidenreich, M., Cong, L., Platt, R. J., Scott, D. A., Church, G. M. & Zhang, F. Optical control of mammalian endogenous transcription and epigenetic states. *Nature*. 500, 472-476, doi:10.1038/nature12466 (2013). PMID:3856241) that represses transcription. To use dCas9 as an activator (CRISPRa), a guide RNA may be engineered to carry RNA binding motifs (e.g., MS2) that recruit effector domains fused to RNA-motif binding proteins, increasing transcription (Konermann, S., Brigham, M. D., Trevino, A. E., Joung, J., Abudayyeh, O. O., Barcena, C., Hsu, P. D., Habib, N., Gootenberg, J. S., Nishimasu, H., Nureki, O. & Zhang, F. Genome-scale transcriptional activation by an engineered CRISPR-Cas9 complex. *Nature*. 517, 583-588, doi:10.1038/nature14136 (2015). PMID:4420636).

[00185] CRISPR-Cas systems may also be used for the deletion of regulatory elements. To target non-coding elements, pairs of guides may be designed and used to delete regions of a defined size, and tile deletions covering sets of regions in pools. The delivery of two sgRNAs may mediate efficient excision of 500 bp genomic fragments.

[00186] CRISPR-Cas systems may also be used for gene editing, e.g. by RNA-templated homologous recombination. Keskin, H., Shen, Y., Huang, F., Patel, M., Yang, T., Ashley, K., Mazin, A. V. & Storici, F. Transcript-RNA-templated DNA recombination and repair. *Nature*. 515, 436-439, doi:10.1038/nature13682 (2014).

[00187] CRISPR transgenic mice may be used to derive 'CRISPR-ready' cells. 'CRISPR-mice' are mice where the mouse germ line is engineered to harbor key elements of a CRISPR system, and cells require only the programmable (sgRNA) element to activate the CRISPR-Cas system. CRISPR mice include Cas9-transgenic mice (Platt, R. J., Chen, S., Zhou, Y., Yim, M. J., Swiech, L., Kempton, H. R., Dahlman, J. E., Parnas, O., Eisenhaure, T. M., Jovanovic, M., Graham, D. B., Jhunjhunwala, S., Heidenreich, M., Xavier, R. J., Langer, R., Anderson, D. G., Hacohen, N., Regev, A., Feng, G., Sharp, P. A. & Zhang, F. CRISPR-Cas9 knockin mice for genome editing and cancer modeling. *Cell*. 159, 440-455, doi:10.1016/j.cell.2014.09.014 (2014). PMID:4265475; Parnas O., Jovanovic M., Eisenhaure TM., Herbst RH., Dixit A., Ye CJ., Przybylski D., Platt RJ., Tirosh I., Sanjana NE., Shalem O., Satija R., Raychowdhury R., Mertins P., Carr SA., Zhang F., Hacohen N., Regev A. A Genome-wide CRISPR Screen in Primary

Immune Cells to Dissect Regulatory Networks. *Cell* Jul 15. (2015) 2015 Jul 30;162(3):675-86. doi: 10.1016/j.cell.2015.06.059. Epub 2015 Jul 16).

[00188] CRISPR-Cas based perturbations, including single order or higher order perturbations, may be implemented in pooled format. The perturbation (screen) may be performed with expression readouts or reporter expression readout (genome-wide reporter-based pooled screens).

[00189] CRISPR-Cas functional genomics assays that may be used to cause sets of genetic perturbations are described in Shalem O., Sanjana NE., Zhang F. High-throughput functional genomics using CRISPR-Cas9. *Nat Rev Genet.* May;16(5):299-311. (2015). doi: 10.1038/nrg3899. Epub 2015 Apr 9.

[00190] sgRNA libraries, including genome-wide libraries of sgRNAs, may be designed as described in Parnas O., Jovanovic M., Eisenhaure TM., Herbst RH., Dixit A., Ye CJ., Przybylski D., Platt RJ., Tirosh I., Sanjana NE., Shalem O., Satija R., Raychowdhury R., Mertins P., Carr SA., Zhang F., Hacohen N., Regev A. A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks. *Cell* Jul 15. (2015) 2015 Jul 30;162(3):675-86. doi: 10.1016/j.cell.2015.06.059. Epub 2015 Jul 16; Sanjana, N. E., Shalem, O. & Zhang, F. Improved vectors and genome-wide libraries for CRISPR screening. *Nat Methods.* 11, 783-784, doi:10.1038/nmeth.3047 (2014); Shalem, O., Sanjana, N. E., Hartenian, E., Shi, X., Scott, D. A., Mikkelsen, T. S., Heckl, D., Ebert, B. L., Root, D. E., Doench, J. G. & Zhang, F. Genome-scale CRISPR-Cas9 knockout screening in human cells. *Science.* 343, 84-87, doi:10.1126/science.1247005 (2014). PMID:4089965; Shalem, O., Sanjana, N. E. & Zhang, F. High-throughput functional genomics using CRISPR-Cas9. *Nat Rev Genet.* 16, 299-311, doi:10.1038/nrg3899 (2015).

[00191] A pooled genome-wide screen for CRISPR-mediated KO (knock-out) may be performed as in Shalem, O., Sanjana, N. E., Hartenian, E., Shi, X., Scott, D. A., Mikkelsen, T. S., Heckl, D., Ebert, B. L., Root, D. E., Doench, J. G. & Zhang, F. Genome-scale CRISPR-Cas9 knockout screening in human cells. *Science.* 343, 84-87, doi:10.1126/science.1247005 (2014). PMID:4089965.

[00192] An expression marker-based genome-wide CRISPR screen may be performed as in Parnas O., Jovanovic M., Eisenhaure TM., Herbst RH., Dixit A., Ye CJ., Przybylski D., Platt RJ., Tirosh I., Sanjana NE., Shalem O., Satija R., Raychowdhury R., Mertins P., Carr SA., Zhang F.,

Hacohen N., Regev A. A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks. *Cell* Jul 15. (2015) 2015 Jul 30;162(3):675-86. doi: 10.1016/j.cell.2015.06.059. Epub 2015 Jul 16.

[00193] A pooled, genome-scale, CRISPRa screen may be performed as in Konermann, S., Brigham, M. D., Trevino, A. E., Joung, J., Abudayyeh, O. O., Barcena, C., Hsu, P. D., Habib, N., Gootenberg, J. S., Nishimasu, H., Nureki, O. & Zhang, F. Genome-scale transcriptional activation by an engineered CRISPR-Cas9 complex. *Nature*. 517, 583-588, doi:10.1038/nature14136 (2015). PMID:4420636.

[00194] Pooled combinatorial perturbations may be performed, where the delivered perturbations and impact (molecular profiling) are determined post hoc, in either a conventional readout (e.g., sorting followed by sequencing) or with high-content single cell genomics.

[00195] In some embodiment, the CRISPR-Cas screen is performed by co-delivering multiple sgRNA using virale vector delivery (eg, sgRNA encoding vectors at a relatively high MOI) into cells pre-expressing the Cas9 enzyme to obtain as many higher order combinations as possible. For small sets of ~5 genes one may generate a combinatorially complete ascertained set of all 32 perturbations.

[00196] To detect which perturbations were co-delivered in pooled bins, several strategies are envisioned: (1) Combining two or more guide barcodes using in situ PCR in PEG hydrogel that restricts the diffusion of double stranded DNA; (2) Split-pool tagging of guide barcodes in hydrogels, such that only guides from the same cell are tagged with the same sequence; (3) FISH of expressed guides for imaging readouts. In each case it is possible to use an error-correction scheme in the barcodes.

[00197] To detect which perturbations were co-delivered with a single cell genomics readout, it is possible to report the (combinatorial) perturbation in a manner compatible with the full genomic readout. For example one may use an sgRNA vector that also highly expresses a synthetic polyadenylated RNA reporter of the sgRNA barcode. This RNA will be captured along with the cellular mRNA in the transcriptome profiling, eg scRNA-seq (Drop-Seq, see below), or reported by FISH hybridization, such that the same assay ascertains the sgRNAs and their impact on expression (Parnas O., Jovanovic M., Eisenhaure TM., Herbst RH., Dixit A., Ye CJ., Przybylski D., Platt RJ., Tirosh I., Sanjana NE., Shalem O., Satija R., Raychowdhury R., Mertins P., Carr SA., Zhang F., Hacohen N., Regev A. A Genome-wide CRISPR Screen in Primary

Immune Cells to Dissect Regulatory Networks. *Cell* Jul 15. (2015) 2015 Jul 30;162(3):675-86. doi: 10.1016/j.cell.2015.06.059. Epub 2015 Jul 16).

[00198] With respect to general information on CRISPR-Cas Systems, components thereof, and delivery of such components, including methods, materials, delivery vehicles, vectors, particles, AAV, and making and using thereof, including as to amounts and formulations, all useful in the practice of the instant invention, reference is made to: US Patents Nos. 8,999,641, 8,993,233, 8,945,839, 8,932,814, 8,906,616, 8,895,308, 8,889,418, 8,889,356, 8,871,445, 8,865,406, 8,795,965, 8,771,945 and 8,697,359; US Patent Publications US 2014-0310830 (US APP. Ser. No. 14/105,031), US 2014-0287938 A1 (U.S. App. Ser. No. 14/213,991), US 2014-0273234 A1 (U.S. App. Ser. No. 14/293,674), US2014-0273232 A1 (U.S. App. Ser. No. 14/290,575), US 2014-0273231 (U.S. App. Ser. No. 14/259,420), US 2014-0256046 A1 (U.S. App. Ser. No. 14/226,274), US 2014-0248702 A1 (U.S. App. Ser. No. 14/258,458), US 2014-0242700 A1 (U.S. App. Ser. No. 14/222,930), US 2014-0242699 A1 (U.S. App. Ser. No. 14/183,512), US 2014-0242664 A1 (U.S. App. Ser. No. 14/104,990), US 2014-0234972 A1 (U.S. App. Ser. No. 14/183,471), US 2014-0227787 A1 (U.S. App. Ser. No. 14/256,912), US 2014-0189896 A1 (U.S. App. Ser. No. 14/105,035), US 2014-0186958 (U.S. App. Ser. No. 14/105,017), US 2014-0186919 A1 (U.S. App. Ser. No. 14/104,977), US 2014-0186843 A1 (U.S. App. Ser. No. 14/104,900), US 2014-0179770 A1 (U.S. App. Ser. No. 14/104,837) and US 2014-0179006 A1 (U.S. App. Ser. No. 14/183,486), US 2014-0170753 (US App Ser No 14/183,429); European Patents EP 2 784 162 B1 and EP 2 771 468 B1; European Patent Applications EP 2 771 468 (EP13818570.7), EP 2 764 103 (EP13824232.6), and EP 2 784 162 (EP14170383.5); and PCT Patent Publications PCT Patent Publications WO 2014/093661 (PCT/US2013/074743), WO 2014/093694 (PCT/US2013/074790), WO 2014/093595 (PCT/US2013/074611), WO 2014/093718 (PCT/US2013/074825), WO 2014/093709 (PCT/US2013/074812), WO 2014/093622 (PCT/US2013/074667), WO 2014/093635 (PCT/US2013/074691), WO 2014/093655 (PCT/US2013/074736), WO 2014/093712 (PCT/US2013/074819), WO2014/093701 (PCT/US2013/074800), WO2014/018423 (PCT/US2013/051418), WO 2014/204723 (PCT/US2014/041790), WO 2014/204724 (PCT/US2014/041800), WO 2014/204725 (PCT/US2014/041803), WO 2014/204726 (PCT/US2014/041804), WO 2014/204727 (PCT/US2014/041806), WO 2014/204728 (PCT/US2014/041808), WO 2014/204729 (PCT/US2014/041809). Reference is also made to US

provisional patent applications 61/758,468; 61/802,174; 61/806,375; 61/814,263; 61/819,803 and 61/828,130, filed on January 30, 2013; March 15, 2013; March 28, 2013; April 20, 2013; May 6, 2013 and May 28, 2013 respectively. Reference is also made to US provisional patent application 61/836,123, filed on June 17, 2013. Reference is additionally made to US provisional patent applications 61/835,931, 61/835,936, 61/836,127, 61/836, 101, 61/836,080 and 61/835,973, each filed June 17, 2013. Further reference is made to US provisional patent applications 61/862,468 and 61/862,355 filed on August 5, 2013; 61/871,301 filed on August 28, 2013; 61/960,777 filed on September 25, 2013 and 61/961,980 filed on October 28, 2013. Reference is yet further made to: PCT Patent applications Nos: PCT/US2014/041803, PCT/US2014/041800, PCT/US2014/041809, PCT/US2014/041804 and PCT/US2014/041806, each filed June 10, 2014 6/10/14; PCT/US2014/041808 filed June 11, 2014; and PCT/US2014/62558 filed October 28, 2014, and US Provisional Patent Applications Serial Nos.: 61/915,150, 61/915,301, 61/915,267 and 61/915,260, each filed December 12, 2013; 61/757,972 and 61/768,959, filed on January 29, 2013 and February 25, 2013; 61/835,936, 61/836,127, 61/836,101, 61/836,080, 61/835,973, and 61/835,931, filed June 17, 2013; 62/010,888 and 62/010,879, both filed June 11, 2014; 62/010,329 and 62/010,441, each filed June 10, 2014; 61/939,228 and 61/939,242, each filed February 12, 2014; 61/980,012, filed April 15, 2014; 62/038,358, filed August 17, 2014; 62/054,490, 62/055,484, 62/055,460 and 62/055,487, each filed September 25, 2014; and 62/069,243, filed October 27, 2014. Reference is also made to US provisional patent applications Nos. 62/055,484, 62/055,460, and 62/055,487, filed September 25, 2014; US provisional patent application 61/980,012, filed April 15, 2014; and US provisional patent application 61/939,242 filed February 12, 2014. Reference is made to PCT application designating, inter alia, the United States, application No. PCT/US14/41806, filed June 10, 2014. Reference is made to US provisional patent application 61/930,214 filed on January 22, 2014. Reference is made to US provisional patent applications 61/915,251; 61/915,260 and 61/915,267, each filed on December 12, 2013. Reference is made to US provisional patent application USSN 61/980,012 filed April 15, 2014. Reference is made to PCT application designating, inter alia, the United States, application No. PCT/US14/41806, filed June 10, 2014. Reference is made to US provisional patent application 61/930,214 filed on January 22, 2014. Reference is made to US provisional patent applications 61/915,251; 61/915,260 and 61/915,267, each filed on December 12, 2013.

[00199] Mention is also made of US application 62/091,455, filed, 12-Dec-14, PROTECTED GUIDE RNAS (PGRNAS); US application 62/096,708, 24-Dec-14, PROTECTED GUIDE RNAS (PGRNAS); US application 62/091,462, 12-Dec-14, DEAD GUIDES FOR CRISPR TRANSCRIPTION FACTORS; US application 62/096,324, 23-Dec-14, DEAD GUIDES FOR CRISPR TRANSCRIPTION FACTORS; US application 62/091,456, 12-Dec-14, ESCORTED AND FUNCTIONALIZED GUIDES FOR CRISPR-CAS SYSTEMS; US application 62/091,461, 12-Dec-14, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR GENOME EDITING AS TO HEMATOPOETIC STEM CELLS (HSCs); US application 62/094,903, 19-Dec-14, UNBIASED IDENTIFICATION OF DOUBLE-STRAND BREAKS AND GENOMIC REARRANGEMENT BY GENOME-WISE INSERT CAPTURE SEQUENCING; US application 62/096,761, 24-Dec-14, ENGINEERING OF SYSTEMS, METHODS AND OPTIMIZED ENZYME AND GUIDE SCAFFOLDS FOR SEQUENCE MANIPULATION; US application 62/098,059, 30-Dec-14, RNA-TARGETING SYSTEM; US application 62/096,656, 24-Dec-14, CRISPR HAVING OR ASSOCIATED WITH DESTABILIZATION DOMAINS; US application 62/096,697, 24-Dec-14, CRISPR HAVING OR ASSOCIATED WITH AAV; US application 62/098,158, 30-Dec-14, ENGINEERED CRISPR COMPLEX INSERTIONAL TARGETING SYSTEMS; US application 62/151,052, 22-Apr-15, CELLULAR TARGETING FOR EXTRACELLULAR EXOSOMAL REPORTING; US application 62/054,490, 24-Sep-14, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR TARGETING DISORDERS AND DISEASES USING PARTICLE DELIVERY COMPONENTS; US application 62/055,484, 25-Sep-14, SYSTEMS, METHODS AND COMPOSITIONS FOR SEQUENCE MANIPULATION WITH OPTIMIZED FUNCTIONAL CRISPR-CAS SYSTEMS; US application 62/087,537, 4-Dec-14, SYSTEMS, METHODS AND COMPOSITIONS FOR SEQUENCE MANIPULATION WITH OPTIMIZED FUNCTIONAL CRISPR-CAS SYSTEMS; US application 62/054,651, 24-Sep-14, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR MODELING COMPETITION OF MULTIPLE CANCER MUTATIONS IN VIVO; US application 62/067,886, 23-Oct-14, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR MODELING COMPETITION OF MULTIPLE CANCER MUTATIONS IN VIVO; US application

62/054,675, 24-Sep-14, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS IN NEURONAL CELLS/TISSUES; US application 62/054,528, 24-Sep-14, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS IN IMMUNE DISEASES OR DISORDERS; US application 62/055,454, 25-Sep-14, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR TARGETING DISORDERS AND DISEASES USING CELL PENETRATION PEPTIDES (CPP); US application 62/055,460, 25-Sep-14, MULTIFUNCTIONAL-CRISPR COMPLEXES AND/OR OPTIMIZED ENZYME LINKED FUNCTIONAL-CRISPR COMPLEXES; US application 62/087,475, 4-Dec-14, FUNCTIONAL SCREENING WITH OPTIMIZED FUNCTIONAL CRISPR-CAS SYSTEMS; US application 62/055,487, 25-Sep-14, FUNCTIONAL SCREENING WITH OPTIMIZED FUNCTIONAL CRISPR-CAS SYSTEMS; US application 62/087,546, 4-Dec-14, MULTIFUNCTIONAL CRISPR COMPLEXES AND/OR OPTIMIZED ENZYME LINKED FUNCTIONAL-CRISPR COMPLEXES; and US application 62/098,285, 30-Dec-14, CRISPR MEDIATED IN VIVO MODELING AND GENETIC SCREENING OF TUMOR GROWTH AND METASTASIS.

[00200] Each of these patents, patent publications, and applications, and all documents cited therein or during their prosecution (“appln cited documents”) and all documents cited or referenced in the appln cited documents, together with any instructions, descriptions, product specifications, and product sheets for any products mentioned therein or in any document therein and incorporated by reference herein, are hereby incorporated herein by reference, and may be employed in the practice of the invention. All documents (e.g., these patents, patent publications and applications and the appln cited documents) are incorporated herein by reference to the same extent as if each individual document was specifically and individually indicated to be incorporated by reference.

[00201] Also with respect to general information on CRISPR-Cas Systems, mention is made of the following (also hereby incorporated herein by reference):

- Multiplex genome engineering using CRISPR/Cas systems. Cong, L., Ran, F.A., Cox, D., Lin, S., Barretto, R., Habib, N., Hsu, P.D., Wu, X., Jiang, W., Marraffini, L.A., & Zhang, F. Science Feb 15;339(6121):819-23 (2013);

- RNA-guided editing of bacterial genomes using CRISPR-Cas systems. Jiang W., Bikard D., Cox D., Zhang F, Marraffini LA. *Nat Biotechnol* Mar;31(3):233-9 (2013);
- One-Step Generation of Mice Carrying Mutations in Multiple Genes by CRISPR/Cas-Mediated Genome Engineering. Wang H., Yang H., Shivalila CS., Dawlaty MM., Cheng AW., Zhang F., Jaenisch R. *Cell* May 9;153(4):910-8 (2013);
- Optical control of mammalian endogenous transcription and epigenetic states. Konermann S, Brigham MD, Trevino AE, Hsu PD, Heidenreich M, Cong L, Platt RJ, Scott DA, Church GM, Zhang F. *Nature*. Aug 22;500(7463):472-6. doi: 10.1038/Nature12466. Epub 2013 Aug 23 (2013);
- Double Nicking by RNA-Guided CRISPR Cas9 for Enhanced Genome Editing Specificity. Ran, FA., Hsu, PD., Lin, CY., Gootenberg, JS., Konermann, S., Trevino, AE., Scott, DA., Inoue, A., Matoba, S., Zhang, Y., & Zhang, F. *Cell* Aug 28. pii: S0092-8674(13)01015-5 (2013-A);
- DNA targeting specificity of RNA-guided Cas9 nucleases. Hsu, P., Scott, D., Weinstein, J., Ran, FA., Konermann, S., Agarwala, V., Li, Y., Fine, E., Wu, X., Shalem, O., Cradick, TJ., Marraffini, LA., Bao, G., & Zhang, F. *Nat Biotechnol* doi:10.1038/nbt.2647 (2013);
- Genome engineering using the CRISPR-Cas9 system. Ran, FA., Hsu, PD., Wright, J., Agarwala, V., Scott, DA., Zhang, F. *Nature Protocols* Nov;8(11):2281-308 (2013-B);
- Genome-Scale CRISPR-Cas9 Knockout Screening in Human Cells. Shalem, O., Sanjana, NE., Hartenian, E., Shi, X., Scott, DA., Mikkelsen, T., Heckl, D., Ebert, BL., Root, DE., Doench, JG., Zhang, F. *Science* Dec 12. (2013). [Epub ahead of print];
- Crystal structure of cas9 in complex with guide RNA and target DNA. Nishimasu, H., Ran, FA., Hsu, PD., Konermann, S., Shehata, SI., Dohmae, N., Ishitani, R., Zhang, F., Nureki, O. *Cell* Feb 27, 156(5):935-49 (2014);
- Genome-wide binding of the CRISPR endonuclease Cas9 in mammalian cells. Wu X., Scott DA., Kriz AJ., Chiu AC., Hsu PD., Dadon DB., Cheng AW., Trevino AE., Konermann S., Chen S., Jaenisch R., Zhang F., Sharp PA. *Nat Biotechnol*. Apr 20. doi: 10.1038/nbt.2889 (2014);
- CRISPR-Cas9 Knockin Mice for Genome Editing and Cancer Modeling. Platt RJ, Chen S, Zhou Y, Yim MJ, Swiech L, Kempton HR, Dahlman JE, Parnas O, Eisenhaure TM, Jovanovic M, Graham DB, Jhunjhunwala S, Heidenreich M, Xavier RJ, Langer R,

- Anderson DG, Hacohen N, Regev A, Feng G, Sharp PA, Zhang F. *Cell* 159(2): 440-455
DOI: 10.1016/j.cell.2014.09.014(2014);
- Development and Applications of CRISPR-Cas9 for Genome Engineering, Hsu PD, Lander ES, Zhang F., *Cell*. Jun 5;157(6):1262-78 (2014).
 - Genetic screens in human cells using the CRISPR/Cas9 system, Wang T, Wei JJ, Sabatini DM, Lander ES., *Science*. January 3; 343(6166): 80–84. doi:10.1126/science.1246981 (2014);
 - Rational design of highly active sgRNAs for CRISPR-Cas9-mediated gene inactivation, Doench JG, Hartenian E, Graham DB, Tothova Z, Hegde M, Smith I, Sullender M, Ebert BL, Xavier RJ, Root DE., (published online 3 September 2014) *Nat Biotechnol*. Dec;32(12):1262-7 (2014);
 - In vivo* interrogation of gene function in the mammalian brain using CRISPR-Cas9, Swiech L, Heidenreich M, Banerjee A, Habib N, Li Y, Trombetta J, Sur M, Zhang F., (published online 19 October 2014) *Nat Biotechnol*. Jan;33(1):102-6 (2015);
 - Genome-scale transcriptional activation by an engineered CRISPR-Cas9 complex, Konermann S, Brigham MD, Trevino AE, Joung J, Abudayyeh OO, Barcena C, Hsu PD, Habib N, Gootenberg JS, Nishimasu H, Nureki O, Zhang F., *Nature*. Jan 29;517(7536):583-8 (2015).
 - A split-Cas9 architecture for inducible genome editing and transcription modulation, Zetsche B, Volz SE, Zhang F., (published online 02 February 2015) *Nat Biotechnol*. Feb;33(2):139-42 (2015);
 - Genome-wide CRISPR Screen in a Mouse Model of Tumor Growth and Metastasis, Chen S, Sanjana NE, Zheng K, Shalem O, Lee K, Shi X, Scott DA, Song J, Pan JQ, Weissleder R, Lee H, Zhang F, Sharp PA. *Cell* 160, 1246–1260, March 12, 2015 (multiplex screen in mouse), and
 - In vivo* genome editing using *Staphylococcus aureus* Cas9, Ran FA, Cong L, Yan WX, Scott DA, Gootenberg JS, Kriz AJ, Zetsche B, Shalem O, Wu X, Makarova KS, Koonin EV, Sharp PA, Zhang F., (published online 01 April 2015), *Nature*. Apr 9;520(7546):186-91 (2015).
 - Shalem et al., "High-throughput functional genomics using CRISPR-Cas9," *Nature Reviews Genetics* 16, 299-311 (May 2015).

- Xu et al., "Sequence determinants of improved CRISPR sgRNA design," *Genome Research* 25, 1147-1157 (August 2015).
- Parnas et al., "A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks," *Cell* 162, 675-686 (July 30, 2015).
- Ramanan et al., "CRISPR/Cas9 cleavage of viral DNA efficiently suppresses hepatitis B virus," *Scientific Reports* 5:10833. doi: 10.1038/srep10833 (June 2, 2015)
- Nishimasu et al., "Crystal Structure of Staphylococcus aureus Cas9," *Cell* 162, 1113-1126 (Aug. 27, 2015)
- Zetsche et al., "Cpf1 Is a Single RNA-Guided Endonuclease of a Class 2 CRISPR-Cas System," *Cell* 163, 1-13 (Oct. 22, 2015)
- Shmakov et al., "Discovery and Functional Characterization of Diverse Class 2 CRISPR-Cas Systems," *Molecular Cell* 60, 1-13 (Available online Oct. 22, 2015)

each of which is incorporated herein by reference, may be considered in the practice of the instant invention, and discussed briefly below:

- Cong et al. engineered type II CRISPR-Cas systems for use in eukaryotic cells based on both *Streptococcus thermophilus* Cas9 and also *Streptococcus pyogenes* Cas9 and demonstrated that Cas9 nucleases can be directed by short RNAs to induce precise cleavage of DNA in human and mouse cells. Their study further showed that Cas9 as converted into a nicking enzyme can be used to facilitate homology-directed repair in eukaryotic cells with minimal mutagenic activity. Additionally, their study demonstrated that multiple guide sequences can be encoded into a single CRISPR array to enable simultaneous editing of several at endogenous genomic loci sites within the mammalian genome, demonstrating easy programmability and wide applicability of the RNA-guided nuclease technology. This ability to use RNA to program sequence specific DNA cleavage in cells defined a new class of genome engineering tools. These studies further showed that other CRISPR loci are likely to be transplantable into mammalian cells and can also mediate mammalian genome cleavage. Importantly, it can be envisaged that several aspects of the CRISPR-Cas system can be further improved to increase its efficiency and versatility.
- Jiang et al. used the clustered, regularly interspaced, short palindromic repeats (CRISPR)-associated Cas9 endonuclease complexed with dual-RNAs to introduce precise mutations

in the genomes of *Streptococcus pneumoniae* and *Escherichia coli*. The approach relied on dual-RNA:Cas9-directed cleavage at the targeted genomic site to kill unmutated cells and circumvents the need for selectable markers or counter-selection systems. The study reported reprogramming dual-RNA:Cas9 specificity by changing the sequence of short CRISPR RNA (crRNA) to make single- and multinucleotide changes carried on editing templates. The study showed that simultaneous use of two crRNAs enabled multiplex mutagenesis. Furthermore, when the approach was used in combination with recombineering, in *S. pneumoniae*, nearly 100% of cells that were recovered using the described approach contained the desired mutation, and in *E. coli*, 65% that were recovered contained the mutation.

- Wang et al. (2013) used the CRISPR/Cas system for the one-step generation of mice carrying mutations in multiple genes which were traditionally generated in multiple steps by sequential recombination in embryonic stem cells and/or time-consuming intercrossing of mice with a single mutation. The CRISPR/Cas system will greatly accelerate the in vivo study of functionally redundant genes and of epistatic gene interactions.
- Konermann et al. (2013) addressed the need in the art for versatile and robust technologies that enable optical and chemical modulation of DNA-binding domains based CRISPR Cas9 enzyme and also Transcriptional Activator Like Effectors
- Ran et al. (2013-A) described an approach that combined a Cas9 nickase mutant with paired guide RNAs to introduce targeted double-strand breaks. This addresses the issue of the Cas9 nuclease from the microbial CRISPR-Cas system being targeted to specific genomic loci by a guide sequence, which can tolerate certain mismatches to the DNA target and thereby promote undesired off-target mutagenesis. Because individual nicks in the genome are repaired with high fidelity, simultaneous nicking via appropriately offset guide RNAs is required for double-stranded breaks and extends the number of specifically recognized bases for target cleavage. The authors demonstrated that using paired nicking can reduce off-target activity by 50- to 1,500-fold in cell lines and to facilitate gene knockout in mouse zygotes without sacrificing on-target cleavage efficiency. This versatile strategy enables a wide variety of genome editing applications that require high specificity.

- Hsu et al. (2013) characterized SpCas9 targeting specificity in human cells to inform the selection of target sites and avoid off-target effects. The study evaluated >700 guide RNA variants and SpCas9-induced indel mutation levels at >100 predicted genomic off-target loci in 293T and 293FT cells. The authors that SpCas9 tolerates mismatches between guide RNA and target DNA at different positions in a sequence-dependent manner, sensitive to the number, position and distribution of mismatches. The authors further showed that SpCas9-mediated cleavage is unaffected by DNA methylation and that the dosage of SpCas9 and sgRNA can be titrated to minimize off-target modification. Additionally, to facilitate mammalian genome engineering applications, the authors reported providing a web-based software tool to guide the selection and validation of target sequences as well as off-target analyses.
- Ran et al. (2013-B) described a set of tools for Cas9-mediated genome editing via non-homologous end joining (NHEJ) or homology-directed repair (HDR) in mammalian cells, as well as generation of modified cell lines for downstream functional studies. To minimize off-target cleavage, the authors further described a double-nicking strategy using the Cas9 nickase mutant with paired guide RNAs. The protocol provided by the authors experimentally derived guidelines for the selection of target sites, evaluation of cleavage efficiency and analysis of off-target activity. The studies showed that beginning with target design, gene modifications can be achieved within as little as 1-2 weeks, and modified clonal cell lines can be derived within 2-3 weeks.
- Shalem et al. described a new way to interrogate gene function on a genome-wide scale. Their studies showed that delivery of a genome-scale CRISPR-Cas9 knockout (GeCKO) library targeted 18,080 genes with 64,751 unique guide sequences enabled both negative and positive selection screening in human cells. First, the authors showed use of the GeCKO library to identify genes essential for cell viability in cancer and pluripotent stem cells. Next, in a melanoma model, the authors screened for genes whose loss is involved in resistance to vemurafenib, a therapeutic that inhibits mutant protein kinase BRAF. Their studies showed that the highest-ranking candidates included previously validated genes NF1 and MED12 as well as novel hits NF2, CUL3, TADA2B, and TADA1. The authors observed a high level of consistency between independent guide RNAs targeting

the same gene and a high rate of hit confirmation, and thus demonstrated the promise of genome-scale screening with Cas9.

- Nishimasu et al. reported the crystal structure of *Streptococcus pyogenes* Cas9 in complex with sgRNA and its target DNA at 2.5 Å resolution. The structure revealed a bilobed architecture composed of target recognition and nuclease lobes, accommodating the sgRNA:DNA heteroduplex in a positively charged groove at their interface. Whereas the recognition lobe is essential for binding sgRNA and DNA, the nuclease lobe contains the HNH and RuvC nuclease domains, which are properly positioned for cleavage of the complementary and non-complementary strands of the target DNA, respectively. The nuclease lobe also contains a carboxyl-terminal domain responsible for the interaction with the protospacer adjacent motif (PAM). This high-resolution structure and accompanying functional analyses have revealed the molecular mechanism of RNA-guided DNA targeting by Cas9, thus paving the way for the rational design of new, versatile genome-editing technologies.
- Wu et al. mapped genome-wide binding sites of a catalytically inactive Cas9 (dCas9) from *Streptococcus pyogenes* loaded with single guide RNAs (sgRNAs) in mouse embryonic stem cells (mESCs). The authors showed that each of the four sgRNAs tested targets dCas9 to between tens and thousands of genomic sites, frequently characterized by a 5-nucleotide seed region in the sgRNA and an NGG protospacer adjacent motif (PAM). Chromatin inaccessibility decreases dCas9 binding to other sites with matching seed sequences; thus 70% of off-target sites are associated with genes. The authors showed that targeted sequencing of 295 dCas9 binding sites in mESCs transfected with catalytically active Cas9 identified only one site mutated above background levels. The authors proposed a two-state model for Cas9 binding and cleavage, in which a seed match triggers binding but extensive pairing with target DNA is required for cleavage.
- Platt et al. established a Cre-dependent Cas9 knockin mouse. The authors demonstrated in vivo as well as ex vivo genome editing using adeno-associated virus (AAV)-, lentivirus-, or particle-mediated delivery of guide RNA in neurons, immune cells, and endothelial cells.
- Hsu et al. (2014) is a review article that discusses generally CRISPR-Cas9 history from yogurt to genome editing, including genetic screening of cells.

- Wang et al. (2014) relates to a pooled, loss-of-function genetic screening approach suitable for both positive and negative selection that uses a genome-scale lentiviral single guide RNA (sgRNA) library.
- Doench et al. created a pool of sgRNAs, tiling across all possible target sites of a panel of six endogenous mouse and three endogenous human genes and quantitatively assessed their ability to produce null alleles of their target gene by antibody staining and flow cytometry. The authors showed that optimization of the PAM improved activity and also provided an on-line tool for designing sgRNAs.
- Swiech et al. demonstrate that AAV-mediated SpCas9 genome editing can enable reverse genetic studies of gene function in the brain.
- Konermann et al. (2015) discusses the ability to attach multiple effector domains, e.g., transcriptional activator, functional and epigenomic regulators at appropriate positions on the guide such as stem or tetraloop with and without linkers.
- Zetsche et al. demonstrates that the Cas9 enzyme can be split into two and hence the assembly of Cas9 for activation can be controlled.
- Chen et al. relates to multiplex screening by demonstrating that a genome-wide in vivo CRISPR-Cas9 screen in mice reveals genes regulating lung metastasis.
- Ran et al. (2015) relates to SaCas9 and its ability to edit genomes and demonstrates that one cannot extrapolate from biochemical assays.
- Shalem et al. (2015) described ways in which catalytically inactive Cas9 (dCas9) fusions are used to synthetically repress (CRISPRi) or activate (CRISPRa) expression, showing advances using Cas9 for genome-scale screens, including arrayed and pooled screens, knockout approaches that inactivate genomic loci and strategies that modulate transcriptional activity.
- Xu et al. (2015) assessed the DNA sequence features that contribute to single guide RNA (sgRNA) efficiency in CRISPR-based screens. The authors explored efficiency of CRISPR/Cas9 knockout and nucleotide preference at the cleavage site. The authors also found that the sequence preference for CRISPRi/a is substantially different from that for CRISPR/Cas9 knockout.
- Parnas et al. (2015) introduced genome-wide pooled CRISPR-Cas9 libraries into dendritic cells (DCs) to identify genes that control the induction of tumor necrosis factor

(Tnf) by bacterial lipopolysaccharide (LPS). Known regulators of Tlr4 signaling and previously unknown candidates were identified and classified into three functional modules with distinct effects on the canonical responses to LPS.

- Ramanan et al (2015) demonstrated cleavage of viral episomal DNA (cccDNA) in infected cells. The HBV genome exists in the nuclei of infected hepatocytes as a 3.2kb double-stranded episomal DNA species called covalently closed circular DNA (cccDNA), which is a key component in the HBV life cycle whose replication is not inhibited by current therapies. The authors showed that sgRNAs specifically targeting highly conserved regions of HBV robustly suppresses viral replication and depleted cccDNA.
- Nishimasu et al. (2015) reported the crystal structures of SaCas9 in complex with a single guide RNA (sgRNA) and its double-stranded DNA targets, containing the 5'-TTGAAT-3' PAM and the 5'-TTGGGT-3' PAM. A structural comparison of SaCas9 with SpCas9 highlighted both structural conservation and divergence, explaining their distinct PAM specificities and orthologous sgRNA recognition.
- Zetsche et al. (2015) reported the characterization of Cpf1, a putative class 2 CRISPR effector. It was demonstrated that Cpf1 mediates robust DNA interference with features distinct from Cas9. Identifying this mechanism of interference broadens our understanding of CRISPR-Cas systems and advances their genome editing applications.
- Shmakov et al. (2015) reported the characterization of three distinct Class 2 CRISPR-Cas systems. The effectors of two of the identified systems, C2c1 and C2c3, contain RuvC like endonuclease domains distantly related to Cpf1. The third system, C2c2, contains an effector with two predicted HEPN RNase domains.

[00202] Also, “Dimeric CRISPR RNA-guided FokI nucleases for highly specific genome editing”, Shengdar Q. Tsai, Nicolas Wyvekens, Cyd Khayter, Jennifer A. Foden, Vishal Thapar, Deepak Reyon, Mathew J. Goodwin, Martin J. Aryee, J. Keith Joung *Nature Biotechnology* 32(6): 569-77 (2014), relates to dimeric RNA-guided FokI Nucleases that recognize extended sequences and can edit endogenous genes with high efficiencies in human cells.

[00203] Useful in the practice of the instant invention, reference is made to the article entitled BCL11A enhancer dissection by Cas9-mediated in situ saturating mutagenesis. Canver, M.C., Smith, E.C., Sher, F., Pinello, L., Sanjana, N.E., Shalem, O., Chen, D.D., Schupp, P.G.,

Vinjamur, D.S., Garcia, S.P., Luc, S., Kurita, R., Nakamura, Y., Fujiwara, Y., Maeda, T., Yuan, G., Zhang, F., Orkin, S.H., & Bauer, D.E. DOI:10.1038/nature15521, published online September 16, 2015, the article is herein incorporated by reference and discussed briefly below:

[00204] Canver *et al.* describes novel pooled CRISPR-Cas9 guide RNA libraries to perform *in situ* saturating mutagenesis of the human and mouse BCL11A erythroid enhancers previously identified as an enhancer associated with fetal hemoglobin (HbF) level and whose mouse ortholog is necessary for erythroid BCL11A expression. This approach revealed critical minimal features and discrete vulnerabilities of these enhancers. Through editing of primary human progenitors and mouse transgenesis, the authors validated the BCL11A erythroid enhancer as a target for HbF reinduction. The authors generated a detailed enhancer map that informs therapeutic genome editing.

[00205] Reference is made to Zetsche et al., "Cpf1 Is a Single RNA-Guided Endonuclease of a Class 2 CRISPR-Cas System," *Cell* 163, 1-13 (Oct. 22, 2015) and Shmakov et al., "Discovery and Functional Characterization of Diverse Class 2 CRISPR-Cas Systems," *Molecular Cell* 60, 1-13 (Available online Oct. 22, 2015). Zetsche et al. (2015) reported the characterization of Cpf1, a putative class 2 CRISPR effector. It was demonstrated that Cpf1 mediates robust DNA interference with features distinct from Cas9. Identifying this mechanism of interference broadens our understanding of CRISPR-Cas systems and advances their genome editing applications. In certain embodiments, perturb-seq uses Cpf1 as the CRISPR enzyme for introducing perturbations. Not being bound by a theory, Cpf1 does not require Tracr RNA and is a smaller enzyme, thus allowing higher combinatorial perturbations to be tested and for more efficient studies.

[00206] Shmakov et al. (2015) reported the characterization of three distinct Class 2 CRISPR-Cas systems. The effectors of two of the identified systems, C2c1 and C2c3, contain RuvC like endonuclease domains distantly related to Cpf1. The third system, C2c2, contains an effector with two predicted HEPN RNase domains. Mentino is also made of "Dimeric CRISPR RNA-guided FokI nucleases for highly specific genome editing", Shengdar Q. Tsai, Nicolas Wyvekens, Cyd Khayter, Jennifer A. Foden, Vishal Thapar, Deepak Reyon, Mathew J. Goodwin, Martin J. Aryee, J. Keith Joung *Nature Biotechnology* 32(6): 569-77 (2014), which relates to dimeric RNA-guided FokI Nucleases that recognize extended sequences and can edit endogenous genes with high efficiencies in human cells. In addition, mention is made of PCT

application PCT/US14/70057, Attorney Reference 47627.99.2060 and BI-2013/107 entitled "DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR TARGETING DISORDERS AND DISEASES USING PARTICLE DELIVERY COMPONENTS (claiming priority from one or more or all of US provisional patent applications: 62/054,490, filed September 24, 2014; 62/010,441, filed June 10, 2014; and 61/915,118, 61/915,215 and 61/915,148, each filed on December 12, 2013) ("the Particle Delivery PCT"), incorporated herein by reference, with respect to a method of preparing an sgRNA-and-Cas9 protein containing particle comprising admixing a mixture comprising an sgRNA and Cas9 protein (and optionally HDR template) with a mixture comprising or consisting essentially of or consisting of surfactant, phospholipid, biodegradable polymer, lipoprotein and alcohol; and particles from such a process. For example, wherein Cas9 protein and sgRNA were mixed together at a suitable, e.g., 3:1 to 1:3 or 2:1 to 1:2 or 1:1 molar ratio, at a suitable temperature, e.g., 15-30C, e.g., 20-25C, e.g., room temperature, for a suitable time, e.g., 15-45, such as 30 minutes, advantageously in sterile, nuclease free buffer, e.g., 1X PBS. Separately, particle components such as or comprising: a surfactant, e.g., cationic lipid, e.g., 1,2-dioleoyl-3-trimethylammonium-propane (DOTAP); phospholipid, e.g., dimyristoylphosphatidylcholine (DMPC); biodegradable polymer, such as an ethylene-glycol polymer or PEG, and a lipoprotein, such as a low-density lipoprotein, e.g., cholesterol were dissolved in an alcohol, advantageously a C1-6 alkyl alcohol, such as methanol, ethanol, isopropanol, e.g., 100% ethanol. The two solutions were mixed together to form particles containing the Cas9-sgRNA complexes. Accordingly, sgRNA may be pre-complexed with the Cas9 protein, before formulating the entire complex in a particle. Formulations may be made with a different molar ratio of different components known to promote delivery of nucleic acids into cells (e.g. 1,2-dioleoyl-3-trimethylammonium-propane (DOTAP), 1,2-ditetradecanoyl-sn-glycero-3-phosphocholine (DMPC), polyethylene glycol (PEG), and cholesterol) For example DOTAP : DMPC : PEG : Cholesterol Molar Ratios may be DOTAP 100, DMPC 0, PEG 0, Cholesterol 0; or DOTAP 90, DMPC 0, PEG 10, Cholesterol 0; or DOTAP 90, DMPC 0, PEG 5, Cholesterol 5. DOTAP 100, DMPC 0, PEG 0, Cholesterol 0. That application accordingly comprehends admixing sgRNA, Cas9 protein and components that form a particle; as well as particles from such admixing. Aspects of the instant invention can involve particles; for example, particles using a process analogous to that of the Particle Delivery PCT, e.g., by admixing a mixture comprising sgRNA

and/or Cas9 as in the instant invention and components that form a particle, e.g., as in the Particle Delivery PCT, to form a particle and particles from such admixing (or, of course, other particles involving sgRNA and/or Cas9 as in the instant invention). These and other CRISPR-Cas or CRISPR systems can be used in the practice of the invention.

Other Perturbations

[00207] The invention also involves perturbing by subjecting the cell to an increase or decrease in temperature. The temperature may range from about 0°C to about 100°C, advantageously about 10°C, 15°C, 20°C, 25°C, 30°C, 35°C, 40°C, 45°C, 50°C, 55°C, 60°C, 65°C, 70°C, 75°C, 80°C, 85°C, 90°C, 95°C or 100°C. In another embodiment, the temperature may be closer to a physiological temperature, e.g., about 30°C, 31°C, 32°C, 33°C, 34°C, 35°C, 36°C, 37°C, 38°C, 39°C or 40°C.

[00208] The invention also involves perturbing by subjecting the cell to a chemical agent. Samples of chemical agents include, but are not limited to, an antibiotic, a small molecule, a hormone, a hormone derivative, a steroid or a steroid derivative.

[00209] The invention also involves perturbing by subjecting the cell to a biological agent. Examples of biological agents may be, but are not limited to cytokines and Toll-like receptor agonists (Amit et al., *Science*. 2009 9;326(5950):257-63, Chevrier et al., *Cell*. 2011 Nov 11;147(4):853-67); Wang et al., *Cell*. 2015 Dec 3;163(6):1413-27; and Gaublomme et al., *Cell*. 2015 Dec 3;163(6):1400-12).

[00210] In one aspect of the invention the perturbing may be with an energy source such as electromagnetic energy or ultrasound. The electromagnetic energy may be a component of visible light having a wavelength in the range of 450nm-700nm. In a preferred embodiment the component of visible light may have a wavelength in the range of 450nm-500nm and may be blue light. The blue light may have an intensity of at least 0.2mW/cm², or more preferably at least 4mW/cm². In another embodiment, the component of visible light may have a wavelength in the range of 620-700nm and is red light.

[00211] The invention also involves perturbing by subjecting the cell to a chemical agent and/or temperature is across a gradient. A biomolecular gradient may be formed, for example, as reviewed in Keenan and Folch, *Lab Chip*. 2008 January; 8(1): doi:10.1039/b711887b. Biomolecule gradients have been shown to play roles in a wide range of biological processes including development, inflammation, wound healing, and cancer metastasis. Elucidation of

these phenomena requires the ability to expose cells to biomolecule gradients that are quantifiable, controllable, and mimic those that are present in vivo.

[00212] A chemical gradient may be formed without requiring fluid flow (see, e.g., Abhyankar et al., *Lab Chip*, 2006, 6, 389–393). This device consists of a membrane-covered source region and a large volume sink region connected by a microfluidic channel. The high fluidic resistance of the membrane limits fluid flow caused by pressure differences in the system, but allows diffusive transport of a chemical species through the membrane and into the channel. The large volume sink region at the end of the microfluidic channel helps to maintain spatial and temporal stability of the gradient. The chemical gradient in a 0.5 mm region near the sink region experiences a maximum of 10 percent change between the 6 and 24 h data points. Abhyankar et al., *Lab Chip*, 2006, 6, 389–393 present the theory, design, and characterization of this device and provide an example of neutrophil chemotaxis as proof of concept for future quantitative cell-signaling applications.

[00213] In another embodiment, a gradient may also be introduced with nanowires. In this embodiment, the nanowires do not necessarily introduce a gradient but may introduce other things into the system. A generalized platform for introducing a diverse range of biomolecules into living cells in high-throughput could transform how complex cellular processes are probed and analyzed. Shalek et al., *PNAS* | February 2, 2010 | vol. 107 | no. 5 demonstrate spatially localized, efficient, and universal delivery of biomolecules into immortalized and primary mammalian cells using surface-modified vertical silicon nanowires. The method relies on the ability of the silicon nanowires to penetrate a cell's membrane and subsequently release surface-bound molecules directly into the cell's cytosol, thus allowing highly efficient delivery of biomolecules without chemical modification or viral packaging. This modality enables one to assess the phenotypic consequences of introducing a broad range of biological effectors (DNAs, RNAs, peptides, proteins, and small molecules) into almost any cell type. Shalek et al., *PNAS* | February 2, 2010 | vol. 107 | no. 5 show that this platform can be used to guide neuronal progenitor growth with small molecules, knock down transcript levels by delivering siRNAs, inhibit apoptosis using peptides, and introduce targeted proteins to specific organelles. Shalek et al., *PNAS* | February 2, 2010 | vol. 107 | no. 5 further demonstrate codelivery of siRNAs and proteins on a single substrate in a microarray format, highlighting this technology's potential as a robust, monolithic platform for high-throughput, miniaturized bioassays.

[00214] A gradient may be established, for example, in a fluidic device, such as a microfluidic device (see, e.g., Tehranirokh et al., *BIOMICROFLUIDICS* 7, 051502 (2013)). Microfluidic technology allows dynamic cell culture in microperfusion systems to deliver continuous nutrient supplies for long term cell culture. It offers many opportunities to mimic the cell-cell and cell-extracellular matrix interactions of tissues by creating gradient concentrations of biochemical signals such as growth factors, chemokines, and hormones. Other applications of cell cultivation in microfluidic systems include high resolution cell patterning on a modified substrate with adhesive patterns and the reconstruction of complicated tissue architectures. In the review of Tehranirokh et al., *BIOMICROFLUIDICS* 7, 051502 (2013), recent advances in microfluidic platforms for cell culturing and proliferation, for both simple monolayer (2D) cell seeding processes and 3D configurations as accurate models of in vivo conditions, are examined. In another embodiment, the fluidic device may be a controlled fluidic device as described below.

Device for Establishing a Gradient

[00215] The invention provides a controlled fluidic device for establishing a gradient, particularly a concentration gradient, which may comprise a closed chamber comprising one or more inlet port(s) that deliver two or more different fluids via separate inlet channels and at least one outlet port wherein the location of the one or more inlet port(s) and the at least one outlet port are located so that the flow of fluids can be controlled within the closed chamber and a gradient of the mixture of the two or more fluids from the two or more inlet ports is established.

[00216] In a related aspect, the invention provides a controlled fluidic device wherein the controlled fluidic device is a polygonal plate having an upper surface and a lower surface and a peripheral plate edge having a pre-determined depth. In an embodiment of the invention, the controlled fluidic device may comprise a device wherein at least one of the two fluids includes at least one component for which a gradient could be established.

[00217] In a related aspect, the invention provides a controlled fluidic device which may comprise at least component which includes at least two subcomponents. In an embodiment of the invention, the controlled fluidic device may comprise a closed chamber comprising a chip.

[00218] The invention provides a method of identifying altered chemical resistance in a bacterial population in the controlled fluidic device as described above, the method which may comprise synthesizing a mutant bacterial strain to express fluorescent proteins; introducing a known concentration of the bacterial strain into the closed chamber; administering the two or more

different fluids into the closed chamber via the two more inlet ports; isolating DNA from a single cell; purifying DNA from bacteria; sequencing DNA from bacteria; preparing and sequencing a single composite sequence library; wherein wherein identification of alteration in level of expression compared to a baseline gene expression measurement of at least one biomarker is indicative of chemical resistance, and wherein the baseline gene expression measurement is the gene expression measured in the microfluidic well prior to administration of the two or more different fluids.

[00219]The present invention also provides a method of evaluating response in a cell population in the controlled fluidic device as described, the method which may comprise introducing a cell population into the closed chamber; administering the two or more different fluids into the closed chamber via the two more inlet ports such that a concentration gradient is established in the closed chamber; and, measuring the response of the cell population at various concentrations across the concentration gradient.

[00220]In another aspect, the invention provides a method of identifying altered bacterial populations according to the method of evaluating a response, the method comprising: a microfluidic device having a closed chamber having an upper surface and a lower surface and a peripheral plate edge having a predetermined depth; a plurality of microfluidic wells extending from the upper surface of the closed chamber, each well connected to adjacent ones of the plurality of wells by microchannels extending from the upper surface of the plate and extending from a first well to a second well such that the first well is in fluid communication with the second well; wherein the microfluidic device and plurality of wells connected by microchannels creates a chemical concentration gradient in adjacent microfluidic wells wherein one microfluidic well has a different chemical concentration than an adjacent microfluidic well; providing a chemical dye via an inlet port of the closed chamber; providing a chemical via an inlet port of the closed chamber; optionally providing a second chemical via an inlet port of the closed chamber; an outlet port of the closed chamber; and, a peripheral flow channel adjacent a portion of the peripheral plate edge and extending from the inlet port to the outlet port.

[00221]The present invention provides a method of identifying a compound associated with an altered bacterial population as described above, the method comprising: designing a combinatorial library wherein each member of the library comprises at least one pharmacophore associated with the altered gene expression; wherein alteration in level of expression compared

to a baseline gene expression measurement of at least one biomarker is indicative of an altered bacterial population; synthesizing a plurality of compounds from said combinatorial library; and, screening said compounds for candidates associated with the altered bacterial population.

[00222] The invention provides an array of controlled microfluidic devices, comprising a plurality of controlled microfluidic devices according to the controlled fluidic device for establishing a gradient, comprising a closed chamber comprising one or more inlet port(s) that deliver two or more different fluids via separate inlet channels and at least one outlet port wherein the location of the one or more inlet port(s) and the at least one outlet port are located so that the flow of fluids can be controlled within the closed chamber and a gradient of the mixture of the two or more fluids from the two or more inlet ports is established.

Labelling

[00223] Methods of international patent publication no. WO2014047561 and US patent publication no. 2015/0259674 are contemplated in the present invention.

[00224] The invention also contemplates a labeling ligand which may comprise a unique perturbation identifier (UPI) sequence attached to a perturbation-sequence-capture sequence, and sequencing includes isolating via microbeads comprising a perturbation-sequence-capture-binding-sequence having specific binding affinity for the perturbation-sequence-capture sequence attached to the UPI sequence.

[00225] The UPI sequence may be attached to a universal ligation handle sequence, whereby a unique source identifier USI may be generated by split-pool ligation. The labeling ligand may comprise an oligonucleotide label which may comprise a regulatory sequence configured for amplification by T7 polymerase. The labeling ligands may comprise oligonucleotide sequences configured to hybridize to a transcript specific region. The labeling ligand may comprise an oligonucleotide label, wherein the oligonucleotide label may further comprise a photocleavable linker.

[00226] The oligonucleotide label may further comprise a restriction enzyme site between the labeling ligand and unique constituent identifier (UCI).

[00227] The method may comprise forming discrete unique-identifier-transfer compositions, each of which may comprise the cell and a transfer particle, wherein: (a) an oligonucleotide label further may comprise a capture sequence, and unique constituent identifier (UCI) and capture sequence are together releasably attached to the labeling ligand; the labelling ligand is bound to

the target cellular constituent; and, (ca) transfer particle may comprise: (i) a capture-binding-sequence having specific binding affinity for the capture sequence attached to the UCI, and, (ii) a unique source identifier (USI) sequence that is unique to each transfer particle.

[00228] In one embodiment, the USI may comprise 4-15 nucleotides.

[00229] In another embodiment, the invention may further comprise releasing the UCI from the labeled ligand, under conditions within the unique-identifier-transfer composition so that the released capture sequence binds to the capture-binding-sequence on the transfer particle, thereby transferring the UCI to the transfer particle.

[00230] In another embodiment, the ligation handle may comprise a restriction site for producing an overhang complementary with a first index sequence overhang, and wherein the method further comprises digestion with a restriction enzyme. In another embodiment, the ligation handle may comprise a nucleotide sequence complementary with a ligation primer sequence and wherein the overhang complementary with a first index sequence overhang is produced by hybridization of the ligation primer to the ligation handle.

[00231] In another embodiment, the invention may further comprise quantitating relative amount of UCI sequence associated with a first cell to the amount of the same UCI sequence associated with a second cell, whereby the relative differences of a cellular constituent between cell(s) are determined.

[00232] In another embodiment, the labeling ligand may comprise an antibody or an antibody fragment, such as but not limited to, a nanobody, Fab, Fab', (Fab')₂, Fv, ScFv, diabody, triabody, tetrabody, Bis-scFv, minibody, Fab₂, or Fab₃ fragment.

[00233] In another embodiment, the labeling ligand may comprise an aptamer.

[00234] In another embodiment, the labeling ligand may comprise a nucleotide sequence complementary to a target sequence.

[00235] In another embodiment, the cell or a population includes wherein the cell(s) are a member of a cell population, and the method further comprises transforming or transducing the cell population with one or more genomic sequence-perturbation constructs that perturb a genomic sequence in the cells, wherein each distinct genomic sequence-perturbation construct comprises a unique-perturbation-identified (UPI) sequence unique to that construct. The genomic sequence-perturbation construct may comprises sequence encoding a guide RNA sequence of a CRISPR-Cas targeting system. The method may further comprise multiplex transformation of

the population of cells with a plurality of genomic sequence-perturbation constructs. The method may further comprise a UPI sequence attached to a perturbation-sequence-capture sequence, and the transfer particle may comprise a perturbation-sequence-capture-binding-sequence having specific binding affinity for the perturbation-sequence-capture sequence attached to the UPI sequence. The UPI sequence is attached to a universal ligation handle sequence, whereby a USI is generated by split-pool ligation.

[00236] In an advantageous embodiment, agents may be uniquely labeled in a dynamic manner (see, e.g., US provisional patent application serial no. 61/703,884 filed September 21, 2012). The unique labels are, at least in part, nucleic acid in nature, and may be generated by sequentially attaching two or more detectable oligonucleotide tags to each other and each unique label may be associated with a separate agent. A detectable oligonucleotide tag may be an oligonucleotide that may be detected by sequencing of its nucleotide sequence and/or by detecting non-nucleic acid detectable moieties to which it may be attached.

[00237] The oligonucleotide tags may be detectable by virtue of their nucleotide sequence, or by virtue of a non-nucleic acid detectable moiety that is attached to the oligonucleotide such as but not limited to a fluorophore, or by virtue of a combination of their nucleotide sequence and the nonnucleic acid detectable moiety.

[00238] In some embodiments, a detectable oligonucleotide tag may comprise one or more nonoligonucleotide detectable moieties. Examples of detectable moieties may include, but are not limited to, fluorophores, microparticles including quantum dots (Empodocles, et al., *Nature* 399:126-130, 1999), gold nanoparticles (Reichert et al., *Anal. Chem.* 72:6025-6029, 2000), microbeads (Lacoste et al., *Proc. Natl. Acad. Sci. USA* 97(17):9461-9466, 2000), biotin, DNP (dinitrophenyl), fucose, digoxigenin, haptens, and other detectable moieties known to those skilled in the art. In some embodiments, the detectable moieties may be quantum dots. Methods for detecting such moieties are described herein and/or are known in the art.

[00239] Thus, detectable oligonucleotide tags may be, but are not limited to, oligonucleotides which may comprise unique nucleotide sequences, oligonucleotides which may comprise detectable moieties, and oligonucleotides which may comprise both unique nucleotide sequences and detectable moieties.

[00240] A unique label may be produced by sequentially attaching two or more detectable oligonucleotide tags to each other. The detectable tags may be present or provided in a plurality

of detectable tags. The same or a different plurality of tags may be used as the source of each detectable tag may be part of a unique label. In other words, a plurality of tags may be subdivided into subsets and single subsets may be used as the source for each tag.

[00241] In some embodiments, one or more other species may be associated with the tags. In particular, nucleic acids released by a lysed cell may be ligated to one or more tags. These may include, for example, chromosomal DNA, RNA transcripts, tRNA, mRNA, mitochondrial DNA, or the like. Such nucleic acids may be sequenced, in addition to sequencing the tags themselves, which may yield information about the nucleic acid profile of the cells, which can be associated with the tags, or the conditions that the corresponding droplet or cell was exposed to.

Drop-Sequence Methods (“Drop-Seq”)

[00242] Cells come in different types, sub-types and activity states, which are classified based on their shape, location, function, or molecular profiles, such as the set of RNAs that they express. RNA profiling is in principle particularly informative, as cells express thousands of different RNAs. Approaches that measure for example the level of every type of RNA have until recently been applied to “homogenized” samples – in which the contents of all the cells are mixed together. Methods to profile the RNA content of tens and hundreds of thousands of individual human cells have been recently developed, including from brain tissues, quickly and inexpensively. To do so, special microfluidic devices have been developed to encapsulate each cell in an individual drop, associate the RNA of each cell with a ‘cell barcode’ unique to that cell/drop, measure the expression level of each RNA with sequencing, and then use the cell barcodes to determine which cell each RNA molecule came from. See, e.g., U.S. 62/048,227 filed September 9, 2014]

[00243] Methods of Macosko et al., 2015, “Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets” *Cell* 161, 1202–1214; International patent application number PCT/US2015/049178, published as WO2016/040476 on March 17, 2016; Klein et al., 2015, “Droplet Barcoding for Single-Cell Transcriptomics Applied to Embryonic Stem Cells” *Cell* 161, 1187–1201; Zheng, et al., 2016, “Haplotyping germline and cancer genomes with high-throughput linked-read sequencing” *Nature Biotechnology* 34, 303–311; and International patent publication number WO 2014210353 A2 are contemplated for the present invention, all the contents and disclosure of each of which are herein incorporated by reference in their entirety.

[00244] Microfluidics involves micro-scale devices that handle small volumes of fluids. Because microfluidics may accurately and reproducibly control and dispense small fluid volumes, in particular volumes less than 1 μ l, application of microfluidics provides significant cost-savings. The use of microfluidics technology reduces cycle times, shortens time-to-results, and increases throughput. Furthermore, incorporation of microfluidics technology enhances system integration and automation. Microfluidic reactions are generally conducted in microdroplets. The ability to conduct reactions in microdroplets depends on being able to merge different sample fluids and different microdroplets. See, e.g., US Patent Publication No. 20120219947. See also international patent application serial no. PCT/US2014/058637 for disclosure regarding a microfluidic laboratory on a chip.

[00245] Droplet microfluidics offers significant advantages for performing high-throughput screens and sensitive assays. Droplets allow sample volumes to be significantly reduced, leading to concomitant reductions in cost. Manipulation and measurement at kilohertz speeds enable up to 10^8 discrete biological entities (including, but not limited to, individual cells or organelles) to be screened in a single day. Compartmentalization in droplets increases assay sensitivity by increasing the effective concentration of rare species and decreasing the time required to reach detection thresholds. Droplet microfluidics combines these powerful features to enable currently inaccessible high-throughput screening applications, including single-cell and single-molecule assays. See, e.g., Guo et al., *Lab Chip*, 2012,12, 2146-2155.

[00246] Drop-Sequence methods and apparatus provides a high-throughput single-cell RNA-Seq and/or targeted nucleic acid profiling (for example, sequencing, quantitative reverse transcription polymerase chain reaction, and the like) where the RNAs from different cells are tagged individually, allowing a single library to be created while retaining the cell identity of each read. A combination of molecular barcoding and emulsion-based microfluidics to isolate, lyse, barcode, and prepare nucleic acids from individual cells in high-throughput is used. Microfluidic devices (for example, fabricated in polydimethylsiloxane), sub-nanoliter reverse emulsion droplets. These droplets are used to co-encapsulate nucleic acids with a barcoded capture bead. Each bead, for example, is uniquely barcoded so that each drop and its contents are distinguishable. The nucleic acids may come from any source known in the art, such as for example, those which come from a single cell, a pair of cells, a cellular lysate, or a solution. The cell is lysed as it is encapsulated in the droplet. To load single cells and barcoded beads into

these droplets with Poisson statistics, 100,000 to 10 million such beads are needed to barcode ~10,000-100,000 cells.

[00247] The invention provides a method for creating a single-cell sequencing library comprising: merging one uniquely barcoded mRNA capture microbead with a single-cell in an emulsion droplet having a diameter of 75-125 μm ; lysing the cell to make its RNA accessible for capturing by hybridization onto RNA capture microbead; performing a reverse transcription either inside or outside the emulsion droplet to convert the cell's mRNA to a first strand cDNA that is covalently linked to the mRNA capture microbead; pooling the cDNA-attached microbeads from all cells; and preparing and sequencing a single composite RNA-Seq library.

[00248] The invention provides a method for preparing uniquely barcoded mRNA capture microbeads, which has a unique barcode and diameter suitable for microfluidic devices comprising: 1) performing reverse phosphoramidite synthesis on the surface of the bead in a pool-and-split fashion, such that in each cycle of synthesis the beads are split into four reactions with one of the four canonical nucleotides (T, C, G, or A) or unique oligonucleotides of length two or more bases; 2) repeating this process a large number of times, at least two, and optimally more than twelve, such that, in the latter, there are more than 16 million unique barcodes on the surface of each bead in the pool. (See www.ncbi.nlm.nih.gov/pmc/articles/PMC206447)

[00249] Generally, the invention provides a method for preparing a large number of beads, particles, microbeads, nanoparticles, or the like with unique nucleic acid barcodes comprising performing polynucleotide synthesis on the surface of the beads in a pool-and-split fashion such that in each cycle of synthesis the beads are split into subsets that are subjected to different chemical reactions; and then repeating this split-pool process in two or more cycles, to produce a combinatorially large number of distinct nucleic acid barcodes. Invention further provides performing a polynucleotide synthesis wherein the synthesis may be any type of synthesis known to one of skill in the art for "building" polynucleotide sequences in a step-wise fashion. Examples include, but are not limited to, reverse direction synthesis with phosphoramidite chemistry or forward direction synthesis with phosphoramidite chemistry. Previous and well-known methods synthesize the oligonucleotides separately then "glue" the entire desired sequence onto the bead enzymatically. Applicants present a complexed bead and a novel process for producing these beads where nucleotides are chemically built onto the bead material in a high-throughput manner. Moreover, Applicants generally describe delivering a "packet" of

beads which allows one to deliver millions of sequences into separate compartments and then screen all at once.

[00250] The invention further provides an apparatus for creating a single-cell sequencing library via a microfluidic system, comprising: a oil-surfactant inlet comprising a filter and a carrier fluid channel, wherein said carrier fluid channel further comprises a resistor; an inlet for an analyte comprising a filter and a carrier fluid channel, wherein said carrier fluid channel further comprises a resistor; an inlet for mRNA capture microbeads and lysis reagent comprising a filter and a carrier fluid channel, wherein said carrier fluid channel further comprises a resistor; said carrier fluid channels have a carrier fluid flowing therein at an adjustable or predetermined flow rate; wherein each said carrier fluid channels merge at a junction; and said junction being connected to a mixer, which contains an outlet for drops.

[00251] A mixture comprising a plurality of microbeads adorned with combinations of the following elements: bead-specific oligonucleotide barcodes created by the described methods; additional oligonucleotide barcode sequences which vary among the oligonucleotides on an individual bead and can therefore be used to differentiate or help identify those individual oligonucleotide molecules; additional oligonucleotide sequences that create substrates for downstream molecular-biological reactions, such as oligo-dT (for reverse transcription of mature mRNAs), specific sequences (for capturing specific portions of the transcriptome, or priming for DNA polymerases and similar enzymes), or random sequences (for priming throughout the transcriptome or genome). In an embodiment, the individual oligonucleotide molecules on the surface of any individual microbead contain all three of these elements, and the third element includes both oligo-dT and a primer sequence.

[00252] Examples of the labeling substance which may be employed include labeling substances known to those skilled in the art, such as fluorescent dyes, enzymes, coenzymes, chemiluminescent substances, and radioactive substances. Specific examples include radioisotopes (e.g., ^{32}P , ^{14}C , ^{125}I , ^3H , and ^{131}I), fluorescein, rhodamine, dansyl chloride, umbelliferone, luciferase, peroxidase, alkaline phosphatase, β -galactosidase, β -glucosidase, horseradish peroxidase, glucoamylase, lysozyme, saccharide oxidase, microperoxidase, biotin, and ruthenium. In the case where biotin is employed as a labeling substance, preferably, after addition of a biotin-labeled antibody, streptavidin bound to an enzyme (e.g., peroxidase) is further added.

[00253] Advantageously, the label is a fluorescent label. Examples of fluorescent labels include, but are not limited to, Atto dyes, 4-acetamido-4'-isothiocyanatostilbene-2,2'-disulfonic acid; acridine and derivatives: acridine, acridine isothiocyanate; 5-(2'-aminoethyl)aminonaphthalene-1-sulfonic acid (EDANS); 4-amino-N-[3-vinylsulfonyl]phenyl]naphthalimide-3,5 disulfonate; N-(4-anilino-1-naphthyl)maleimide; anthranilamide; BODIPY; Brilliant Yellow; coumarin and derivatives; coumarin, 7-amino-4-methylcoumarin (AMC, Coumarin 120), 7-amino-4-trifluoromethylcoumarin (Coumarin 151); cyanine dyes; cyanosine; 4',6-diaminidino-2-phenylindole (DAPI); 5''-dibromopyrogallol-sulfonaphthalein (Bromopyrogallol Red); 7-diethylamino-3-(4'-isothiocyanatophenyl)-4-methylcoumarin; diethylenetriamine pentaacetate; 4,4'-diisothiocyanatodihydro-stilbene-2,2'-disulfonic acid; 4,4'-diisothiocyanatostilbene-2,2'-disulfonic acid; 5-[dimethylamino]naphthalene-1-sulfonyl chloride (DNS, dansylchloride); 4-dimethylaminophenylazophenyl-4'-isothiocyanate (DABITC); eosin and derivatives; eosin, eosin isothiocyanate, erythrosin and derivatives; erythrosin B, erythrosin, isothiocyanate; ethidium; fluorescein and derivatives; 5-carboxyfluorescein (FAM), 5-(4,6-dichlorotriazin-2-yl)aminofluorescein (DTAF), 2',7'-dimethoxy-4'5'-dichloro-6-carboxyfluorescein, fluorescein, fluorescein isothiocyanate, QFITC, (XRITC); fluorescamine; IR144; IR1446; Malachite Green isothiocyanate; 4-methylumbelliferoneortho cresolphthalein; nitrotyrosine; pararosaniline; Phenol Red; B-phycoerythrin; o-phthaldialdehyde; pyrene and derivatives: pyrene, pyrene butyrate, succinimidyl 1-pyrene; butyrate quantum dots; Reactive Red 4 (Cibacron.TM. Brilliant Red 3B-A) rhodamine and derivatives: 6-carboxy-X-rhodamine (ROX), 6-carboxyrhodamine (R6G), lissamine rhodamine B sulfonyl chloride rhodamine (Rhod), rhodamine B, rhodamine 123, rhodamine X isothiocyanate, sulforhodamine B, sulforhodamine 101, sulfonyl chloride derivative of sulforhodamine 101 (Texas Red); N,N,N',N' tetramethyl-6-carboxyrhodamine (TAMRA); tetramethyl rhodamine; tetramethyl rhodamine isothiocyanate (TRITC); riboflavin; rosolic acid; terbium chelate derivatives; Cy3; Cy5; Cy5.5; Cy7; IRD 700; IRD 800; La Jolla Blue; phthalocyanine; and naphthalocyanine.

[00254] The fluorescent label may be a fluorescent protein, such as blue fluorescent protein, cyan fluorescent protein, green fluorescent protein, red fluorescent protein, yellow fluorescent protein or any photoconvertible protein. Colormetric labeling, bioluminescent labeling and/or chemiluminescent labeling may further accomplish labeling. Labeling further may include

energy transfer between molecules in the hybridization complex by perturbation analysis, quenching, or electron transport between donor and acceptor molecules, the latter of which may be facilitated by double stranded match hybridization complexes. The fluorescent label may be a perylene or a terrylene. In the alternative, the fluorescent label may be a fluorescent bar code.

[00255] In an advantageous embodiment, the label may be light sensitive, wherein the label is light-activated and/or light cleaves the one or more linkers to release the molecular cargo. The light-activated molecular cargo may be a major light-harvesting complex (LHCII). In another embodiment, the fluorescent label may induce free radical formation.

[00256] The invention described herein enables high throughput and high resolution delivery of reagents to individual emulsion droplets that may contain cells, organelles, nucleic acids, proteins, etc. through the use of monodisperse aqueous droplets that are generated by a microfluidic device as a water-in-oil emulsion. The droplets are carried in a flowing oil phase and stabilized by a surfactant. In one aspect single cells or single organelles or single molecules (proteins, RNA, DNA) are encapsulated into uniform droplets from an aqueous solution/dispersion. In a related aspect, multiple cells or multiple molecules may take the place of single cells or single molecules. The aqueous droplets of volume ranging from 1 pL to 10 nL work as individual reactors. Disclosed embodiments provide 10^4 to 10^5 single cells in droplets which can be processed and analyzed in a single run.

[00257] To utilize microdroplets for rapid large-scale chemical screening or complex biological library identification, different species of microdroplets, each containing the specific chemical compounds or biological probes cells or molecular barcodes of interest, have to be generated and combined at the preferred conditions, e.g., mixing ratio, concentration, and order of combination.

[00258] Each species of droplet is introduced at a confluence point in a main microfluidic channel from separate inlet microfluidic channels. Preferably, droplet volumes are chosen by design such that one species is larger than others and moves at a different speed, usually slower than the other species, in the carrier fluid, as disclosed in U.S. Publication No. US 2007/0195127 and International Publication No. WO 2007/089541, each of which are incorporated herein by reference in their entirety. The channel width and length is selected such that faster species of droplets catch up to the slowest species. Size constraints of the channel prevent the faster moving droplets from passing the slower moving droplets resulting in a train of droplets entering a merge

zone. Multi-step chemical reactions, biochemical reactions, or assay detection chemistries often require a fixed reaction time before species of different type are added to a reaction. Multi-step reactions are achieved by repeating the process multiple times with a second, third or more confluence points each with a separate merge point. Highly efficient and precise reactions and analysis of reactions are achieved when the frequencies of droplets from the inlet channels are matched to an optimized ratio and the volumes of the species are matched to provide optimized reaction conditions in the combined droplets.

[00259] Fluidic droplets may be screened or sorted within a fluidic system of the invention by altering the flow of the liquid containing the droplets. For instance, in one set of embodiments, a fluidic droplet may be steered or sorted by directing the liquid surrounding the fluidic droplet into a first channel, a second channel, etc. In another set of embodiments, pressure within a fluidic system, for example, within different channels or within different portions of a channel, can be controlled to direct the flow of fluidic droplets. For example, a droplet can be directed toward a channel junction including multiple options for further direction of flow (e.g., directed toward a branch, or fork, in a channel defining optional downstream flow channels). Pressure within one or more of the optional downstream flow channels can be controlled to direct the droplet selectively into one of the channels, and changes in pressure can be effected on the order of the time required for successive droplets to reach the junction, such that the downstream flow path of each successive droplet can be independently controlled. In one arrangement, the expansion and/or contraction of liquid reservoirs may be used to steer or sort a fluidic droplet into a channel, e.g., by causing directed movement of the liquid containing the fluidic droplet. In another embodiment, the expansion and/or contraction of the liquid reservoir may be combined with other flow-controlling devices and methods, e.g., as described herein. Non-limiting examples of devices able to cause the expansion and/or contraction of a liquid reservoir include pistons.

[00260] Key elements for using microfluidic channels to process droplets include: (1) producing droplet of the correct volume, (2) producing droplets at the correct frequency and (3) bringing together a first stream of sample droplets with a second stream of sample droplets in such a way that the frequency of the first stream of sample droplets matches the frequency of the second stream of sample droplets. Preferably, bringing together a stream of sample droplets with

a stream of premade library droplets in such a way that the frequency of the library droplets matches the frequency of the sample droplets.

[00261] Methods for producing droplets of a uniform volume at a regular frequency are well known in the art. One method is to generate droplets using hydrodynamic focusing of a dispersed phase fluid and immiscible carrier fluid, such as disclosed in U.S. Publication No. US 2005/0172476 and International Publication No. WO 2004/002627. It is desirable for one of the species introduced at the confluence to be a pre-made library of droplets where the library contains a plurality of reaction conditions, e.g., a library may contain plurality of different compounds at a range of concentrations encapsulated as separate library elements for screening their effect on cells or enzymes, alternatively a library could be composed of a plurality of different primer pairs encapsulated as different library elements for targeted amplification of a collection of loci, alternatively a library could contain a plurality of different antibody species encapsulated as different library elements to perform a plurality of binding assays. The introduction of a library of reaction conditions onto a substrate is achieved by pushing a premade collection of library droplets out of a vial with a drive fluid. The drive fluid is a continuous fluid. The drive fluid may comprise the same substance as the carrier fluid (e.g., a fluorocarbon oil). For example, if a library consists of ten pico-liter droplets is driven into an inlet channel on a microfluidic substrate with a drive fluid at a rate of 10,000 pico-liters per second, then nominally the frequency at which the droplets are expected to enter the confluence point is 1000 per second. However, in practice droplets pack with oil between them that slowly drains. Over time the carrier fluid drains from the library droplets and the number density of the droplets (number/mL) increases. Hence, a simple fixed rate of infusion for the drive fluid does not provide a uniform rate of introduction of the droplets into the microfluidic channel in the substrate. Moreover, library-to-library variations in the mean library droplet volume result in a shift in the frequency of droplet introduction at the confluence point. Thus, the lack of uniformity of droplets that results from sample variation and oil drainage provides another problem to be solved. For example if the nominal droplet volume is expected to be 10 pico-liters in the library, but varies from 9 to 11 pico-liters from library-to-library then a 10,000 pico-liter/second infusion rate will nominally produce a range in frequencies from 900 to 1,100 droplet per second. In short, sample to sample variation in the composition of dispersed phase for droplets made on chip, a tendency for the number density of library droplets to increase over time and library-to-

library variations in mean droplet volume severely limit the extent to which frequencies of droplets may be reliably matched at a confluence by simply using fixed infusion rates. In addition, these limitations also have an impact on the extent to which volumes may be reproducibly combined. Combined with typical variations in pump flow rate precision and variations in channel dimensions, systems are severely limited without a means to compensate on a run-to-run basis. The foregoing facts not only illustrate a problem to be solved, but also demonstrate a need for a method of instantaneous regulation of microfluidic control over microdroplets within a microfluidic channel.

[00262] Combinations of surfactant(s) and oils must be developed to facilitate generation, storage, and manipulation of droplets to maintain the unique chemical/biochemical/biological environment within each droplet of a diverse library. Therefore, the surfactant and oil combination must (1) stabilize droplets against uncontrolled coalescence during the drop forming process and subsequent collection and storage, (2) minimize transport of any droplet contents to the oil phase and/or between droplets, and (3) maintain chemical and biological inertness with contents of each droplet (e.g., no adsorption or reaction of encapsulated contents at the oil-water interface, and no adverse effects on biological or chemical constituents in the droplets). In addition to the requirements on the droplet library function and stability, the surfactant-in-oil solution must be coupled with the fluid physics and materials associated with the platform. Specifically, the oil solution must not swell, dissolve, or degrade the materials used to construct the microfluidic chip, and the physical properties of the oil (e.g., viscosity, boiling point, etc.) must be suited for the flow and operating conditions of the platform.

[00263] Droplets formed in oil without surfactant are not stable to permit coalescence, so surfactants must be dissolved in the oil that is used as the continuous phase for the emulsion library. Surfactant molecules are amphiphilic--part of the molecule is oil soluble, and part of the molecule is water soluble. When a water-oil interface is formed at the nozzle of a microfluidic chip for example in the inlet module described herein, surfactant molecules that are dissolved in the oil phase adsorb to the interface. The hydrophilic portion of the molecule resides inside the droplet and the fluorophilic portion of the molecule decorates the exterior of the droplet. The surface tension of a droplet is reduced when the interface is populated with surfactant, so the stability of an emulsion is improved. In addition to stabilizing the droplets against coalescence,

the surfactant should be inert to the contents of each droplet and the surfactant should not promote transport of encapsulated components to the oil or other droplets.

[00264] A droplet library may be made up of a number of library elements that are pooled together in a single collection (see, e.g., US Patent Publication No. 2010002241). Libraries may vary in complexity from a single library element to 10¹⁵ library elements or more. Each library element may be one or more given components at a fixed concentration. The element may be, but is not limited to, cells, organelles, virus, bacteria, yeast, beads, amino acids, proteins, polypeptides, nucleic acids, polynucleotides or small molecule chemical compounds. The element may contain an identifier such as a label. The terms "droplet library" or "droplet libraries" are also referred to herein as an "emulsion library" or "emulsion libraries." These terms are used interchangeably throughout the specification.

[00265] A cell library element may include, but is not limited to, hybridomas, B-cells, primary cells, cultured cell lines, cancer cells, stem cells, cells obtained from tissue, or any other cell type. Cellular library elements are prepared by encapsulating a number of cells from one to hundreds of thousands in individual droplets. The number of cells encapsulated is usually given by Poisson statistics from the number density of cells and volume of the droplet. However, in some cases the number deviates from Poisson statistics as described in Edd et al., "Controlled encapsulation of single-cells into monodisperse picolitre drops." *Lab Chip*, 8(8): 1262-1264, 2008. The discrete nature of cells allows for libraries to be prepared in mass with a plurality of cellular variants all present in a single starting media and then that media is broken up into individual droplet capsules that contain at most one cell. These individual droplets capsules are then combined or pooled to form a library consisting of unique library elements. Cell division subsequent to, or in some embodiments following, encapsulation produces a clonal library element.

[00266] A bead based library element may contain one or more beads, of a given type and may also contain other reagents, such as antibodies, enzymes or other proteins. In the case where all library elements contain different types of beads, but the same surrounding media, the library elements may all be prepared from a single starting fluid or have a variety of starting fluids. In the case of cellular libraries prepared in mass from a collection of variants, such as genomically modified, yeast or bacteria cells, the library elements will be prepared from a variety of starting fluids.

[00267] Often it is desirable to have exactly one cell per droplet with only a few droplets containing more than one cell when starting with a plurality of cells or yeast or bacteria, engineered to produce variants on a protein. In some cases, variations from Poisson statistics may be achieved to provide an enhanced loading of droplets such that there are more droplets with exactly one cell per droplet and few exceptions of empty droplets or droplets containing more than one cell.

[00268] Examples of droplet libraries are collections of droplets that have different contents, ranging from beads, cells, small molecules, DNA, primers, antibodies. Smaller droplets may be in the order of femtoliter (fL) volume drops, which are especially contemplated with the droplet dispensors. The volume may range from about 5 to about 600 fL. The larger droplets range in size from roughly 0.5 micron to 500 micron in diameter, which corresponds to about 1 pico liter to 1 nano liter. However, droplets may be as small as 5 microns and as large as 500 microns. Preferably, the droplets are at less than 100 microns, about 1 micron to about 100 microns in diameter. The most preferred size is about 20 to 40 microns in diameter (10 to 100 picoliters). The preferred properties examined of droplet libraries include osmotic pressure balance, uniform size, and size ranges.

[00269] The droplets comprised within the emulsion libraries of the present invention may be contained within an immiscible oil which may comprise at least one fluorosurfactant. In some embodiments, the fluorosurfactant comprised within immiscible fluorocarbon oil is a block copolymer consisting of one or more perfluorinated polyether (PFPE) blocks and one or more polyethylene glycol (PEG) blocks. In other embodiments, the fluorosurfactant is a triblock copolymer consisting of a PEG center block covalently bound to two PFPE blocks by amide linking groups. The presence of the fluorosurfactant (similar to uniform size of the droplets in the library) is critical to maintain the stability and integrity of the droplets and is also essential for the subsequent use of the droplets within the library for the various biological and chemical assays described herein. Fluids (e.g., aqueous fluids, immiscible oils, etc.) and other surfactants that may be utilized in the droplet libraries of the present invention are described in greater detail herein.

[00270] The present invention provides an emulsion library which may comprise a plurality of aqueous droplets within an immiscible oil (e.g., fluorocarbon oil) which may comprise at least one fluorosurfactant, wherein each droplet is uniform in size and may comprise the same

aqueous fluid and may comprise a different library element. The present invention also provides a method for forming the emulsion library which may comprise providing a single aqueous fluid which may comprise different library elements, encapsulating each library element into an aqueous droplet within an immiscible fluorocarbon oil which may comprise at least one fluorosurfactant, wherein each droplet is uniform in size and may comprise the same aqueous fluid and may comprise a different library element, and pooling the aqueous droplets within an immiscible fluorocarbon oil which may comprise at least one fluorosurfactant, thereby forming an emulsion library.

[00271] For example, in one type of emulsion library, all different types of elements (e.g., cells or beads), may be pooled in a single source contained in the same medium. After the initial pooling, the cells or beads are then encapsulated in droplets to generate a library of droplets wherein each droplet with a different type of bead or cell is a different library element. The dilution of the initial solution enables the encapsulation process. In some embodiments, the droplets formed will either contain a single cell or bead or will not contain anything, i.e., be empty. In other embodiments, the droplets formed will contain multiple copies of a library element. The cells or beads being encapsulated are generally variants on the same type of cell or bead. In one example, the cells may comprise cancer cells of a tissue biopsy, and each cell type is encapsulated to be screened for genomic data or against different drug therapies. Another example is that 10^{11} or 10^{15} different type of bacteria; each having a different plasmid spliced therein, are encapsulated. One example is a bacterial library where each library element grows into a clonal population that secretes a variant on an enzyme.

[00272] In another example, the emulsion library may comprise a plurality of aqueous droplets within an immiscible fluorocarbon oil, wherein a single molecule may be encapsulated, such that there is a single molecule contained within a droplet for every 20-60 droplets produced (e.g., 20, 25, 30, 35, 40, 45, 50, 55, 60 droplets, or any integer in between). Single molecules may be encapsulated by diluting the solution containing the molecules to such a low concentration that the encapsulation of single molecules is enabled. In one specific example, a LacZ plasmid DNA was encapsulated at a concentration of 20 fM after two hours of incubation such that there was about one gene in 40 droplets, where 10 μ m droplets were made at 10 kHz per second. Formation of these libraries rely on limiting dilutions.

[00273] Methods of the invention involve forming sample droplets. The droplets are aqueous droplets that are surrounded by an immiscible carrier fluid. Methods of forming such droplets are shown for example in Link et al. (U.S. patent application numbers 2008/0014589, 2008/0003142, and 2010/0137163), Stone et al. (U.S. Pat. No. 7,708,949 and U.S. patent application number 2010/0172803), Anderson et al. (U.S. Pat. No. 7,041,481 and which reissued as RE41,780) and European publication number EP2047910 to Raindance Technologies Inc. The content of each of which is incorporated by reference herein in its entirety.

[00274] In certain embodiments, the carrier fluid may contain one or more additives, such as agents which reduce surface tensions (surfactants). Surfactants can include Tween, Span, fluorosurfactants, and other agents that are soluble in oil relative to water. In some applications, performance is improved by adding a second surfactant to the sample fluid. Surfactants can aid in controlling or optimizing droplet size, flow and uniformity, for example by reducing the shear force needed to extrude or inject droplets into an intersecting channel. This can affect droplet volume and periodicity, or the rate or frequency at which droplets break off into an intersecting channel. Furthermore, the surfactant can serve to stabilize aqueous emulsions in fluorinated oils from coalescing.

[00275] In certain embodiments, the droplets may be surrounded by a surfactant which stabilizes the droplets by reducing the surface tension at the aqueous oil interface. Preferred surfactants that may be added to the carrier fluid include, but are not limited to, surfactants such as sorbitan-based carboxylic acid esters (e.g., the "Span" surfactants, Fluka Chemika), including sorbitan monolaurate (Span 20), sorbitan monopalmitate (Span 40), sorbitan monostearate (Span 60) and sorbitan monooleate (Span 80), and perfluorinated polyethers (e.g., DuPont Krytox 157 FSL, FSM, and/or FSH). Other non-limiting examples of non-ionic surfactants which may be used include polyoxyethylenated alkylphenols (for example, nonyl-, p-dodecyl-, and dinonylphenols), polyoxyethylenated straight chain alcohols, polyoxyethylenated polyoxypropylene glycols, polyoxyethylenated mercaptans, long chain carboxylic acid esters (for example, glyceryl and polyglyceryl esters of natural fatty acids, propylene glycol, sorbitol, polyoxyethylenated sorbitol esters, polyoxyethylene glycol esters, etc.) and alkanolamines (e.g., diethanolamine-fatty acid condensates and isopropanolamine-fatty acid condensates).

[00276] By incorporating a plurality of unique tags into the additional droplets and joining the tags to a solid support designed to be specific to the primary droplet, the conditions that the

primary droplet is exposed to may be encoded and recorded. For example, nucleic acid tags can be sequentially ligated to create a sequence reflecting conditions and order of same. Alternatively, the tags can be added independently appended to solid support. Non-limiting examples of a dynamic labeling system that may be used to bioninformatically record information can be found at US Provisional Patent Application entitled "Compositions and Methods for Unique Labeling of Agents" filed September 21, 2012 and November 29, 2012. In this way, two or more droplets may be exposed to a variety of different conditions, where each time a droplet is exposed to a condition, a nucleic acid encoding the condition is added to the droplet each ligated together or to a unique solid support associated with the droplet such that, even if the droplets with different histories are later combined, the conditions of each of the droplets are remain available through the different nucleic acids. Non-limiting examples of methods to evaluate response to exposure to a plurality of conditions can be found at US Provisional Patent Application entitled "Systems and Methods for Droplet Tagging" filed September 21, 2012.

[00277] Applications of the disclosed device may include use for the dynamic generation of molecular barcodes (e.g., DNA oligonucleotides, flurophores, etc.) either independent from or in concert with the controlled delivery of various compounds of interest (drugs, small molecules, siRNA, CRISPR guide RNAs, reagents, etc.). For example, unique molecular barcodes can be created in one array of nozzles while individual compounds or combinations of compounds can be generated by another nozzle array. Barcodes/compounds of interest can then be merged with cell-containing droplets. An electronic record in the form of a computer log file is kept to associate the barcode delivered with the downstream reagent(s) delivered. This methodology makes it possible to efficiently screen a large population of cells for applications such as single-cell drug screening, controlled perturbation of regulatory pathways, etc. The device and techniques of the disclosed invention facilitate efforts to perform studies that require data resolution at the single cell (or single molecule) level and in a cost effective manner. Disclosed embodiments provide a high throughput and high resolution delivery of reagents to individual emulsion droplets that may contain cells, nucleic acids, proteins, etc. through the use of monodisperse aqueous droplets that are generated one by one in a microfluidic chip as a water-in-oil emulsion. Hence, the invention proves advantageous over prior art systems by being able to dynamically track individual cells and droplet treatments/combinations during life cycle

experiments. Additional advantages of the disclosed invention provides an ability to create a library of emulsion droplets on demand with the further capability of manipulating the droplets through the disclosed process(es). Disclosed embodiments may, thereby, provide dynamic tracking of the droplets and create a history of droplet deployment and application in a single cell based environment.

[00278] Droplet generation and deployment is produced via a dynamic indexing strategy and in a controlled fashion in accordance with disclosed embodiments of the present invention. Disclosed embodiments of the microfluidic device described herein provides the capability of microdroplets that be processed, analyzed and sorted at a highly efficient rate of several thousand droplets per second, providing a powerful platform which allows rapid screening of millions of distinct compounds, biological probes, proteins or cells either in cellular models of biological mechanisms of disease, or in biochemical, or pharmacological assays.

Undersampling - A sampling based framework for genetic interactions

[00279] According to the invention, random sampling may comprise matrix completion, tensor completion, compressed sensing, or kernel learning.

[00280] In some aspects, where random sampling comprises matrix completion, tensor completion, or compressed sensing, π may be of the order of $\log P$.

[00281] The invention relies on a random sampling assumption, e.g. that the combinatorial space is sparse and/or of low rank. This assumption is generic and advantageously does not rely on the pre-determination of a (known) set of genetic interactions. This assumption constrains the range or complexity of models, and thus can be used to restrict sampling size (undersampling). Further, as detailed below, the invention relies on the following: (1) Given a limited number of assays, if one wishes to infer interactions up to an order j , it is advantageous to randomly sample interactions at a higher order $k > j$, because higher order perturbations maximize the information that can be recovered; and (2) in such a method, one may use a model that accounts for higher order interactions when analyzing lower order ones. For example, it is possible to aim for each perturbation to target $k \sim 5-7$ genes at once to estimate/model interactions at lower order $j \sim 3-5$.

[00282] Although some experimental methods open the way to test non-linear interactions by high order combinatorial genetic perturbations, exhaustive combinatorial exploration is intractable for anything but 2- or 3-way interactions for a few genes.

[00283] According to the invention, random matrix theory and compressive sensing may be used to re-formulate this as a random sampling problem, developing a new framework from experimental design to model inference, testing and refinement.

[00284] To infer combinatorial models from a dramatic under-sampling of the full high-order combinatorial space with massively combinatorial molecular perturbations (MCP), one may rely on random matrix theory, compressive sensing and kernel learning.

[00285] According to the invention, it is made possible to model non-linear regulatory functions from genetic manipulations (perturbations).

[00286] One may learn models of higher-order genetic interactions from combinatorial perturbations with single cell profiling. Although the learning problem is underdetermined due to combinatorial explosion (2^m possible interaction terms among m genes), it can become tractable in the presence of additional structure, including sparsity and smoothness, that constrains the range or complexity of models. One may thus rely on the following: (1) Given a limited number of assays, if one wishes to infer interactions up to an order j , it is advantageous to randomly sample interactions at a higher order $k > j$, because higher order perturbations maximize the information that can be recovered; and (2) in such a design, one can use a model that accounts for higher order interactions when analyzing lower order ones. One may for example aim for each perturbation to target $k \sim 5-7$ genes at once to estimate interactions at $j \sim 3-5$.

[00287] Thus the present invention relies on a learning approach that takes multiplex perturbations at a high order n and a complex readout data (e.g., RNA profile) and infers a model of genetic interactions at a lower order ($m < n$), as well as strategies for experimental design, model testing and refinement.

[00288] If one assumes that genetic interactions are low rank, sparse, or both, then the true number of degrees of freedom is small relative to the complete combinatorial expansion, so that one can infer the full nonlinear landscape with a relatively small random sampling of high-order perturbations, without specific knowledge of which genes are likely to interact. Analysis of prior studies supports the sparsity assumption in yeast (for fitness: Costanzo, M., Baryshnikova, A., Bellay, J., Kim, Y., Spear, E. D., Sevier, C. S., Ding, H., Koh, J. L., Toufighi, K., Mostafavi, S., Prinz, J., St Onge, R. P., VanderSluis, B., Makhnevych, T., Vizeacoumar, F. J., Alizadeh, S., Bahr, S., Brost, R. L., Chen, Y., Cokol, M., Deshpande, R., Li, Z., Lin, Z. Y., Liang, W., Marback, M., Paw, J., San Luis, B. J., Shuteriqi, E., Tong, A. H., van Dyk, N., Wallace, I. M.,

Whitney, J. A., Weirauch, M. T., Zhong, G., Zhu, H., Houry, W. A., Brudno, M., Ragibizadeh, S., Papp, B., Pal, C., Roth, F. P., Giaever, G., Nislow, C., Troyanskaya, O. G., Bussey, H., Bader, G. D., Gingras, A. C., Morris, Q. D., Kim, P. M., Kaiser, C. A., Myers, C. L., Andrews, B. J. & Boone, C. The genetic landscape of a cell. *Science*. 327, 425-431, doi:10.1126/science.1180823 (2010)), and fly (for 11 imaging phenotypes: Laufer, C., Fischer, B., Billmann, M., Huber, W. & Boutros, M. Mapping genetic interactions in human cancer cells with RNAi and multiparametric phenotyping. *Nat Methods*. 10, 427-431, doi:10.1038/nmeth.2436 (2013)), and to the limited tested extent, mammals (for 60 genes: Bassik, M. C., Kampmann, M., Lebbink, R. J., Wang, S., Hein, M. Y., Poser, I., Weibezahn, J., Horlbeck, M. A., Chen, S., Mann, M., Hyman, A. A., Leproust, E. M., McManus, M. T. & Weissman, J. S. A systematic mammalian genetic interaction map reveals pathways underlying ricin susceptibility. *Cell*. 152, 909-922, doi:10.1016/j.cell.2013.01.030 (2013). PMID:3652613).

Matrix (tensor) completion.

[00289] All the values of a matrix (tensor) are filled in using a small collection of sampled entries. Applicants hypothesize that the rank of a tensor of higher-order interactions is a fraction of the number of tested genes which is tested by by calculating the rank from a dense sampling of second or third order knockouts from a small collection of genes. If the rank of interactions is limited, then Applicants randomly sample sets of genes to knockout from a larger collection, and fill in the remaining values via nuclear norm regularized least-squares optimization (Candes, E. J. & Plan, Y. Matrix Completion With Noise. *Proceedings of the IEEE*. 98, 925-936, doi:Doi 10.1109/Jproc.2009.2035722 (2010)). Provable guarantees suggest that if the rank, r , is small relative to the number of genes, n , then $m \geq O(n^{6/5} r \log n)$ sampled entries are sufficient. However, since these guarantees assume rough uniformity in the loadings of interaction singular vectors, this assumption is unlikely to hold if the interaction matrix is very sparse. In this case, Applicants perform the same random sampling, and simultaneously regularize over both the nuclear norm and the L1 norm of the matrix (Richard, E., Savalle, P. & Vayatis, N. Estimation of Simultaneously Sparse and Low Rank Matrices. *arXiv*. doi:arXiv:1206.6474).

Compressed sensing

[00290] Here, instead of working with a tensor of interaction terms, Applicants work with a basis that spans all higher order interactions. Each single quantitative phenotype is a real-valued

function $f(g)$ on possible genotypes g (the 2^m possible allelic or knockout states), represented as binary strings of length m . Applicants analyze such Boolean functions using *Fourier decomposition* (O'Donnell, R. *Analysis of boolean functions*. (Cambridge University Press, 2014)) $f(g) = \sum_{b \in \{0,1\}^m} \hat{f}_b (-1)^{b \cdot g}$, $\hat{f}_b = \frac{1}{2^m} \sum_{g \in \{0,1\}^m} f(g) (-1)^{g \cdot b}$, where f is an orthogonal basis indexed by binary strings b , and each Fourier coefficient \hat{f}_b precisely quantifies the effect of one possible multi-gene interaction. For example with $m = 2$, \hat{f}_{00} is the average phenotype; \hat{f}_{10} is the effect of the first gene KO, marginalized over the genetic background of the second; similarly for \hat{f}_{01} ; and \hat{f}_{11} quantifies the two-way interaction (the extent to which the double KO phenotype differs from that predicted by the sum of the effects of the single KOs). Applicants hypothesize that such genotype-phenotype maps are approximately sparse in the Fourier basis, such that there is a small number, s , of nonzero Fourier coefficients (not known *a priori*). With perturbations generated only up to a limited order, Applicants obtain a truncated Fourier model, which is a general linear model: the genetic interactions are in the basis functions (encoded into a design matrix), and the response is linear in the unknown Fourier coefficients. Applicants assume most truncated coefficients are negligible. Assuming that the genotype-phenotype maps are approximately sparse in the Fourier basis, Applicants use L1-penalized regression to learn the coefficients of the map from paired genotype-phenotype observations $g_i, f(g_i)$ (with uncertainty or noise in both).

[00291] Compressed sensing posits that if Applicants' perturbations are de-coherent under the given basis, then exact recovery is possible with dramatic under-sampling (in the noiseless case) (Candes, E. Compressive sampling. *Proceedings of the International Congress of Mathematicians Madrid, August 22--30, 2006*. 3, 19, doi:10.4171/022 (2006)), such that a sample size $n = C s \log p$ will suffice, where s is the number of effectively nonzero coefficients, p is the magnitude of combinatorial expansion and C depends on noise and experimental design (how the g_i are sampled) (Candes, E. Mathematics of sparsity (and few other things). *ICM 2014 Proceedings, to appear*. (2014)). By varying the penalization parameter, Applicants learn sparse structures at different levels of thresholding, and find the level below which the data become insufficient to capture the signal (Hastie, T., Friedman, J. & Tibshirani, R. *The elements of statistical learning*. Vol. 2 (Springer, 2009)). Applicants explore using a larger penalization parameter on the higher order interaction coefficients, and, with good estimates of single perturbations, even no penalty on the linear terms, or regressing those out first. If each

experiment is a Poisson random sampling of KOs, Applicants expect the measurements to have good de-coherence under the Fourier basis, provided the mean number of KO experiments per gene is not too low. If Applicants' assumptions are correct, a soft phase transition in performance as the number of observations crosses a threshold should be observed. Applicants use a small complete dataset or downsampling of a larger more random dataset, to assess if the appropriate transition is observed.

Kernel learning.

[00292] If there is no strict sparsity in the rank or in the coefficients, Applicants build predictive functions of the effects of combinatorial perturbations, using a kernel of experimental similarity. Given m experiments, Applicants define an $m \times m$ polynomial kernel, for example, based on the overlap in knockouts between any pair of experiments. Applicants learn a weighted combination of kernel vectors that fits a collection of training data, and use the coefficients to predict the outcome of new experiments. Here, the density of nonlinear interaction terms can be much greater, since Applicants do not directly learn any particular interaction coefficient, but rather a kernelized version of the entire polynomial. Indeed, if the interaction terms are too sparse, kernel learning is unlikely to be successful with under-sampling.

[00293] Applicants analyzed 3-way interaction data measured by overexpression of every 3-way combination of 39 miRNAs and a phenotype of drug resistance, and confirmed substantial sparsity in the data. Applicants analyzed the 5-way interactions affecting expression profiles in response to salt in yeast between the MAPK Hog1 (p38 ortholog) and 4 TFs (1, 2, 3, 4, and 5 KO: 32 perturbations). Using a (non-regularized) linear model, Applicants quantified 1- and 2-way interactions, finding diverse non-linearities.

Analyzing a cell population at the single cell level

[00294] The method according to the invention may comprise a step for single-cell molecular profiling. In some embodiments the step may comprise processing said cell population in order to physically separate cells. In some embodiments the step may comprise single-cell manipulation, e.g. using microfluidics based techniques. In some embodiments the step may comprise reverse emulsion droplet-based single-cell analysis or hydrogel droplet-based single-cell analysis.

[00295] The method of the invention may use microfluidics, e.g. to culture cells in specific combinations, control the spatiotemporal signals they receive, and/or trace and sample them as desired.

Molecular profiling at the single cell level

[00296] The method according to the invention may comprise a step for single-cell molecular profiling. This step may involve analyzing biomolecules quantitatively or semi-quantitatively. The biomolecules may include RNA, mRNA, pre-mRNA, proteins, chromatin or DNA. Said analysis may be performed genome-wide. Said analysis may be coupled (dual or sequential analysis of two or more types of biomolecules).

[00297] In some embodiments the step may comprise single-cell genomic profiling, single-cell RNA profiling, single-cell DNA profiling, single-cell epigenomic profiling, single-cell protein profiling, or single-cell reporter gene expression profiling. Proteins that may be used to alter genomic and epigenomic state are described in Shmakov et al., 2015, *Molecular Cell* 60, 1–13 and Zetsche et al., 2015, *Cell* 163, 759–771.

[00298] In some embodiments the step may comprise single-cell RNA abundance analysis, single-cell transcriptome analysis, single-cell exome analysis, single-cell transcription rate analysis, or single-cell RNA degradation rate analysis.

[00299] In some embodiments the step may comprise single-cell DNA abundance analysis, single-cell DNA methylation profiling, single-cell chromatin profiling, single-cell chromatin accessibility profiling, single-cell histone modification profiling, or single-cell chromatin indexing.

[00300] In some embodiments the step may comprise single-cell protein abundance analysis, single-cell post-translational protein modification analysis, or single-cell proteome analysis.

[00301] In some embodiments the step may comprise single-cell mRNA reporter analysis, detection or quantification.

[00302] In some embodiments the step may comprise single-cell dual molecular profiling, such any combination of two amongst single-cell RNA profiling, single-cell DNA profiling, single-cell protein profiling, mRNA reporter analysis.

[00303] The method of the invention may include at the step determining single cell RNA levels. For single cell RNA-Seq (scRNA-Seq), one may use Drop-Seq (Macosko, E. Z., Basu, A., Satija, R., Nemesh, J., Goldman, M., Tirosh, I., Bialas, A. R., Kamitaki, N., Sanes, J. R.,

Weitz, D. A., Shalek, A. K., Regev, A. & McCarroll, S. A. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell*. 2015 May 21;161(5):1202-14. doi: 10.1016/j.cell.2015.05.002. PMID:4481139) and variants thereof. This technique relies on reverse-emulsion, early barcoding for analyzing 10^4 – 10^6 cells/experiment at very low cost. Drop-Seq enables to co-encapsulate individual cells with uniquely barcoded mRNA capture beads in reverse emulsion droplets. After lysis and mRNA capture, the emulsion is broken and all beads/cells are processed (RT, library prep) together, deconvolving each cell's profile from bead barcodes. In some embodiments, droplets can compartmentalize hundreds of cells/sec, are stable over time and to heat, and can serve as micro-vessels to add reagents; after RT, barcoded beads are stable and can be sorted or subselected. Sampling noise from shallow read depth is substantially lower than the technical variability between cells (Shalek, A. K., Satija, R., Shuga, J., Trombetta, J. J., Gennert, D., Lu, D., Chen, P., Gertner, R. S., Gaublomme, J. T., Yosef, N., Schwartz, S., Fowler, B., Weaver, S., Wang, J., Wang, X., Ding, R., Raychowdhury, R., Friedman, N., Hacohen, N., Park, H., May, A. P. & Regev, A. Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. *Nature*. 510, 363-369, doi:10.1038/nature13437 (2014). PMID:4193940.), so one may sufficiently estimate expression with ~100,000 reads per cell for many applications (especially with a 5' or 3'-end protocol, Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nature biotechnology*. 33, 495-502, doi:10.1038/nbt.3192 (2015)).

[00304] Single cell RNA may also be analyzed as described in Klein, A. M., Mazutis, L., Akartuna, I., Tallapragada, N., Veres, A., Li, V., Peshkin, L., Weitz, D.A., Kirschner, M. W. Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell*. 2015 May 21;161(5):1187-201. doi: 10.1016/j.cell.2015.04.044. PMID: 4441768.

[00305] The method of the invention may include determining RNA transcription and degradation rates. One may use RNA metabolically labeled with 4-thiouridine, to measure RNA transcription and degradation rates (Rabani, M., Raychowdhury, R., Jovanovic, M., Rooney, M., Stumpo, D. J., Pauli, A., Hacohen, N., Schier, A. F., Blackshear, P. J., Friedman, N., Amit, I. & Regev, A. High-resolution sequencing and modeling identifies distinct dynamic RNA regulatory strategies. *Cell*. 159, 1698-1710, doi:10.1016/j.cell.2014.11.015 (2014). PMID:4272607; Rabani, M., Levin, J. Z., Fan, L., Adiconis, X., Raychowdhury, R., Garber, M., Gnirke, A., Nusbaum, C., Hacohen, N., Friedman, N., Amit, I. & Regev, A. Metabolic labeling of RNA

uncovers principles of RNA production and degradation dynamics in mammalian cells. *Nature biotechnology*. 29, 436-442, doi:10.1038/nbt.1861 (2011). PMID:3114636).

[00306] The method of the invention may include a step of determining DNA methylation. One may apply methods for reduced representation (RRBS), targeted capture, and whole genome bisulfite sequencing of DNA methylation from bulk to ultra-low inputs (Chan, M. M., Smith, Z. D., Egli, D., Regev, A. & Meissner, A. Mouse ooplasm confers context-specific reprogramming capacity. *Nature genetics*. 44, 978-980, doi:10.1038/ng.2382 (2012). PMID:3432711; Smith, Z. D., Chan, M. M., Humm, K. C., Karnik, R., Mekhoubad, S., Regev, A., Eggan, K. & Meissner, A. DNA methylation dynamics of the human preimplantation embryo. *Nature*. 511, 611-615, doi:10.1038/nature13581 (2014). PMID:4178976; Smith, Z. D., Chan, M. M., Mikkelsen, T. S., Gu, H., Gnirke, A., Regev, A. & Meissner, A. A unique regulatory phase of DNA methylation in the early mammalian embryo. *Nature*. 484, 339-344, doi:10.1038/nature10960 (2012). PMID:3331945) to single cells.

[00307] The method of the invention may include a step determining Chromatin accessibility. This may be performed by ATAC-Seq. For massively parallel single cell ATAC-Seq one may implement a droplet-based assay. First, in-tube, one may use Tn5 transposase to fragment chromatin inside isolated intact nuclei and add universal primers at cutting sites. Next, in-drop, one may use a high diversity library of barcoded primers to uniquely tag all DNA that originated from the same single cell. Alternatively, one may perform all steps in drop. One may also use a strategy that relies on split pooled nuclei barcoding in plates (Cusanovich, D. A., Daza, R., Adey, A., Pliner, H., Christiansen, L., Gunderson, K. L., Steemers, F. J., Trapnell, C. & Shendure, J. Multiplex single-cell profiling of chromatin accessibility by combinatorial cellular indexing. *Science*. 2015 May 22;348(6237):910-4. doi: 10.1126/science.aab1601. Epub 2015 May 7). Applicants have optimized key steps in a mixture of human and mouse cells, with specificity that exceeds the initial performance of mRNA Drop-Seq. Applicants have also used a Fluidigm C1 protocol (see www.fluidigm.com/products/c1-system) to analyze ~100 single DCs, closely reproducing ensemble measures, high enrichment in TSSs, and nucleosome-like periodicity.

[00308] ATAC-seq (assay for transposase-accessible chromatin) identifies regions of open chromatin using a hyperactive prokaryotic Tn5-transposase, which preferentially inserts into accessible chromatin and tags the sites with sequencing adaptors (Pott and Lieb *Genome Biology* (2015) 16:172 DOI 10.1186/s13059-015-0737-7 and Buenrostro JD, Giresi PG, Zaba LC, Chang

HY, Greenleaf WJ. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat Methods*. 2013;10:1213–128). Two very different approaches were used: one relied on physical isolation of single cells (Buenrostro JD, Wu B, Litzenburger UM, Ruff D, Gonzales ML, Snyder MP, et al. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*. 2015;523:486–90), and the other avoided single-cell reaction volumes by using a two-step combinatorial indexing strategy (Cusanovich DA, Daza R, Adey A, Pliner HA, Christiansen L, Gunderson KL, et al. Epigenetics. Multiplex single-cell profiling of chromatin accessibility by combinatorial cellular indexing. *Science*. 2015;348:910–4).

[00309] In the indexing scheme, Cusanovich et al. [Cusanovich DA, Daza R, Adey A, Pliner HA, Christiansen L, Gunderson KL, et al. Epigenetics. Multiplex single-cell profiling of chromatin accessibility by combinatorial cellular indexing. *Science*. 2015;348:910–4] lysed cells, and 2500 nuclei were placed into each well of a 96-well plate. Transposases loaded with unique adaptors were added to each well, creating 96 pools of approximately 2500 nuclei, each pool with distinct barcodes. Nuclei from all of the transposition reactions were mixed, and using a fluorescence-activated cell sorter (FACS) 15–25 nuclei were deposited into each well of a second 96-well plate. Nuclei in each well of this second plate were lysed, and the DNA was amplified using a primer containing a second barcode. The low number of nuclei per well ensured that about 90 % of the resulting barcode combinations were unique to a single cell. This combinatorial indexing strategy enabled the recovery of 500–1500 cells with unique tags per experiment. Overall Cusanovich et al. obtained scATAC-seq data from over 15,000 individual cells from mixtures of GM12878 lymphoblastoid cells with HEK293, HL-60, or mouse Patski cells. The number of reads associated with any single cell was very low, varying from 500 to about 70,000 with a median of fewer than 3000 reads per cell.

[00310] Buenrostro et al. [Buenrostro JD, Wu B, Litzenburger UM, Ruff D, Gonzales ML, Snyder MP, et al. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*. 2015;523:486–90] used a programmable microfluidic device (C1, Fluidigm) to isolate single cells and perform ATAC-seq on them in nanoliter reaction chambers. Each nanochamber was analyzed under a microscope to ensure that a single viable cell had been captured. This approach is simple and has the significant advantage of a carefully monitored reaction environment for each individual cell, although the throughput was limited to processing 96 cells

in parallel. Buenrostro et al. sampled 1632 cells from eight different cell lines, including GM12878, K562, and H1 cells, and obtained an average of 73,000 reads per cell, about 20 times the number of reads per cell obtained using the barcoding strategy.

[00311] The method of the invention may include a step of determining histone modifications and protein-DNA interactions. One may apply tools that use genomic barcoding to index chromatin prior to immunoprecipitation to enable multiplexed analysis of limited samples and individual cells in a single reaction. For single-cell chromatin profiling, one may use Drop-ChIP where the chromatin of individual cells is barcoded in droplets. Based on the Drop-Seq technique, one may encapsulate single cells, lyse and MNase-digest chromatin, then fuse a second droplet with barcoded oligos, ligate them to the fragmented chromatin, break the emulsion, add carrier chromatin, and carry out ChIP-Seq. this may be performed using a protocol with split-pool barcoding to collect 10^4 – 10^5 single cells/assay.

[00312] ChIP-sequencing, also known as ChIP-seq, is a method used to analyze protein interactions with DNA which may be used with perturbation. ChIP-seq combines chromatin immunoprecipitation (ChIP) with massively parallel DNA sequencing to identify the binding sites of DNA-associated proteins. It can be used to map global binding sites precisely for any protein of interest. ChIP-seq is used primarily to determine how transcription factors and other chromatin-associated proteins influence phenotype-affecting mechanisms. Determining how proteins interact with DNA to regulate gene expression is for understanding many biological processes and disease states. This epigenetic information is complementary to genotype and expression analysis. ChIP-seq technology is as an alternative to ChIP-chip which requires a hybridization array. Specific DNA sites in direct physical interaction with transcription factors and other proteins can be isolated by chromatin immunoprecipitation. ChIP produces a library of target DNA sites bound to a protein of interest in vivo. Massively parallel sequence analyses are used in conjunction with whole-genome sequence databases to analyze the interaction pattern of any protein with DNA, see, e.g., Johnson DS, Mortazavi A et al. (2007) Genome-wide mapping of in vivo protein–DNA interactions. *Science* 316: 1497–1502, or the pattern of any epigenetic chromatin modifications. This can be applied to the set of ChIP-able proteins and modifications, such as transcription factors, polymerases and transcriptional machinery, structural proteins, protein modifications, and DNA modifications. See, e.g., “Whole Genome Chromatin IP Sequencing,” Illumina, Inc (2010), available at

www.illumina.com/Documents/products/datasheets/datasheet_chip_sequence.pdf (Chromatin Immunoprecipitation with massively parallel sequencing).

[00313] For multiplex analysis of (limited) bulk samples, one may rely on chromatin indexing (MINT-ChIP; iChIP), where MNase-fragmented chromatin is indexed by ligation to a uniquely barcoded adaptor and then pooled and processed in multiplex through all subsequent phases, either with (MINT-ChIP) or without (iChIP: Lara-Astiaso, D., Weiner, A., Lorenzo-Vivas, E., Zaretsky, I., Jaitin, D. A., David, E., Keren-Shaul, H., Mildner, A., Winter, D., Jung, S., Friedman, N. & Amit, I. Immunogenetics. Chromatin state dynamics during blood formation. *Science*. 345, 943-949, doi:10.1126/science.1256271 (2014). PMID:4412442) carrier chromatin (without adaptors).

[00314] The method of the invention may include a step of determining proteins. Recently developed assays (e.g., CyTOF: Bendall, S. C., Simonds, E. F., Qiu, P., Amir el, A. D., Krutzik, P. O., Finck, R., Bruggner, R. V., Melamed, R., Trejo, A., Ornatsky, O. I., Balderas, R. S., Plevritis, S. K., Sachs, K., Pe'er, D., Tanner, S. D. & Nolan, G. P. Single-cell mass cytometry of differential immune and drug responses across a human hematopoietic continuum. *Science*. 332, 687-696, doi:10.1126/science.1198704 (2011). PMID:3273988), allow multiplexed, single cell detection of dozens of proteins in millions of cells, but rely on antibodies and cannot yet be combined with DNA readout. Conversely, mass spectrometry (LC-MS/MS) allows quantitative analysis of entire proteomes, but deep analysis requires large amounts of protein/cells. To measure single cell protein levels and post-translational modifications (PTMs), one may use one of three complementary antibody-based assays: (1) standard flow cytometry with a few proteins/PTMs, $>10^6$ single cells); (2) CyTOF (Bendall, S. C., Simonds, E. F., Qiu, P., Amir el, A. D., Krutzik, P. O., Finck, R., Bruggner, R. V., Melamed, R., Trejo, A., Ornatsky, O. I., Balderas, R. S., Plevritis, S. K., Sachs, K., Pe'er, D., Tanner, S. D. & Nolan, G. P. Single-cell mass cytometry of differential immune and drug responses across a human hematopoietic continuum. *Science*. 332, 687-696, doi:10.1126/science.1198704 (2011). PMID:3273988) (heavy metal labeling with multiplex barcoding; ~ 30 – 50 proteins/PTMs, 10^5 – 10^6 single cells); and (3) novel, highly multiplexed, DNA sequencing-based readouts of protein levels (100s proteins/PTMs; 10^6 cells). For sequencing based readouts, one may use one of two approaches, geared at detecting hundreds of proteins in single cells: Immuno-Seq (when antibodies can be washed out: Niemeyer, C. M., Adler, M. & Wacker, R. Detecting antigens by quantitative

immuno-PCR. *Nat Protoc.* 2, 1918-1930, doi:10.1038/nprot.2007.267 (2007)) and proximity extension assays (PEA, when antibodies cannot be washed away: Hammond, M., Nong, R. Y., Ericsson, O., Pardali, K. & Landegren, U. Profiling cellular protein complexes by proximity ligation with dual tag microarray readout. *PLoS One.* 7, e40405, doi:10.1371/journal.pone.0040405 (2012). PMID:3393744; Nong, R. Y., Wu, D., Yan, J., Hammond, M., Gu, G. J., Kamali-Moghaddam, M., Landegren, U. & Darmanis, S. Solid-phase proximity ligation assays for individual or parallel protein analyses with readout via real-time PCR or sequencing. *Nat Protoc.* 8, 1234-1248, doi:10.1038/nprot.2013.070 (2013); Stahlberg, A., Thomsen, C., Ruff, D. & Aman, P. Quantitative PCR analysis of DNA, RNAs, and proteins in the same single cell. *Clin Chem.* 58, 1682-1691, doi:10.1373/clinchem.2012.191445 (2012). These use DNA-sequence based encoding, and are compatible with other genomic readouts (e.g., sgRNA barcodes). DNA-sequence tags can be conjugated to antibodies (Janssen, K. P., Knez, K., Spasic, D. & Lammertyn, J. Nucleic acids for ultra-sensitive protein detection. *Sensors (Basel).* 13, 1353-1384, doi:10.3390/s130101353 (2013). PMID:3574740), nanobodies (Pardon, E., Laeremans, T., Triest, S., Rasmussen, S. G., Wohlkonig, A., Ruf, A., Muyldermans, S., Hol, W. G., Kobilka, B. K. & Steyaert, J. A general protocol for the generation of Nanobodies for structural biology. *Nat Protoc.* 9, 674-693, doi:10.1038/nprot.2014.039 (2014). PMID:4297639; Theile, C. S., Witte, M. D., Blom, A. E., Kundrat, L., Ploegh, H. L. & Guimaraes, C. P. Site-specific N-terminal labeling of proteins using sortase-mediated reactions. *Nat Protoc.* 8, 1800-1807, doi:10.1038/nprot.2013.102 (2013). PMID:3941705) or aptamers (Janssen, K. P., Knez, K., Spasic, D. & Lammertyn, J. Nucleic acids for ultra-sensitive protein detection. *Sensors (Basel).* 13, 1353-1384, doi:10.3390/s130101353 (2013). PMID:3574740). Detection of specific transcripts and proteins in single cells may be performed as well (Frei et al., Highly multiplexed simultaneous detection of RNAs and proteins in single cells. *Nat Methods.* 2016 Mar;13(3):269-75.)

[00315] In certain embodiments of the present invention, quantitative measurements of both the copy number and spatial distribution of large fractions of the transcriptome in single cells is measured by multiplexed error-robust fluorescence in situ hybridization (MERFISH) (Moffitt and Zhuang, Chapter One – RNA Imaging with Multiplexed Error-Robust Fluorescence In Situ Hybridization (MERFISH), *Methods in Enzymology*. Volume 572, 2016, Pages 1–49).

Multiplex Analysis of Single Cell Constituents

[00316] It is an object of the present invention to provide a method for the high multiplex analysis of cellular constituents by linking nucleic acids tags to existing ligand binding and/or antibody technologies to enable proteomic or cellular constituent detection and relative quantification by next-generation sequencing (NGS) in single cells or isolated aggregations of cellular constituents. The present invention may combine perturbation of single cells followed by protein analysis in the single cells, Thus, protein expression may be linked to a perturbation.

[00317] It is a further object of the present invention to provide for comparing high multiplex protein data variation between single cells or isolated aggregation of cellular constituents and between different biological conditions (e.g. healthy vs. diseased states; one genetic perturbation vs. another, different genetic backgrounds).

[00318] It is a further object of the present invention to provide massively parallel profiling of all circuit aspects in single cells or isolated aggregations of cellular constituents: from RNA to chromatin organization to protein levels.

[00319] In a first aspect, the present invention provides a method of assaying segregated cellular constituents, comprising: admixing at least one isolated aggregation of cellular constituents with monomers of a polymerizable gel; polymerizing the gel, to embed the cellular constituents in discrete polymer matrices; incubating the cellular constituents embedded in the polymer matrices with one or more labeling ligands with specific binding affinity for one or more target cellular constituents to produce one or more labeled cellular constituents in the polymer matrices, wherein each of the one or more labeling ligand comprises a bound oligonucleotide label comprising a unique constituent identifier (UCI) sequence, and wherein the incubation comprises binding conditions under which the labeling ligand will bind to the cellular constituent within the polymer matrix, and the incubation further comprises washing conditions under which unbound labeling ligands will be washed out of the polymer matrix; and sequencing the oligonucleotide label, whereby detecting the UCI by sequencing indicates the presence of the target cellular constituent.

[00320] Cellular constituents may include any molecule within a cell; i.e. proteins, nucleic acids, or post translational modifications (PTM). The cellular constituent may be a protein, RNA transcript, metabolite, or a DNA molecule. Specific cellular constituents may be proteins, modified proteins, hormones, cytokines, cellular metabolites, or carbohydrates. The isolated aggregation of cellular constituents may be a cell, an extracellular vesicle, an organelle, or an

organized subcomponent thereof, including molecular complexes. Isolated aggregations of cellular constituents may include separate organelles of a single cell or separate organelles isolated from a population of cells. Organelles may be for example, mitochondria, nuclei, or cellular vesicles. In one embodiment, a specific type of single cells may be isolated. In one embodiment, immune cells are isolated from a population of cells. Not being bound by a theory, single mitochondria can be purified from a population of cells and the relative amounts of constituents present in each individual mitochondrion may be analyzed. Not being bound by a theory, immune cells may be isolated by a method such as cell sorting and the relative representation of cellular constituents may be determined for each individual cell.

[00321] The step of admixing the isolated aggregation of cellular constituents with monomers may be carried out in an aqueous solution, or in an aqueous aliquot or droplet present in an oil emulsion. The polymer matrix may be a hydrogel. The polymer matrix may be any hydrogel capable of polymerization to create a solid matrix that fixes the cellular constituents and provides a porosity capable of allowing labeling ligands to freely diffuse through the network of pores. The cellular constituents may be further fixed by treating with an aldehyde. The aldehyde may be formaldehyde, paraformaldehyde, or glutaraldehyde. Not being bound by a theory the fixation in a solid matrix prevents the mixing of the cellular constituents between the isolated aggregations of cellular constituents. Not being bound by a theory, capturing cellular constituents in a solid polymer mesh insures that they are physical units that can be ligand and/or antibody stained as a pool and isolated as single cells or isolated aggregates of cellular constituents subsequently. Not being bound by a theory, the fixing of cellular constituents in the polymer matrix allows access to the labeling ligands to intracellular constituents.

[00322] The physical units formed by the polymer matrix may be particles, droplets, or a continuous polymer matrix with discrete regions comprising the isolated aggregates of cellular constituents. Therefore, the polymer matrix may include more than one isolated aggregate of cellular constituents. The polymer matrix may be divided such that isolated aggregates of cellular constituents are separable. The polymer matrix may be separable in that individual particles, droplets, or sections can be isolated. They may be isolated by physical manipulation using a sorting device. The sorting device may use microfluidics. They may be separated by use of dilution or manual manipulation by a user. They may be separated by use any kind of (micro) dissection. The cellular constituents within the polymer matrix may be stained with a dye, or a

dye-conjugated ligand indicating the location of individual cellular constituents or cells. The polymer matrix may be punched to isolate a core, wherein each core from the polymer matrix contains a single isolated aggregate of cellular constituents. Not being bound by a theory, the fixation of isolated aggregates of cellular constituents in a matrix allows each isolated aggregate of cellular constituents to be compartmentalized wherein the separate compartments can be treated in a single experimental vessel or container and separated subsequently.

[00323] The labeling ligands are linked with an oligonucleotide label that can be used to determine the identity of the ligand. Each oligonucleotide label may comprise a unique constituent identifier (UCI) that can be used to determine the presence of a cellular constituent. Not being bound by a theory, the availability of unique sequences allows the labeling and detection of a plurality of ligands each for a specific constituent. Not being bound by a theory, the UCI allows a DNA readout for detection of a cellular constituent. The DNA readout may be by any sequencing method or method of amplification, such as by PCR or next generation sequencing. The oligonucleotide label may additionally include a promoter for amplification by an RNA polymerase, such as T7 polymerase. Not being bound by a theory, amplification by T7 polymerase allows amplification of low represented sequences, whereas such sequences may be diluted out by domination of a higher represented sequence during PCR. Not being bound by a theory, the labeling of each labeling ligand with a unique UCI allows the identification of more than ten, or hundred, or thousands of cellular constituents in an isolated aggregation of cellular constituents.

[00324] The method may further comprise segregating the discrete polymer matrices comprising the labeled constituents before the step of sequencing. Segregating the discrete polymer matrices may be by sorting single discrete matrices into separate reaction vessels. Segregating the discrete polymer matrices may be by forming discrete unique-identifier-transfer compositions, each comprising the cellular constituents embedded in a discrete polymer matrix and a transfer particle, wherein: the ligand oligonucleotide label further comprises a capture sequence, and the UCI and capture sequence are together releasably attached to the labeling ligand; the labelling ligand is bound to the target cellular constituent; and, the transfer particle comprises: a capture-binding-sequence having specific binding affinity for the capture sequence attached to the UCI, and, a unique source identifier (USI) sequence that is unique to each transfer particle. The USI of each transfer particle preferably comprises 4-15 nucleotides. The method

may further comprise releasing the UCI from the labeled ligand, under conditions within the unique-identifier-transfer composition so that the released capture sequence binds to the capture-binding-sequence on the transfer particle, thereby transferring the UCI to the transfer particle. The transfer particle may be a solid bead. The transfer particle may be a hydrogel bead. The transfer particle may also be used to capture nucleic acids present in a discrete polymer matrix. The nucleic acids may be RNA and/or DNA. Not being bound by a theory the transfer particle may be used to capture both the UCI and the nucleic acids, whereby the source of the bound cellular constituents and nucleic acids can be determined after sequencing.

[00325] The method may further comprise, before the sequencing step, generating a USI for each discrete polymer matrix by a split pool ligation method, wherein the oligonucleotide label further comprises a universal ligation handle (ULH) sequence configured to produce a DNA overhang capable of hybridization to a complementary overhang on a first index nucleotide sequence, wherein the first index nucleotide sequence comprises an overhang complementary to a final index sequence or optionally a middle index sequence, wherein the middle index sequence comprises overhangs complementary to the first index sequence and to the final index sequence or optionally to another middle index sequence and final index sequence, wherein the final index sequence has a single overhang complementary to the preceding index sequence, and wherein the first, middle, and final index sequences are selected from a plurality of unique sequences comprising compatible DNA overhangs and 10 to 30 base pairs of unique sequence.

[00326] The split pool ligation method may comprise: splitting the pool of discrete polymer matrices into separate pools of polymer matrices, each containing a unique first index sequence; ligating the first index sequence to the ligation handle; pooling the discrete polymer matrices; optionally, splitting the pool of discrete polymer matrices into separate pools each containing a unique middle index sequence; ligating the middle index sequence to the first index sequence; and pooling the discrete polymer matrices; optionally, repeating the steps with another middle index sequence; splitting the pool of discrete polymer matrices into pools containing a unique final index sequence; and ligating the final index sequence to the preceding index sequence, whereby each discrete polymer matrix comprises a USI. The USI may have no middle index sequence, one middle index sequence, two middle index sequences, preferably the USI has a first, middle, and final index sequence. Not being bound by a theory, the size of the unique sequences in each index determines the amount included. Not being bound by a theory, the

number of indices selected is the amount necessary such that the probability of having identical USI sequences on sparse polymer matrices is approaching zero. In an exemplary embodiment, each index includes 192 unique sequences.

[00327] The ligation handle may comprise a restriction site for producing an overhang complementary with a first index sequence overhang, and wherein the method further comprises digestion with a restriction enzyme. The ligation handle may comprise a nucleotide sequence complementary with a ligation primer sequence and wherein the overhang complementary with a first index sequence overhang is produced by hybridization of the ligation primer to the ligation handle. Additionally, the ULH may comprise a dsDNA part that already includes the overhang needed for index ligation.

[00328] The UCI may comprise 4 to 30 nucleotides or 7 to 30 nucleotides, preferably about 21 nucleotides. The oligonucleotide label may further comprise a unique molecular identifier (UMI) sequence. The first, middle, or final index sequence may further comprises a unique molecular identifier (UMI) sequence. The UMI may comprise 4-20 nucleotides. The UMI may comprise 8 to 16 nucleotides.

[00329] The isolated aggregation of cellular constituents may be a cell, an extracellular vesicle, an organelle, or an organized subcomponent thereof.

[00330] The sequencing may comprise combining a primer having a unique source identifier (USI) sequence with UCI, so that the USI and UCI sequences are sequenced together, and the USI preferably comprises 20 to 120 nucleotides.

[00331] The step of admixing the isolated aggregation of cellular constituents with monomers may be carried out in an aqueous aliquot or in a droplet formed by an aqueous solution in oil emulsion. The aqueous aliquot may be a separate reaction vessel such as a well in a plate. The droplet may be formed by a microfluidic device. The polymer matrix may be a hydrogel. The method may be a multiplex assay with a plurality of labeling ligands, each labeling ligand have a distinct UCI. The labeling ligand may be non-covalently bound to the target cellular constituent.

[00332] The method may further comprise pooling the oligonucleotide labels comprising a USI from a plurality of polymer matrices and sequencing the pooled UCI sequences and USI sequences. The method may further comprise pooling the oligonucleotide labels comprising a USI and UMI from a plurality of polymer matrices and sequencing the pooled UCI sequences, USI sequences, and UMI sequences.

[00333] The method may further comprise washing the cellular constituents embedded in the polymer matrices to remove selected cellular components from the polymer matrices before incubating the cellular constituents with the labeling ligand. The washing may comprise treating the cellular constituents embedded in the polymer matrices with a detergent so as to remove lipids from the polymer matrices before incubating the cellular constituents with the labeling ligand. The detergent may be an anionic detergent or nonionic detergent. The detergent may be SDS, NP-40, triton X-100, or any other detergent known in the art capable of removing lipids.

[00334] The method may further comprise quantitating the relative amount of the UCI sequence associated with a first aggregation of cellular constituents to the amount of the same UCI sequence associated with a second aggregation of cellular constituents, whereby the relative differences of a cellular constituent between aggregations of cellular constituents are determined. The relative amount may be compared to a control sample. The control sample may have predetermined amounts of cellular constituents. There may be more than one control sample. There may be at least three control samples. The at least three control samples can be used to generate a standard curve upon which all of the other cellular constituents within discrete polymer matrices are compared. The control sample may comprise isolated aggregations of cellular constituents that were untreated as compared relative to isolated aggregations of cellular constituents that were treated with a different condition. Cells may be treated with drugs, small molecules, pathogens, hormones, cytokines, proteins, nucleic acids, virus particles, or grown in different cellular environments. Cells may be isolated from a diseased tissue. The cells from the diseased tissue may be compared to cells from non-diseased tissue. Cells may be treated with systems that knockout, decrease or increase expression of a gene. Cells may be treated with systems that knockout functional elements of a genome. Functional elements include, but are not limited to promoters, enhancers, repressors, centromeres, or telomeres. CRISPR systems may be used.

[00335] The labeling ligand may be an antibody or an antibody fragment. The antibody fragment may be a nanobody, Fab, Fab', (Fab')₂, Fv, ScFv, diabody, triabody, tetrabody, Bis-scFv, minibody, Fab₂, or Fab₃ fragment. The labeling ligand may be an aptamer. The labeling ligand may be a nucleotide sequence complementary to a target sequence.

[00336] The method may comprise multiplex binding of two or more labeling ligands to each aggregation of cellular constituents. The two or more distinct labeling ligands may comprise

complementary oligonucleotide sequences, so that binding of the labeling ligands to respective target cellular constituents that are in proximity permits the complementary sequences of the distinct ligands to hybridize, forming an amplifiable polynucleotide duplex. The method may further comprise amplifying the polynucleotide duplex to provide an amplified sequence that is a detectable signal that target cellular constituents are in proximity. The complementary oligonucleotide sequences, which serve as a start site for polymerase extension, can either be designed to query proximity of two specific cellular constituents, or it can be designed to be universal, thereby querying interactions between all members of the labeling ligand panel.

[00337] In one embodiment, at least two distinct labeling ligands comprise oligonucleotide sequences configured to be ligated, so that binding of the labeling ligands to respective target cellular constituents that are in proximity permits the oligonucleotide sequences of the distinct ligands to ligate, forming an amplifiable polynucleotide duplex. The method may further comprise amplifying the polynucleotide duplex to provide an amplified sequence that is a detectable signal that target cellular constituents are in proximity.

[00338] One of the labeling ligands may comprise an oligonucleotide label with a restriction enzyme site between the labeling ligand and the UCI, and wherein the method may further comprise treating with a restriction enzyme, whereby the UCI from the labeling ligand is transferred to the oligonucleotide label of the labeling ligand in proximity.

[00339] The method may further comprise labeling the aggregation of cellular constituents by fluorescent in situ hybridization.

[00340] The aggregation of cellular constituents may be a cell that is a member of a cell population. The cell may be transformed or transduced with one or more genomic sequence-perturbation constructs that perturb a genomic sequence in the cells, wherein each distinct genomic sequence-perturbation construct comprises a unique-perturbation-identifier (UPI) sequence unique to that construct. The genomic sequence-perturbation construct may comprise a sequence encoding a guide RNA sequence of a CRISPR-Cas targeting system. The method may further comprise multiplex transformation of the population of cells with a plurality of genomic sequence-perturbation constructs. The UPI sequence may be attached to a perturbation-sequence-capture sequence, and the microbeads may comprise a perturbation-sequence-capture-binding-sequence having specific binding affinity for the perturbation-sequence-capture sequence attached to the UPI sequence. The UPI sequence may be attached to a universal ligation handle

sequence, whereby a USI may be generated by split-pool ligation. The method may further comprise multiplex sequencing of the pooled UCI sequences, USI sequences, and UPI sequences.

[00341] The oligonucleotide label may comprise a regulatory sequence configured for amplification by an RNA polymerase, such as T7 polymerase. The labeling ligands may comprise oligonucleotide sequences configured to hybridize to a transcript specific region. The oligonucleotide label may further comprise attachment chemistry, such as an acrylic phosphoramidite modification, whereby the modification allows for incorporation into the polymer matrices upon polymerization. The acrylic phosphoramidite may be Acrydite™ (Eurofins Scientific, Luxembourg). The method may further comprise amplification of the oligonucleotide label and USI by PCR or T7 amplification before sequencing. T7 amplification may be followed by cDNA generation and optionally amplification by PCR. The oligonucleotide label may further comprise at least one spacer sequence, preferably two spacer sequences. The oligonucleotide label may further comprise a photocleavable linker. The oligonucleotide label may further comprise a restriction enzyme site between the labeling ligand and UCI.

[00342] The discrete polymer matrices may be labeled and washed more than once. Discrete polymer matrices may be labeled with a marker specific for a cell type or cell cycle marker or developmental marker, or differentiation marker, or disease marker. The label may be a fluorescent label. The fluorescent label may be used to separate the discrete polymer matrices into distinct groups. The label may be used to identify a certain cell type prior to embedding it into a discrete polymer matrix. The discrete polymer matrices of a distinct group may then be labeled again with labeling ligands that contain an oligonucleotide label of the present invention. After novel information is obtained from the multiplex assay of the present invention, a 'banked' population of polymer matrices can be stained for newly identified markers and the population of interest can be sorted (enriched) for, and investigated more deeply.

[00343] In another aspect, the present invention provides a method of determining open chromatin in individual cells comprising: isolating single cells into droplets formed by an aqueous solution in oil emulsions, wherein the droplets further comprise Tn5-transposase loaded with two tagmentation adapters, wherein one adapter is configured for incorporation into a polymer matrix and the second adapter is configured with a ligation handle for generating a USI; incubating the droplets to allow cell lysis and tagmentation of open chromatin; infusing

monomers of a polymerizable gel into the droplets; polymerizing the gel, to embed the cellular constituents in discrete polymer matrices; optionally incubating the polymer matrices with one or more labeling ligands with specific binding affinity for one or more target cellular constituents to produce one or more labeled cellular constituents in the polymer matrices, wherein each of the one or more labeling ligand comprises a bound oligonucleotide label comprising a unique constituent identifier (UCI) sequence and a sequence capable of hybridization to the tagmentation adapter configured for incorporation into a polymer matrix, and wherein the incubation comprises binding conditions under which the labeling ligand will bind to the cellular constituent within the polymer matrix and the oligonucleotide label will hybridize to said tagmentation adapter, and wherein the incubation further comprises washing conditions under which unbound labeling ligands will be washed out of the polymer matrix; and extending the genomic DNA and adapter DNA, whereby a continuous DNA strand is generated comprising the adapters, genomic DNA, and DNA overhang; optionally the oligonucleotide label bound to a labeling ligand; generating a USI at the DNA overhang by split-pool ligation; sequencing the continuous DNA strand, whereby open chromatin is determined and optionally the presence of a cellular constituent at a site of open chromatin is determined.

[00344] In another aspect, the present invention provides a method of measuring RNA levels in individual cells comprising: isolating single cells into droplets formed by an aqueous solution in oil emulsions, wherein the droplets comprise at least one labeling ligands specific for binding at one or more target RNA transcripts, wherein the labeling ligands are configured for incorporation into a polymer matrix and comprise a ligation handle for generating a USI; lysing the cells in the droplets under conditions wherein the labeling ligands will bind to the target RNA transcripts; injecting monomers of a polymerizable gel into the droplets; polymerizing the gel, to embed the labeling ligands in discrete polymer matrices; optionally, staining the discrete polymer matrices with at least one additional labeling ligand; generating a USI by split-pool ligation; and sequencing the resulting DNA, whereby RNA levels and optionally protein levels are determined in single cells. The droplets may comprise at least one pair of labeling ligands specific for binding at adjacent sites of one or more target RNA transcripts, wherein each pair of labeling ligands comprises one labeling ligand configured for incorporation into a polymer matrix and one labeling ligand comprising a ligation handle for generating a USI, and wherein the method may further comprise injecting a ligation reaction buffer comprising a ligase that is

configured to allow ligation of the pair of labeling ligands if they are hybridized adjacently with single nucleotide resolution on the target RNA transcript, such that off target binding of labeling ligand does not get ligated, and will not be amplified in subsequent steps.

[00345] In another aspect, the present invention provides a method of assaying segregated cellular constituents, comprising: fixing and permeablizing at least one cell; incubating the fixed and permeablized cell(s) with one or more labeling ligands with specific binding affinity for one or more target cellular constituents to produce one or more labeled cell(s), wherein each of the one or more labeling ligands comprise a bound oligonucleotide label comprising a unique constituent identifier (UCI) sequence, and wherein the incubation comprises binding conditions under which the labeling ligand will bind to the cellular constituent within the cell(s), and the incubation further comprises washing conditions under which unbound labeling ligands will be washed from the cell(s); admixing the cell(s) with monomers of a polymerizable gel; isolating single cells into droplets formed by an aqueous solution in oil emulsions; polymerizing the gel, to embed the labeling ligands and other cellular constituents in discrete polymer matrices; optionally, staining the discrete polymer matrices with at least one additional labeling ligand; generating a USI by split-pool ligation; and sequencing the oligonucleotide label, whereby detecting the UCI by sequencing indicates the presence of the target cellular constituent. The labeling ligands in step (b) may comprise at least one pair of labeling ligands specific for binding at adjacent sites of one or more target RNA transcripts, wherein each pair of labeling ligands comprises one labeling ligand configured for incorporation into a polymer matrix and one labeling ligand comprising a ligation handle for generating a USI, and wherein the method further comprises ligating the pair of labeling ligands if they are within proximity after binding to the target RNA transcripts. Any of the preceding methods may comprise polymer matrices wherein they further comprise magnetic particles. In one embodiment, any hydrogel droplet encapsulated aggregations of cellular constituents may further comprise magnetic particles embedded into the droplets. Not being bound by a theory, magnetic particles enable magnetic separation, aiding in clean up and washing steps in multiple reactions. Not being bound by a theory, the use of magnetic particles greatly enhances automation and therefore throughput.

[00346] In another aspect, the present invention provides a method of assaying segregated cellular constituents, comprising: fixing and permeablizing at least one cell; incubating the fixed and permeablized cell(s) with one or more labeling ligands with specific binding affinity for one

or more target cellular constituents to produce one or more labeled cell(s), wherein each of the one or more labeling ligands comprise a bound oligonucleotide label comprising a unique constituent identifier (UCI) sequence, and wherein the incubation comprises binding conditions under which the labeling ligand will bind to the cellular constituent within the cell(s), and the incubation further comprises washing conditions under which unbound labeling ligands will be washed from the cell(s); and sequencing the oligonucleotide label, whereby detecting the UCI by sequencing indicates the presence of the target cellular constituent. The cellular constituent may comprise a protein, RNA transcript, or a DNA molecule. The method may further comprise segregating the cell(s) before sequencing. The segregating the cell(s) may comprise sorting the single cell(s) into a separate reaction vessel(s). The segregating the cell(s) may comprise forming discrete unique-identifier-transfer compositions, each comprising a cell and a transfer particle, wherein: the oligonucleotide label further comprises a capture sequence, and the UCI and capture sequence are together releasably attached to the labeling ligand; the labelling ligand is bound to the target cellular constituent; and, the transfer particle comprises: a capture-binding-sequence having specific binding affinity for the capture sequence attached to the UCI, and, a unique source identifier (USI) sequence that is unique to each transfer particle, and the USI preferably comprises 4-15 nucleotides. The method may further comprise releasing the UCI from the labeled ligand, under conditions within the unique-identifier-transfer composition so that the released capture sequence binds to the capture-binding-sequence on the transfer particle, thereby transferring the UCI to the transfer particle. The method may further comprise, before sequencing in step, generating a USI for each cell(s) by a split pool ligation method, wherein the oligonucleotide label further comprises a universal ligation handle (ULH) sequence configured to produce a DNA overhang capable of hybridization to a complementary over hang on a first index nucleotide sequence, wherein the first index nucleotide sequence comprises an overhang complementary to a final index sequence or optionally a middle index sequence, wherein the middle index sequence comprises overhangs complementary to the first index sequence and to the final index sequence or optionally to another middle index sequence and final index sequence, wherein the final index sequence has a single overhang complementary to the preceding index sequence, and wherein the first, middle, and final index sequences are selected from a plurality of unique sequences comprising compatible DNA overhangs and 10 to 30 base pairs of unique sequence. The split pool ligation method may comprise: splitting the pool of

cell(s) into separate pools of cell(s), each containing a unique first index sequence; ligating the first index sequence to the ligation handle; pooling the cell(s); optionally, splitting the pool of cell(s) into separate pools each containing a unique middle index sequence; ligating the middle index sequence to the first index sequence; and pooling the cell(s); optionally, repeating with another middle index sequence; splitting the pool of cell(s) into pools containing a unique final index sequence; and ligating the final index sequence to the preceding index sequence, whereby each cell comprises a USI.

[00347] The ligation handle may comprise a restriction site for producing an overhang complementary with a first index sequence overhang, and wherein the method further comprises digestion with a restriction enzyme. The ligation handle may comprise a nucleotide sequence complementary with a ligation primer sequence and wherein the overhang complementary with a first index sequence overhang is produced by hybridization of the ligation primer to the ligation handle.

[00348] The UCI may comprise 4 to 30 nucleotides, or 7 to 30 nucleotides, or about 21 nucleotides. The oligonucleotide label may further comprise a unique molecular identifier (UMI) sequence. The first, middle, or final index sequence may further comprise a unique molecular identifier (UMI) sequence. The UMI may be 4-20 nucleotides. The UMI may be 8 to 16 nucleotides.

[00349] The sequencing may comprise combining a primer having a unique source identifier (USI) sequence with UCI, so that the USI and UCI sequences are sequenced together, and the USI preferably comprises 20 to 120 nucleotides.

[00350] The method may comprise a multiplex assay with a plurality of labeling ligands, each labeling ligand have a distinct UCI. The labeling ligand may be non-covalently bound to the target cellular constituent. The method may further comprise pooling the oligonucleotide labels comprising a USI from a plurality of cells and sequencing the pooled UCI sequences and USI sequences. The method may further comprise pooling the oligonucleotide labels comprising a USI and UMI from a plurality of cells and sequencing the pooled UCI sequences, USI sequences, and UMI sequences. The method may further comprise quantitating the relative amount of the UCI sequence associated with a first cell to the amount of the same UCI sequence associated with a second cell, whereby the relative differences of a cellular constituent between cell(s) are determined.

[00351] The labeling ligand may be an antibody or an antibody fragment. The antibody fragment may be a nanobody, Fab, Fab', (Fab')₂, Fv, ScFv, diabody, triabody, tetrabody, Bis-scFv, minibody, Fab₂, or Fab₃ fragment. The labeling ligand may be an aptamer. The labeling ligand may be a nucleotide sequence complementary to a target sequence.

[00352] The method may comprise multiplex binding of two or more labeling ligands to the cellular constituents. At least two distinct labeling ligands may comprise complementary oligonucleotide sequences, so that binding of the labeling ligands to respective target cellular constituents that are in proximity permits the complementary sequences of the distinct ligands to hybridize, forming an amplifiable polynucleotide duplex. The method may further comprise amplifying the polynucleotide duplex to provide an amplified sequence that is a detectable signal that target cellular constituents are in proximity. At least two distinct labeling ligands may comprise oligonucleotide sequences configured to be ligated, so that binding of the labeling ligands to respective target cellular constituents that are in proximity permits the oligonucleotide sequences of the distinct ligands to ligate, forming an amplifiable polynucleotide duplex. The method may further comprise amplifying the polynucleotide duplex to provide an amplified sequence that is a detectable signal that target cellular constituents are in proximity. One of the labeling ligands may comprise a restriction enzyme site between the labeling ligand and the oligonucleotide label, and wherein the method further comprises treating with a restriction enzyme, whereby the UCI from said labeling ligand is transferred to the oligonucleotide label of the labeling ligand in proximity.

[00353] The method may further comprise labeling the cell(s) by fluorescent in situ hybridization.

[00354] The cell(s) may be a member of a cell population, further comprising transforming or transducing the cell population with one or more genomic sequence-perturbation constructs that perturb a genomic sequence in the cells, wherein each distinct genomic sequence-perturbation construct comprises a unique-perturbation-identified (UPI) sequence unique to that construct. The genomic sequence-perturbation construct may comprise a sequence encoding a guide RNA sequence of a CRISPR-Cas targeting system. The method may further comprise multiplex transformation of the population of cells with a plurality of genomic sequence-perturbation constructs. The UPI sequence may be attached to a perturbation-sequence-capture sequence, and the transfer particle may comprise a perturbation-sequence-capture-binding-sequence having

specific binding affinity for the perturbation-sequence-capture sequence attached to the UPI sequence. The UPI sequence may be attached to a universal ligation handle sequence, whereby a USI may be generated by split-pool ligation. The method may further comprise multiplex sequencing of the pooled UCI sequences, USI sequences, and UPI sequences.

[00355] In another aspect, the present invention provides a method of determining interactions between 2 or more cellular constituents, comprising: admixing at least one isolated aggregation of cellular constituents with monomers of a polymerizable gel; polymerizing the gel, to embed the cellular constituents in discrete polymer matrices; incubating the cellular constituents embedded in the polymer matrices with one or more labeling ligands with specific binding affinity for one or more target cellular constituents to produce one or more labeled cellular constituents in the polymer matrices, wherein each of the one or more labeling ligands comprise a bound oligonucleotide label comprising a unique constituent identifier (UCI) sequence and a universal hybridization nucleotide sequence, and wherein the incubation comprises binding conditions under which the labeling ligand will bind to the cellular constituent within the polymer matrix, and the incubation further comprises washing conditions under which unbound labeling ligands will be washed out of the polymer matrix; incubating the polymer matrices with at least one Unique Location Index probe, wherein the probe comprises at least two repeating nucleotide sequences, each repeat comprising a restriction enzyme site, a Unique Location Index (ULI) sequence, and a complementary universal hybridization nucleotide sequence, and wherein the incubation comprises binding conditions under which the universal hybridization sequence will hybridize the complementary universal hybridization sequence; extending the oligonucleotide label hybridized to the probe such that the oligo bound to the affinity ligand incorporates the ULI sequence that is unique to that Unique Location Index probe; digestion with a restriction enzyme specific for the site on the probe, sequencing the oligonucleotide label, whereby detecting the same ULI with two or more UCIs indicates that the cellular constituents were interacting. The ULI sequence may be randomly generated, such that no two ULI sequences are the same. Methods of generating a barcode sequence described herein may be used to generate a ULI. The ULI will be detected with the UCI, such that when multiple cellular constituents are in proximity oligonucleotide labels comprising each UCI and the ULI from a single probe will be generated. Not being bound by a theory, using a plurality of labeling ligands with specificity for a plurality of cellular constituents will allow novel interactions to be

determined. The use of polymer matrices allows a stable platform for washing out the unbound labeling ligands before staining with the ULI probes. The cellular constituent may comprise a protein, RNA transcript, or a DNA molecule. The ULI may be 4-30 nucleotides. The ULI may be 8-20 nucleotides.

[00356] The method may further comprise segregating the discrete polymer matrices comprising the labeled constituents before sequencing. The segregating of the discrete polymer matrices may comprise sorting single discrete matrices into separate reaction vessels.

[00357] The method may further comprise, before sequencing, generating a USI for each discrete polymer matrix by a split pool ligation method, wherein the restriction site on the ULI probe is a universal ligation handle (ULH) sequence configured to produce a DNA overhang capable of hybridization to a complementary overhang on a first index nucleotide sequence, wherein the first index nucleotide sequence comprises an overhang complementary to a final index sequence or optionally a middle index sequence, wherein the middle index sequence comprises overhangs complementary to the first index sequence and to the final index sequence or optionally to another middle index sequence and final index sequence, wherein the final index sequence has a single overhang complementary to the preceding index sequence, and wherein the first, middle, and final index sequences are selected from a plurality of unique sequences comprising compatible DNA overhangs and 10 to 30 base pairs of unique sequence. The split pool ligation method may comprise: splitting the pool of discrete polymer matrices into separate pools of polymer matrices, each containing a unique first index sequence; ligating the first index sequence to the ligation handle; pooling the discrete polymer matrices; optionally, splitting the pool of discrete polymer matrices into separate pools each containing a unique middle index sequence; ligating the middle index sequence to the first index sequence; and pooling the discrete polymer matrices; optionally, repeating step (d) with another middle index sequence; splitting the pool of discrete polymer matrices into pools containing a unique final index sequence; and ligating the final index sequence to the preceding index sequence, whereby each discrete polymer matrix comprises a USI.

[00358] The oligonucleotide label may further comprise a unique molecular identifier (UMI) sequence. The first, middle, or final index sequence may further comprise a unique molecular identifier (UMI) sequence. The method may further comprise pooling the oligonucleotide labels

comprising a USI, ULI and UMI from a plurality of polymer matrices and sequencing the pooled UCI sequences, USI sequences, ULI sequences, and UMI sequences.

[00359] The aggregation of cellular constituents may be a cell that is a member of a cell population, further comprising transforming or transducing the cell population with one or more genomic sequence-perturbation constructs that perturb a genomic sequence in the cells, wherein each distinct genomic sequence-perturbation construct comprises a unique-perturbation-identified (UPI) sequence unique to that construct.

[00360] In another aspect, the present invention provides a method of determining interactions between 2 or more cellular constituents, comprising: fixing and permeablizing at least one cell; incubating the fixed and permeablized cell(s) with one or more labeling ligands with specific binding affinity for one or more target cellular constituents to produce one or more labeled cell(s), wherein each of the one or more labeling ligands comprise a bound oligonucleotide label comprising a unique constituent identifier (UCI) sequence and a universal hybridization nucleotide sequence, and wherein the incubation comprises binding conditions under which the labeling ligand will bind to the cellular constituent within the cell(s), and the incubation further comprises washing conditions under which unbound labeling ligands will be washed from the polymer cell(s); incubating the cell(s) with at least one Unique Location Index probe, wherein the probe comprises at least two repeating nucleotide sequences, each repeat comprising a restriction enzyme site, a Unique Location Index (ULI) sequence, and a complementary universal hybridization nucleotide sequence, and wherein the incubation comprises binding conditions under which the universal hybridization sequence will hybridize to the complementary universal hybridization sequence; extending the oligonucleotide label hybridized to the probe; digesting with a restriction enzyme specific for the site on the probe; and sequencing the oligonucleotide label, whereby detecting the same ULI with two or more UCI's indicates that the cellular constituents were interacting. The cellular constituent may comprise a protein, RNA transcript, or a DNA molecule. The ULI may be 4-30 nucleotides. The ULI may be 8-20 nucleotides.

[00361] The method may further comprise segregating the cell(s) comprising the labeled constituents before sequencing. The segregating of the cell(s) may comprise sorting single discrete matrices into separate reaction vessels. The method may further comprise, before sequencing, generating a USI for each cell by a split pool ligation method, wherein the restriction

site on the ULI probe is a universal ligation handle (ULH) sequence configured to produce a DNA overhang capable of hybridization to a complementary overhang on a first index nucleotide sequence, wherein the first index nucleotide sequence comprises an overhang complementary to a final index sequence or optionally a middle index sequence, wherein the middle index sequence comprises overhangs complementary to the first index sequence and to the final index sequence or optionally to another middle index sequence and final index sequence, wherein the final index sequence has a single overhang complementary to the preceding index sequence, and wherein the first, middle, and final index sequences are selected from a plurality of unique sequences comprising compatible DNA overhangs and 10 to 30 base pairs of unique sequence. The split pool ligation method may comprise: splitting the pool of cells into separate pools of cells, each containing a unique first index sequence; ligating the first index sequence to the ligation handle; pooling the cells; optionally, splitting the pool of cells into separate pools each containing a unique middle index sequence; ligating the middle index sequence to the first index sequence; and pooling the cells; optionally, repeating with another middle index sequence; splitting the pool of cells into pools containing a unique final index sequence; and ligating the final index sequence to the preceding index sequence, whereby each cell comprises a USI.

[00362] The oligonucleotide label may further comprise a unique molecular identifier (UMI) sequence. The first, middle, or final index sequence may further comprise a unique molecular identifier (UMI) sequence. The method may further comprise pooling the oligonucleotide labels comprising a USI, ULI and UMI from a plurality of polymer matrices and sequencing the pooled UCI sequences, USI sequences, ULI sequences, and UMI sequences.

[00363] The cells may be a member of a cell population, further comprising transforming or transducing the cell population with one or more genomic sequence-perturbation constructs that perturb a genomic sequence in the cells, wherein each distinct genomic sequence-perturbation construct comprises a unique-perturbation-identified (UPI) sequence unique to that construct. The perturbation constructs may be any as described herein.

[00364] The oligonucleotide label may comprise a regulatory sequence configured for amplification by T7 polymerase.

[00365] The labeling ligands may comprise oligonucleotide sequences configured to hybridize to a transcript specific region.

[00366] Before sequencing, the method may further comprise: amplification of the oligonucleotide label by PCR; or T7 amplification of the oligonucleotide label followed by subsequent cDNA generation, and optionally amplification by PCR.

[00367] The oligonucleotide label may further comprise at least one spacer sequence. The oligonucleotide label may further comprise a photocleavable linker. The oligonucleotide label may further comprise a restriction enzyme site between the labeling ligand and UCI.

[00368] The oligonucleotide label may comprise one or more iso-dG and/or iso-dC nucleotides. The oligonucleotide labels for hybridization in a proximity assay may comprise one or more iso-dG and/or iso-dC nucleotides. The universal hybridization sequences may comprise one or more iso-dG and/or iso-dC nucleotides. Not being bound by a theory the one or more iso-dG and/or iso-dC nucleotides will increase specificity of hybridization.

[00369] In one embodiment, the oligonucleotide label of any of the methods described herein may comprise one or more iso-dG and/or iso-dC nucleotides. Two complementary sequences may comprise one sequence with iso-dG and the other complementary sequence with iso-dC, whereby the two sequences are capable of hybridizing with each other, but not with sequences containing only dG, dC, dA, and/or dT. The sequence of the oligonucleotide labels for hybridization in a proximity assay may advantageously comprise one or more iso-dG and/or iso-dC nucleotides.

[00370] Any of the methods of the present invention may advantageously be combined for determining any combination of protein detection, RNA detection, open chromatin detection, protein-protein interactions, protein-RNA interactions, or protein-DNA interactions.

[00371] The terms “isolated aggregation of cellular constituents” or “single aggregations of cellular constituents” or “aggregations of cellular constituents” or “aggregations of biologically connected cellular constituents” are used interchangeably and refer to any group of cellular constituents that originate from the same source, that are functionally connected biologically, and that can be isolated individually. Examples may be a cell, an extracellular vesicle, an organelle, or an organized subcomponent thereof. Specific examples may be a nucleus or a mitochondria.

[00372] The term “cellular constituent” refers to any cellular molecule, including but not limited to a protein, nucleic acid, RNA molecule, DNA molecule, or carbohydrate.

[00373] The term “unique molecular identifiers” (UMI) refers to a sequencing linker used in a method that uses molecular tags to detect and quantify unique amplified products. A UMI is used

to distinguish effects through a single clone from multiple clones. In preferred embodiments, the amplification is by PCR. A sequencer linker with a random sequence of between 4 and 20 basepairs and an index sequence is added to the 5' end of the template, which is amplified and sequenced. Sequencing allows for high resolution reads, enabling accurate detection of true variants. As used herein, a "true variant" will be present in every amplified product originating from the original clone as identified by aligning all products with a UMI. Each clone amplified will have a different random UMI that will indicate that the amplified product originated from that clone. Background caused by the fidelity of the amplification process can be eliminated because true variants will be present in all amplified products and background representing random error will only be present in single amplification products (See e.g., Islam S. et al., 2014. Nature Methods No:11, 163-166). Not being bound by a theory, the UMI and UCI's are designed such that assignment to the original can take place despite up to 4-7 errors during amplification or sequencing.

[00374] The term "unique constituent identifier" (UCI) refers to any unique nucleotide sequence linked to a labeling ligand, such that the presence of the sequence indicates the presence of the cellular constituent that the labeling ligand specifically binds. In an exemplary embodiment, the UCI is linked to an antibody for a specific cellular constituent. If the cellular constituent is present in a sample, the antibody will bind and the UCI can be detected. If the cellular constituent is not present in a sample, the antibody will not bind and the UCI will not be detected above background. In another exemplary embodiment, the labeling ligand is an oligonucleotide probe and the cellular constituent is an RNA transcript molecule complementary to the sequence of the oligonucleotide probe. The sequence of the oligonucleotide probe may be the UCI or may comprise an additional UCI sequence to identify the RNA transcript.

[00375] The term "unique source identifier" (USI) refers to a unique nucleotide sequence that is associated with the nucleic acids from a single cell or single isolated aggregation of cellular constituents (source), such that upon sequencing a pool of nucleic acid sequences from more than one cell or isolated aggregation of cellular constituents, the presence of a USI in the sequenced product indicates that a product originated from a single source. USI may be used interchangeably with the term "barcode."

[00376] The term "unique-amplification-identifier" (UAI) refers to a nucleotide sequence that is only formed only when two or more nucleotide sequences are in close proximity to each other

such that they can be ligated. The UAI can be generated using methods described for the proximity ligation assay (PLA) or proximity extension assay (PEA) (Fredriksson S, et al. (2002) Protein detection using proximity-dependent DNA ligation assays. *Nature biotechnology* 20: 473–477; Gullberg M, et al. (2004) Cytokine detection by antibody-based proximity ligation. *Proceedings of the National Academy of Sciences of the United States of America* 101: 8420–8424; and Lundberg M, et al. (2011) Homogeneous antibody-based proximity extension assays provide sensitive and specific detection of low-abundant proteins in human blood. *Nucleic acids research* 39(15): e102). PEA is based on pairs of antibodies that are linked to oligonucleotides having slight affinity to one another (PEA probes). Upon target binding the probes are brought in proximity, and the two oligonucleotides are extended by a DNA polymerase forming the UAI that now acts as a unique surrogate marker for the specific antigen.

[00377] The terms “sticky end,” “overhang” and “DNA overhang” refer to a double stranded DNA having either a 3’ or 5’ single stranded DNA overhang capable of hybridization to another complementary sticky end or DNA overhang.

[00378] The term “hydrogel” refers to any network of polymer chains that are hydrophilic, and sometimes found as a colloidal gel, in which water is the dispersion medium. Hydrogels are highly absorbent (they can contain over 90% water) natural or synthetic polymeric networks. Hydrogels also possess a degree of flexibility very similar to natural tissue, due to their significant water content. Hydrogel may include polyvinyl alcohol, sodium polyacrylate, acrylate polymers, copolymers with an abundance of hydrophilic groups, agarose, methylcellulose, hyaluronan, and other naturally derived polymers.

[00379] The term “tagmentation” refers to a step in the Assay for Transposase Accessible Chromatin using sequencing (ATAC-seq) as described. (See, Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y., Greenleaf, W. J., Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nature methods* 2013; 10 (12): 1213-1218). Specifically, a hyperactive Tn5 transposase loaded in vitro with adapters for high-throughput DNA sequencing, can simultaneously fragment and tag a genome with sequencing adapters. In one embodiment the adapters are compatible with the methods described herein.

[00380] The present invention may also include barcoding. Barcoding may be performed based on any of the compositions or methods disclosed in patent publication WO 2014047561

A1, Compositions and methods for labeling of agents, incorporated herein in its entirety. In one embodiment each labeling ligand has a barcode (UCI). In one embodiment, a sgRNA has a barcode. In one embodiment the UCI is captured on a bead that includes a barcode sequence (USI). Not being bound by a theory, amplified sequences from single cells or isolated aggregations of cellular constituents can be sequenced together and resolved based on the barcode associated with each USI. Not being bound by a theory, the presence of a labeling ligand can be determined by sequencing of the UCI.

[00381] Barcoding may be performed based on any of the compositions or methods disclosed in patent publication WO 2014047561 A1, Compositions and methods for labeling of agents, incorporated herein in its entirety. In certain embodiments barcoding uses an error correcting scheme (T. K. Moon, *Error Correction Coding: Mathematical Methods and Algorithms* (Wiley, New York, ed. 1, 2005)). Not being bound by a theory, amplified sequences from single cells can be sequenced together and resolved based on the barcode associated with each cell.

[00382] The term “barcode” as used herein, refers to any unique, non-naturally occurring, nucleic acid sequence that may be used to identify the originating source of a nucleic acid fragment. Such barcodes may be sequences including but not limited to about 20 base pair sequences. Although it is not necessary to understand the mechanism of an invention, it is believed that the barcode sequence provides a high-quality individual read of a barcode associated with a viral vector, labeling ligand, shRNA, sgRNA or cDNA such that multiple species can be sequenced together.

[00383] DNA barcoding is also a taxonomic method that uses a short genetic marker in an organism's DNA to identify it as belonging to a particular species. It differs from molecular phylogeny in that the main goal is not to determine classification but to identify an unknown sample in terms of a known classification. Kress et al., “Use of DNA barcodes to identify flowering plants” *Proc. Natl. Acad. Sci. U.S.A.* 102(23):8369-8374 (2005). Barcodes are sometimes used in an effort to identify unknown species or assess whether species should be combined or separated. Koch H., “Combining morphology and DNA barcoding resolves the taxonomy of Western Malagasy *Liotrigona* Moure, 1961” *African Invertebrates* 51(2): 413-421 (2010); and Seberg et al., “How many loci does it take to DNA barcode a crocus?” *PLoS One* 4(2):e4598 (2009). Barcoding has been used, for example, for identifying plant leaves even when flowers or fruit are not available, identifying the diet of an animal based on stomach contents or

feces, and/or identifying products in commerce (for example, herbal supplements or wood). Soininen et al., “Analysing diet of small herbivores: the efficiency of DNA barcoding coupled with high-throughput pyrosequencing for deciphering the composition of complex plant mixtures” *Frontiers in Zoology* 6:16 (2009).

[00384] It has been suggested that a desirable locus for DNA barcoding should be standardized so that large databases of sequences for that locus can be developed. Most of the taxa of interest have loci that are sequencable without species-specific PCR primers. CBOL Plant Working Group, “A DNA barcode for land plants” *PNAS* 106(31):12794-12797 (2009). Further, these putative barcode loci are believed short enough to be easily sequenced with current technology. Kress et al., “DNA barcodes: Genes, genomics, and bioinformatics” *PNAS* 105(8):2761-2762 (2008). Consequently, these loci would provide a large variation between species in combination with a relatively small amount of variation within a species. Lahaye et al., “DNA barcoding the floras of biodiversity hotspots” *Proc Natl Acad Sci USA* 105(8):2923-2928 (2008).

[00385] DNA barcoding is based on a relatively simple concept. For example, most eukaryote cells contain mitochondria, and mitochondrial DNA (mtDNA) has a relatively fast mutation rate, which results in significant variation in mtDNA sequences between species and, in principle, a comparatively small variance within species. A 648-bp region of the mitochondrial cytochrome c oxidase subunit 1 (CO1) gene was proposed as a potential ‘barcode’. As of 2009, databases of CO1 sequences included at least 620,000 specimens from over 58,000 species of animals, larger than databases available for any other gene. Ausubel, J., “A botanical microscope” *Proceedings of the National Academy of Sciences* 106(31):12569 (2009).

[00386] Software for DNA barcoding requires integration of a field information management system (FIMS), laboratory information management system (LIMS), sequence analysis tools, workflow tracking to connect field data and laboratory data, database submission tools and pipeline automation for scaling up to eco-system scale projects. Geneious Pro can be used for the sequence analysis components, and the two plugins made freely available through the Moorea Biocode Project, the Biocode LIMS and Genbank Submission plugins handle integration with the FIMS, the LIMS, workflow tracking and database submission.

[00387] Additionally, other barcoding designs and tools have been described (see e.g., Birrell et al., (2001) *Proc. Natl Acad. Sci. USA* 98, 12608–12613; Giaever, et al., (2002) *Nature* 418,

387–391; Winzeler et al., (1999) *Science* 285, 901–906; and Xu et al., (2009) *Proc Natl Acad Sci U S A*. Feb 17;106(7):2289-94).

[00388] In certain embodiments, sequencing is performed using unique molecular identifiers (UMI). The term “unique molecular identifiers” (UMI) refers to a sequencing linker used in a method that uses molecular tags to detect and quantify unique amplified products. A UMI is used to distinguish effects through a single clone from multiple clones. In preferred embodiments, the amplification is by PCR. A sequencer linker with a random sequence of between 4 and 20 base pairs is added to the 5' end of the template, which is amplified and sequenced. Sequencing allows for high resolution reads, enabling accurate detection of true variants. As used herein, a “true variant” will be present in every amplified product originating from the original clone as identified by aligning all products with a UMI. Each clone amplified will have a different random UMI that will indicate that the amplified product originated from that clone. Background caused by the fidelity of the amplification process can be eliminated because true variants will be present in all amplified products and background representing random error will only be present in single amplification products (See e.g., Islam S. et al., 2014. *Nature Methods* No:11, 163-166). Not being bound by a theory, the UMI's are designed such that assignment to the original can take place despite up to 4-7 errors during amplification or sequencing.

[00389] Unique molecular identifiers are a subtype of nucleic acid barcode that can be used, for example, to normalize samples for variable amplification efficiency. For example, in various embodiments, featuring a solid or semisolid support (for example a hydrogel bead), to which nucleic acid barcodes (for example a plurality of barcode sharing the same sequence) are attached, each of the barcodes may be further coupled to a unique molecular identifier, such that every barcode on the particular solid or semisolid support receives a distinct unique molecule identifier. A unique molecular identifier can then be, for example, transferred to a target molecule with the associated barcode, such that the target molecule receives not only a nucleic acid barcode, but also an identifier unique among the identifiers originating from that solid or semisolid support.

[00390] In certain embodiments, multiple displacement amplification (MDA) is used. Multiple displacement amplification (MDA, is a non-PCR-based isothermal method based on the annealing of random hexamers to denatured DNA, followed by strand-displacement synthesis at constant temperature (Blanco et al. *J. Biol. Chem.* 1989, 264, 8935-8940). It has been applied to

samples with small quantities of genomic DNA, leading to the synthesis of high molecular weight DNA with limited sequence representation bias (Lizardi et al. *Nature Genetics* 1998, 19, 225-232; Dean et al., *Proc. Natl. Acad. Sci. U. S. A.* 2002, 99, 5261-5266). As DNA is synthesized by strand displacement, a gradually increasing number of priming events occur, forming a network of hyper- branched DNA structures. The reaction can be catalyzed by enzymes such as the Phi29 DNA polymerase or the large fragment of the Bst DNA polymerase. The Phi29 DNA polymerase possesses a proofreading activity resulting in error rates 100 times lower than Taq polymerase (Lasken et al. *Trends Biotech.* 2003, 21, 531-535).

[00391] The invention provides a method for preparing uniquely barcoded particles. Unique barcode sequences may be generated by a split pool method. The split pool method may include sticky end ligation. Sticky end ligation may include a sticky end ligation handle and separate indexes containing unique sequences capable of hybridizing to a sticky end (see examples). The sticky end may comprise a ssDNA overhang. The over-hang may be 2, 3, 4, 5, 6, 7, 8, preferably 4 bases. The overhang may be generated by a restriction enzyme. Each index may contain a plurality of unique sequences. Each index may contain 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 200, preferably 192 sequences. In one embodiment there are 2, 3, 4, preferably 3 indexes. A unique barcode sequence is generated by ligation of the first index to the ligation handle, splitting and pooling of the ligated samples, and then addition of the next index also containing sticky ends. The last index preferably has a sticky end for ligation to the previous index. The last index may advantageously include a primer sequence for priming of PCR. Methods of split pooling have been described. In one embodiment the ligation handle is digested with a restriction enzyme to produce a four base overhang. In another embodiment, a ligation primer is hybridized to the ligation handle to generate an at least 4 base overhang that is complementary to an index in the split pool method.

[00392] In one exemplary embodiment, the hydrogel particles or polymer matrices are split into pools, each pool containing a unique index A and each ligation handle is ligated to a sequence in index A. All particles are then pooled and re-split into new pools containing a unique index B. After ligation, all of the particles are pooled again and re-split into new pools containing a unique index C. If each index has 100 unique sequences and for each cycle the particles are split into 100 pools each containing a unique sequence, then after 3 cycles of split

and pool ligation, the barcode on any given particle possess the same one of $100^3 = 1,000,000$ possible barcodes, but different particles have different sequences.

[00393] In another embodiment, single cell or single isolated aggregation of cellular constituent analysis is performed by digital polymerase chain reactions (PCR), e.g., Fluidigm C. Digital polymerase chain reaction (digital PCR, DigitalPCR, dPCR, or dePCR) is a refinement of conventional polymerase chain reaction methods that can be used to directly quantify and clonally amplify nucleic acids including DNA, cDNA or RNA. The key difference between dPCR and traditional PCR lies in that PCR carries out one reaction per single sample and dPCR carries out a single reaction within samples separated into a large number of partitions wherein the reactions are carried out in each partition individually. A sample is partitioned so that individual nucleic acid molecules within the sample are localized and concentrated within many separate regions. The capture or isolation of individual nucleic acid molecules may be effected in micro well plates, capillaries, the dispersed phase of an emulsion, and arrays of miniaturized chambers, as well as on nucleic acid binding surfaces.

[00394] In a preferred embodiment, single cell or single aggregation of cellular constituent analysis is performed using microfluidics. Microfluidics involves micro-scale devices that handle small volumes of fluids. Because microfluidics may accurately and reproducibly control and dispense small fluid volumes, in particular volumes less than 1 μ l, application of microfluidics provides significant cost-savings. The use of microfluidics technology reduces cycle times, shortens time-to-results, and increases throughput. Furthermore, incorporation of microfluidics technology enhances system integration and automation. Microfluidic reactions are generally conducted in microdroplets. The ability to conduct reactions in microdroplets depends on being able to merge different sample fluids and different microdroplets. See, e.g., US Patent Publication No. 20120219947 and PCT publication No. WO2014085802 A1.

[00395] Droplet microfluidics offers significant advantages for performing high-throughput screens and sensitive assays. Droplets allow sample volumes to be significantly reduced, leading to concomitant reductions in cost. Manipulation and measurement at kilohertz speeds enable up to 10^8 samples to be screened in a single day. Compartmentalization in droplets increases assay sensitivity by increasing the effective concentration of rare species and decreasing the time required to reach detection thresholds. Droplet microfluidics combines these powerful features to

enable currently inaccessible high-throughput screening applications, including single-cell and single-molecule assays. See, e.g., Guo et al., *Lab Chip*, 2012,12, 2146-2155.

[00396] The manipulation of fluids to form fluid streams of desired configuration, discontinuous fluid streams, droplets, particles, dispersions, etc., for purposes of fluid delivery, product manufacture, analysis, and the like, is a relatively well-studied art. Microfluidic systems have been described in a variety of contexts, typically in the context of miniaturized laboratory (e.g., clinical) analysis. Other uses have been described as well. For example, WO 2001/89788; WO 2006/040551 ; U.S. Patent Application Publication No. 2009/0005254; WO 2006/040554; U.S. Patent Application Publication No. 2007/0184489; WO 2004/002627; U.S. Patent No. 7,708,949; WO 2008/063227; U.S. Patent Application Publication No. 2008/0003142; WO 2004/091763; U.S. Patent Application Publication No. 2006/0163385; WO 2005/021151 ; U.S. Patent Application Publication No. 2007/0003442; WO 2006/096571 ; U.S. Patent Application Publication No. 2009/0131543; WO 2007/089541; U.S. Patent Application Publication No. 2007/0195127; WO 2007/081385; U.S. Patent Application Publication No. 2010/0137163; WO 2007/133710; U.S. Patent Application Publication No. 2008/0014589; U.S. Patent Application Publication No. 2014/0256595; and WO 2011/079176. In a preferred embodiment single cell analysis is performed in droplets using methods according to WO 2014085802. Each of these patents and publications is herein incorporated by reference in their entireties for all purposes.

[00397] Single cells or isolated aggregations of cellular constituents may be sorted into separate vessels by dilution of the sample and physical movement, such as pipetting. A machine can control the pipetting and separation. The machine may be a computer controlled robot.

[00398] Microfluidics may also be used to separate the single cells and/or isolated aggregations of cellular constituents. Single cells and/or isolated aggregations of cellular constituents can be separated using microfluidic devices. Microfluidics involves micro-scale devices that handle small volumes of fluids. Because microfluidics may accurately and reproducibly control and dispense small fluid volumes, in particular volumes less than 1 μ l, application of microfluidics provides significant cost-savings. The use of microfluidics technology reduces cycle times, shortens time-to-results, and increases throughput. The small volume of microfluidics technology improves amplification and construction of DNA libraries made from single cells and single isolated aggregations of cellular constituents. Furthermore, incorporation of microfluidics technology enhances system integration and automation.

[00399] Single cells and/or single isolated aggregations of cellular constituents of the present invention may be divided into single droplets using a microfluidic device. The single cells and/or single isolated aggregations of cellular constituents in such droplets may be further labeled with a barcode. In this regard reference is made to Macosko et al., 2015, “Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets” Cell 161, 1202–1214 and Klein et al., 2015, “Droplet Barcoding for Single-Cell Transcriptomics Applied to Embryonic Stem Cells” Cell 161, 1187–1201, all the contents and disclosure of each of which are herein incorporated by reference in their entirety. Not being bound by a theory, the volume size of an aliquot within a droplet may be as small as 1 fL.

[00400] Single cells and/or single aggregations of cellular constituents may be diluted into a physical multi-well plate or a plate free environment. The multi-well assay modules (e.g., plates) may have any number of wells and/or chambers of any size or shape, arranged in any pattern or configuration, and be composed of a variety of different materials. Preferred embodiments of the invention are multi-well assay plates that use industry standard multi-well plate formats for the number, size, shape and configuration of the plate and wells. Examples of standard formats include 96-, 384-, 1536- and 9600-well plates, with the wells configured in two-dimensional arrays. Other formats include single well, two well, six well and twenty-four well and 6144 well plates. Plate free environments of the present invention utilize a single polymerizable gel containing compartmentalized cells and/or isolated aggregations of cellular constituents. In one embodiment, extraction of single cells and/or single isolated aggregations of cellular constituents may be by a mechanical punch. Single cells and/or single isolated aggregations of cellular constituents may be visualized in the gel before a punch.

[00401] In one embodiment, a DNA tag including a protein specific barcode (UCI) is conjugated to detection biomolecules or labeling ligands with high target affinity and low unspecific binding, such as antibodies (Janssen et al., 2013) or nanobodies (Pardon et al., 2014; Theile et al., 2013) or aptamers (Janssen et al., 2013).

[00402] In one embodiment, to ensure proper staining of intracellular and cell surface proteins with, for instance, DNA-tagged antibodies, single cells are embedded in hydrogel droplets. Not being bound by a theory, the hydrogel mesh provides a physical framework, chemically incorporates biomolecules and is permeable to macromolecules such as antibodies (Chung et al., 2013). In one embodiment, to further improve permeability and staining efficiency, lipids are

cleared ([Chung et al., 2013](#)). Not being bound by a theory, the clearance of the lipids and the porosity of the hydrogel allow for more efficient washing and removal of unspecific antibodies. This higher accuracy of measurement is important for the high multiplex measurements and computational inference of regulatory mechanisms.

[00403] In one embodiment, cells embedded in a hydrogel mesh can be stained with the DNA-tagged antibodies and washed in bulk before isolating the single cells. Once isolated, a cell specific oligonucleotide barcode (USI) can be introduced before subsequent DNA amplification and library preparation steps. Isolating single cells into individual reaction chambers to perform PCR amplification or a proximity ligation/extension assay ([Assarsson et al., 2014](#)) can be achieved at modest throughput either by FACS sorting into multi-well plates or microfluidic capture using the Fluidigm C1 ([Shalek et al., 2014](#)).

[00404] In one embodiment, for more high throughput processing, a microfluidic chip can be used to capture the hydrogel embedded cells or cellular constituents in nanoliter-sized aqueous droplets (Macosko et al., 2015, “Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets” *Cell* 161, 1202–1214). In one embodiment, the hydrogel embedded cells or cellular constituents are poisson loaded into microwells ([Fan et al., 2015](#)). The aqueous droplets or microwells may be simultaneously loaded with barcoded beads, each of which has oligonucleotides including; a “cell barcode” that is the same across all the primers on the surface of any one bead, but different from the cell barcodes on all other beads; a Unique Molecular Identifier (UMI), different on each primer, that enables sequence reads derived from the same original DNA tag (amplification and PCR duplicates) to be identified computationally (Kivioja et al., 2012); and a capture sequence to bind the oligos (either amplified PCR products or original DNA tags released by proteinase K treatment, or enzymatic/photonic oligo cleavage). Once the beads are loaded, they can be pooled for amplification and library preparation, and sequencing. These beads can take multiple forms, the preferred drop-seq beads are polystyrene, oligo functionalized beads, but alternative beads are possible, such as soft beads (polymer gel based beads), that allow for one on one pairing with cells, as to avoid the poisson loading needed in the described drop-seq scheme. This reduces the amount of cells one needs, and makes it possible to analyze rare cell types or clinical samples only available in low amounts of cells.

[00405] In one embodiment, the present invention provides for the simultaneous detection of proteins and nucleic acids. Nucleic acids can be reverse cross-linked after separation of discrete polymer matrices into separate wells or droplets. The contents of individual wells or droplets may then be sequenced. In one embodiment, crosslinking is reversed by incubating the cross-linked sample in high salt (approximately 200 mM NaCl) at 65°C for at least 4h.

[00406] In one embodiment, Drop-Seq (Macosko et al., 2015) is used to analyze RNA or DNA in single cells in parallel to the detection of cellular constituents. Drop-Seq is a reverse emulsion, early barcoding method for analyzing 10^4 - 10^6 cells/experiment at very low cost (\$0.06/cell). The Drop-seq method may be used to encapsulate discrete hydrogel matrices in a droplet. The RNA and/or DNA can be reverse cross-linked and the oligonucleotide labels can be removed from the labelling ligand. Capture of RNA, DNA, and oligonucleotide labels on barcoded beads, library preparation, and sequencing is performed as described previously.

[00407] In one embodiment, the detection of proteins or post translational modifications (PTM) is determined by sequencing based readouts. In some embodiments, Immuno-Seq is used when antibodies can be washed out (Niemeyer, C. M., et al., Nat Protoc. 2, 1918-1930 (2007)) and proximity extension assays (PEA) is used when antibodies cannot be washed away (Hammond, M., et al. PLoS One. 7, e40405, (2012); and Stahlberg, A. , et al. Clin Chem. 58, 1682-1691 (2012)). These methods use DNA-sequence based encoding , and are compatible with other genomic readouts (e.g., sgRNA barcodes).

[00408] In another embodiment, the detection of proteins embedded in a hydrogel matrix is determined by FACS. Not being bound by a theory, the encapsulation of cellular constituents in a hydrogel matrix and removing lipids provides for improved binding of antibodies to intracellular targets as compared to regular fixation and permeabilization protocols for FACS alone.

[00409] In one embodiment, PEA methods are used for profiling protein-protein or protein-nucleic acid interactions by, respectively, using antibodies against two protein targets (Leuchowius, K. J., et al. Cytometry A. 75, 833-839 (2009)). or replacing one antibody with an oligonucleotide complementary to a sequence of interest (Gustafsdottir, S. M., et al. Proceedings of the National Academy of Sciences of the United States of America. 104, 3067-3072, (2007)).

[00410] In another aspect, the present invention provides screening methods to determine the effect on protein, post translational modifications and cellular constituents of single cells or

isolated aggregations of cellular constituents in response to the perturbation of genes or cellular circuits. Perturbation may be knocking down a gene, increasing expression of a gene, mutating a gene, mutating a regulatory sequence, or deleting non-protein-coding DNA.

[00411] In one embodiment, CRISPR/Cas9 may be used to perturb protein-coding genes or non-protein-coding DNA. CRISPR/Cas9 may be used to knockout protein-coding genes by frameshifts, point mutations, inserts, or deletions. An extensive toolbox may be used for efficient and specific CRISPR/Cas9 mediated knockout as described herein, including a double-nicking CRISPR to efficiently modify both alleles of a target gene or multiple target loci and a smaller Cas9 protein for delivery on smaller vectors (Ran, F. A. , et al., *In vivo genome editing using Staphylococcus aureus Cas9*. *Nature*. 520, 186-191 (2015)). A genome-wide sgRNA mouse library (10 sgRNAs/gene) may also be used in a mouse that expresses a Cas9 protein. The cells of the mouse can then be analyzed using the methods of the present invention.

[00412] In one embodiment, a CRISPR system may be used to activate gene transcription. A nuclease-dead RNA-guided DNA binding domain, dCas9, tethered to transcriptional repressor domains that promote epigenetic silencing (e.g., KRAB) may be used for "CRISPRi" that represses transcription. To use dCas9 as an activator (CRISPRa), a guide RNA is engineered to carry RNA binding motifs (e.g., MS2) that recruit effector domains fused to RNA-motif binding proteins, increasing transcription. A key dendritic cell molecule, p65, may be used as a signal amplifier, but is not required.

[00413] In one embodiment, perturbation is by deletion of regulatory elements. Non-coding elements may be targeted by using pairs of guide RNAs to delete regions of a defined size, and by tiling deletions covering sets of regions in pools.

[00414] In one embodiment, perturbation of genes is by RNAi. The RNAi may be shRNA's targeting genes. The shRNA's may be delivered by any methods known in the art. In one embodiment the shRNA's may be delivered by a viral vector. The viral vector may be a lentivirus.

[00415] In one embodiment, a CRISPR based pooled screen is used. Perturbation may rely on sgRNA expression cassettes that are stably integrated into the genome. The expressed sgRNA may serve as a molecular barcode, reporting the loss of function of the target in a cell. Alternatively, optimized separate barcodes may be co-expressed with the sgRNA, should sgRNAs not be ideal as barcodes. Transduction of cells at a higher multiplicity of infection

(MOI) or delivering vectors by transfection at a higher MOI would result in any given cell receiving multiple sgRNA's and allow combinatorial perturbations. In one embodiment, 2, or 3, or 4, or 5, or up to 10 genes, preferably 5-7 genes are perturbed in a single cell.

[00416] In one embodiment, recombinant Cas9 protein and sgRNA is delivered simultaneously to cells with nanowires or the recently developed 'Cellsqueeze' (Sharei, A., et al. Proceedings of the National Academy of Sciences of the United States of America. 110, 2082-2087, (2013)). Applicants have shown that nanowires can deliver functional proteins, RNA and small molecules alone and in combinations into the cell's cytoplasm, and do not cause toxicity or inappropriate activation and allow the cells to respond normally to signals (Shalek, A. K., et al. Nano Lett. 12, 6498-6504, (2012); Yosef, N., et al. Nature. 496, 461-468, (2013); and Shalek, A. K., et al. Proceedings of the National Academy of Sciences of the United States of America. 107, 1870-1875, (2010)).

[00417] In one embodiment, hybrid measurements or alternative readouts are measured. The alternative readouts may either be stand alone, or hybrid measurements. One alternative readout may be epigenetic measurements. Not being bound by a theory, when biomolecules with functional groups are formaldehyde fixed and bound to the polymer mesh, and membrane and nuclear lipids are cleared, chromosomal DNA is preserved and is accessible for further interrogation. Epigenetic assays that have been applied to single cells may be combined with a perturbation and protein level readout. Not being bound by a theory, the new layers of information aid in understanding of the regulatory mechanisms underpinning cellular behavior. Histone modifications have been measured at specific gene loci at the single cell level (Gomez et al., 2013). This publication uses ISH-PLA (in situ hybridization (ISH), proximity ligation assay (PLA)). They use a biotin modified ISH probe, by binding with streptavidin and an oligo bound anti-streptavidin antibody. As antibodies against multiple histone modifications are readily available, the PLA scheme is applicable to the present invention. Not being bound by a theory, a histone code based on the combination of a plurality of histone modifications determines gene expression at a given locus. Many histone modifications at many genetic loci can be determined simultaneously by replacing the biotin-streptavidin construct by an ISH probe conjugated to a linker (peptide, DNA or nanoparticles,...), followed by another DNA barcode reporting on the genetic locus, and including a binding sequence to the oligo conjugated to the histone modification antibody.

[00418] In one embodiment, chromatin accessibility is determined using a single cell ATAC-seq assay. ATAC-seq offers genome-wide chromatin accessibility of regulatory elements, transcription factor binding and nucleosome positioning.

[00419] In one embodiment, DNA methylation analysis is determined. Cytosine methylation analysis has been analyzed at the single cell level (Kantlehner et al., 2011), as has adenine methylation (Lorthongpanich et al., 2013).

[00420] In one embodiment, the spatial organization of chromosomes is determined. The spatial organization of chromosomes has been found to have fundamental effects on gene expression and cellular function. Single cell measurements (Hi-C) have revealed extensive cell-to-cell heterogeneity in chromosome structure (Nagano et al., 2013). This method can be incorporated into the present invention.

[00421] In one embodiment, protein-protein interactions are measured. In addition to assessing presence and abundance of individual proteins, assays such as Proximity Extension Assay (PEA) allow for assaying the proximity of two proteins. In particular, the present invention allows for probing protein-protein interactions by designing pairs of antibodies for the interacting proteins of interest, such that the oligos conjugated to these antibodies have a binding region, which only bind when the two proteins are in near proximity, and therefore only PCR amplify in this case.

[00422] In one embodiment, protein-DNA interaction measurements are determined. Similar to the modified ISH-PLA described herein, instead of probing histone modifications, one could probe protein (transcription factor) proximity to many specific genetic loci, in a multiplex fashion.

[00423] In one embodiment, fluorescent in situ hybridization methods are used in the present invention. The present invention allows a combined approach where cells can be fluorescently labeled by methods known in the art, and cells of interest can be selected for downstream profiling of cellular constituents. In addition, the assays of the present invention can be combined with in situ hybridization methods such as RNA and DNA FISH.

[00424] In another embodiment, the gelled and cleared cells offer a platform in which any biological agent that is able to be detected by a high affinity and specific counterpart or ligand that can directly or indirectly be conjugated to a DNA molecule could be detected and quantified using the methods of the present invention.

[00425] Releasing the oligo's to be sequenced from their antibody can take a multitude of forms; i.e. in one embodiment, oligo's could be released from their antibodies by digesting all proteins (for instance proteinase K), alternatively, photocleavable linkers could be used, or restriction sites could be included in the oligo sequence to allow for enzymatic restriction and release. In another embodiment, the oligo can stay bound to the antibody, and in situ amplified (i.e. either by PCR, rolling circle amplification or T7 polymerase amplification) and the products of this reaction could be captured and sequenced.

[00426] Similarly, capturing the released oligo's could take a number of forms: in a drop based approach, beads can be loaded with capture oligo's as described herein. Microwells could either be loaded with beads, or their surface could be functionalized with capture oligos from which further amplification could take place. Alternatively, in the scenario where drops are sorted into multiwell plates, or microfluidic reaction chambers such as the Fluidigm C1 system, oligos can be amplified linearly or exponentially, and cellular barcodes and library adapters can be added on during these amplification steps.

[00427] Many different assays have been developed for oligo-barcode based detection of proteins (Janssen et al., 2013) and may be used in the present invention.

[00428] In one embodiment, cells are fixed and monomer infused before capturing them in a droplet. Alternatively, cells or aggregations of constituents are co-flowed with a lysis/monomer solution into a larger diameter drop. In this embodiment, biomolecules from a single cell or isolated aggregation of constituents are spread over a larger volume, which with similar polymer density could increase accessibility for staining.

[00429] The present invention also provides for cell handling before hydrogel polymerization. In one embodiment, cells are fixed and infused with polymer monomers in bulk. Cells may then be segregated and polymerization initiated. Segregation can be by any means described herein. In preferred embodiments, segregation is performed by making single cell drops.

[00430] In another embodiment, biochemical, thermal, or optical treatment on chip of individual cells in reverse emulsion droplets is performed. In this embodiment, polymer monomers may be spiked in microfluidically and optionally fixation reagents. Polymerization of the monomers may then be performed. This allows biochemical, thermal, or optical treatments at the single-cell level. Examples include, but are not limited to: lysis, DNA/RNA

fragmentation/tagmentation, dosing with drugs, enzymatic reactions, or any perturbation of the sample before fixation and/or anchoring biomolecules to the polymer mesh upon polymerization.

[00431] In one embodiment, the oligonucleotide label may comprise Iso-deoxyguanosine (iso-dG) and 5-methyl iso-dC (iso-dC). Iso-deoxyguanosine forms a Watson-Crick base pair with 5-methyl iso-dC, but has a different type of hydrogen bonding pattern than those observed for the natural base pairs A:T and C:G. Substitution of a iso-dG: 5-Me-iso-dC base pair for a C:G pair increases the T_m of the resulting duplex by ~2 deg C per base pair substitution (Switzer, C., et al., Enzymatic incorporation of a new base pair into DNA and RNA. *J. Am Chem. Soc.* (1989), **111**: 8322-8323; and Horn, T., et al., Hybridization properties of the 5-methylisocytidine/isoguanosine base pair in synthetic oligodeoxynucleotides. *Tetrahedron Lett.* (1995), **36**: 2033-2036). Furthermore, since iso-dG does not pair with dC, iso-dG: 5-Me-iso-dC can function as a stable unnatural base pair that can be used to expand the genetic code. The combination of iso-dG's high selectivity for 5-Me-iso-dC, and the resulting base pair's high thermodynamic stability, make this modified base pair particular attractive in embodiments of the present invention.

[00432] In one embodiment, iso-dG:5-Me-iso-dC base pairing is used for molecular recognition. The 5-Me-iso-dC:iso-dG base pair may be incorporated into hybridization assays to enhance probe-target specificity and reduce spurious hybridization to non-target sequences. For example, Collins and co-workers significantly improved the sensitivity of a branched DNA quantitative hybridization assay for detecting the HIV POL sequence by incorporating ~30% 5-Me-iso-dC and iso-dG into the pre-amplifier, branched DNA (bDNA) amplifier and alkaline phosphate probe sequences used in the assay (Collins, M.L., et al. A branched DNA signal amplification assay for quantification of nucleic acid targets below 100 molecules/ml. *Nucleic Acids Res.* (1997), **25**: 2979-2984). Use of this strategy resulted in a significant reduction in non-specific hybridization of the above three sequence types to non-target nucleic acid sequences, and thus less amplification of background. The limits of detection of the assay were improved 10-fold, from < 500 HIV molecules/mL to < 50 molecules/mL. In preferred embodiments, the present invention utilizes the 5-Me-iso-dC:iso-dG base pair to ensure the correct sequences base pair during hybridization of ligation handle primers and during hybridization of two oligonucleotide labels in proximity assays.

[00433] In another embodiment, iso-dG:5-Me-iso-dC base pairing is used for qPCR and artificially expanded genetic systems. A number of research groups have been working on optimizing PCR amplification on templates containing 5-Me-iso-dC. Such optimization is necessary to enable the full development of artificially expanded genetic systems utilizing an expanded genetic code, thereby allowing for the site-specific incorporation of novel functional components (such as unnatural amino acids) into proteins. In 2004, Johnson and co-workers observed that, by using the Klenow fragment of Taq polymerase (KF-Taq) in PCR, the fidelity of the 5-Me-iso-dC:iso-dG base pair was about 96% per amplification cycle (Johnson, S.C., et al., A third base pair for the polymerase chain reaction: inserting isoC and isoG. *Nucleic Acids Resl.* (2004), 32: 1937-1941). The limit in fidelity is chiefly due to the ability of iso-dG's 1,2 tautomer to mis-pair with dT. More recently, Sismour and Benner solved this problem by using 2-thio-dT (dT*) in place of dT. dT* pairs with dA, but not with iso-dG (Sismour, A.M.; Benner, S.A. The use of thymidine analogs to improve the replication of an extra DNA base pair: a synthetic biological system. *Nucleic Acids Resl.* (2005), 33: 5640-5646). Using this artificial base pair system (5-Me-iso-dC:iso-dG, dA:dT*, dC:dG) with KF-Taq, the fidelity in PCR was increased to about 98% per amplification cycle.

[00434] The practice of the present invention employs, unless otherwise indicated, conventional techniques of immunology, biochemistry, chemistry, molecular biology, microbiology, cell biology, genomics and recombinant DNA, which are within the skill of the art. See MOLECULAR CLONING: A LABORATORY MANUAL, 2nd edition (1989) (Sambrook, Fritsch and Maniatis); MOLECULAR CLONING: A LABORATORY MANUAL, 4th edition (2012) (Green and Sambrook); CURRENT PROTOCOLS IN MOLECULAR BIOLOGY (1987) (F. M. Ausubel, et al. eds.); the series METHODS IN ENZYMOLOGY (Academic Press, Inc.); PCR 2: A PRACTICAL APPROACH (1995) (M.J. MacPherson, B.D. Hames and G.R. Taylor eds.); ANTIBODIES, A LABORATORY MANUAL (1988) (Harlow and Lane, eds.); ANTIBODIES A LABORATORY MANUAL, 2nd edition (2013) (E.A. Greenfield ed.); and ANIMAL CELL CULTURE (1987) (R.I. Freshney, ed.).

[00435] The practice of the present invention employs, unless otherwise indicated, conventional techniques for generation of genetically modified mice. See Marten H. Hofker and Jan van Deursen, TRANSGENIC MOUSE METHODS AND PROTOCOLS, 2nd edition (2011).

[00436] The present invention also provides methods applicable to the study of bulk cells and is not limited to single cells. Moreover, the assays described herein are also amenable to regularly fixed and permeabilized cells (i.e. not using polymerization). The proximity assays described herein may be performed on cells without generating discrete polymer matrices. Additionally, detection of cellular constituents utilizing labeling ligands and a sequencing readout may be used to detect low abundant cellular constituents. Not being bound by a theory, the oligonucleotide label may be amplified and increase the signal as compared to antibody readouts known in the art. Moreover, determination of proteins in relation to open chromatin need not be performed in a polymer matrix.

[00437] The present inventions provides advantages over prior assays for detecting proteins and post translation modifications (PTM) in single cells or isolated aggregations of cellular constituents. Standard flow cytometry can be used to detect a few proteins/PTMs in greater than 10^6 single cells; and CyTOF (heavy metal labeling with multiplex barcoding) can be used to detect ~30-50 proteins/PTMs in 10^5 - 10^6 single cells. The present invention provides highly multiplexed, DNA sequencing-based readouts of protein/PTM levels of greater than 100's of proteins/PTMs in greater than 10^6 cells.

[00438] The present invention advantageously provides a Massively Combinatorial Perturbation Profiling (MCP) approach. Applicants can perturb vast numbers of combinations of genes, each targeting many circuit components at once. Applicants can use massively-parallel single cell genomics to measure genomic profiles and single cell proteomics to measure protein profiles after each perturbation. Applicants can infer the individual and combinatorial effects at each order, relying on random matrix theory, compressive sensing and kernel learning .

[00439] Biological systems are not linear: the combined effect of multiple factors is not simply the sum of their individual effects. This is a direct outcome of the biochemistry underlying molecular biology, from allosteric protein changes to cooperative binding, and is essential for cells to process complex signals. However, it has remained an insurmountable stumbling block to achieving a quantitative and predictive understanding of circuits on a genomic scale, with far-reaching implications for basic and translational science. Thus, the present invention provides a powerful combination by being able to measure transcriptional chromatin, epigenetic and proteomic changes as a function of genetic perturbation at the single cell level.

[00440] Combinatorial perturbation analyses have measured important genetic interactions, mainly from growth phenotypes in yeast. Mammalian studies have used ricin susceptibility and cell count phenotypes, but none combined large-scale, combinatorial genetic manipulation with complex, quantitative phenotypes, such as proteomic profiles. The single cell resolution readout of both response and perturbation, across many cells, serves as an improved starting point to unravel the function and interaction of the perturbed genes.

Sparse Coding

[00441] A simple, flexible, and cost-effective, transcriptome-wide gene-expression profiling solution that does not require measuring individual genes or single cell profiling is desired. This would greatly accelerate the rate of discovery of medically-relevant connections encoded therein by leveraging knowledge of relative abundances of genes to extrapolate underlying cell circuitry

[00442] The present invention relates to genomic informatics and gene-expression profiling. Gene-expression profiles provide complex molecular fingerprints regarding the relative state of a cell or tissue. Similarities in gene-expression profiles between organic states (i.e., for example, normal and diseased cells and/or tissues) provide molecular taxonomies, classification, and diagnostics. Similarities in gene-expression profiles resulting from various external perturbations (i.e., for example, ablation or enforced expression of specific genes, and/or small molecules, and/or environmental changes) reveal functional similarities between these perturbagens, of value in pathway and mechanism-of-action elucidation. Similarities in gene-expression profiles between organic (e.g. disease) and induced (e.g. by small molecule) states may identify clinically-effective therapies.

[00443] To achieve these and other advances and in accordance with the purpose of the present invention, as embodied and broadly described, in one aspect of the present invention, a probe set comprising 100 or more molecules assembled according to a set of random measurement values forming at least one measurement vector, where each molecule comprises a tag for the at least one measurement vector operably linked to a probe for one type of transcript of a plurality of types of transcripts.

[00444] In an aspect of the present invention, the probe set corresponds to a Design Matrix comprising $m \times n$ measurement values, where m is a number of measurement vectors and n is the number of types of transcripts. In another aspect, the tag uniquely corresponds to one of the

measurement vectors. In an aspect, the molecules of the probe set are single-stranded DNA. In an aspect, the tag is a barcode. In another aspect, the transcript is a gene.

[00445] In an aspect, the number of measurement vectors is approximately 100-30,000. In an aspect, the number of measurement vectors is based on an estimate of system sparsity. In an aspect, the number of measurement vectors is based on the log of the number of types of transcripts. In an aspect the number of measurement vectors is approximately $k \log(n)$, where k is an estimate of sparsity. In an aspect, the Design Matrix may be adjusted according to a basis. In an aspect, k is approximately equal to 10. In another aspect, m is less than n . In another aspect, n is greater than 10.

[00446] In another aspect of the invention, a method of measuring relative abundances of transcripts in a pool of samples, comprises generating a Design Matrix comprising $m \times n$ measurement values, where m is a number of measurement vectors and n is the number of types of transcripts; generating a probe library corresponding to the Design Matrix, wherein the probe library comprises a collection of molecules assembled according to the measurement values, where each molecule has a tag for one of the measurement vectors operably linked to a probe for one of the types of transcripts; contacting the probe library to the pool of samples, resulting in m measurement results for each sample of the pool of samples; generating an Observed Measurement Matrix M comprising the measurement results for each sample of the pool of samples; and applying a sparse coding solving process to the Observed Measurement Matrix M to learn system matrix S as indicative of relative abundance of the transcripts in each of the samples.

[00447] In another aspect of the present invention, the measurement values in the Design Matrix are independent. In another aspect, the measurement values in the Design Matrix are random.

[00448] In an aspect, the number of measurement vectors is approximately 100-30,000. In an aspect, the number of measurement vectors is based on an estimate of system sparsity. In an aspect, the number of measurement vectors is based on the log of the number of types of transcripts. In an aspect the number of measurement vectors is approximately $k \log(n)$, where k is an estimate of sparsity. In an aspect, the Design Matrix may be adjusted according to a basis. In an aspect, k is approximately equal to 10. In another aspect, m is less than n . In another aspect, n is greater than 10.

[00449] In another aspect of the present invention, the tag uniquely corresponds to one of the measurement vectors. In an aspect, the tag is a barcode.

[00450] In another aspect of the present invention, the molecules of the probe library are single-stranded DNA. In an aspect, each molecule in the probe library further comprises a tag for one of the samples.

[00451] In another aspect of the present invention, the contacting includes binding and the contacting may, in some aspects, include hybridization.

[00452] In another aspect of the present invention, the sample is a cell.

[00453] In another aspect of the present invention, generating the Observed Measurement Matrix include hybrid selection and tag quantification. In an aspect, tag quantification includes sequencing.

[00454] In another aspect of the present invention, a method for measuring relative abundances of n biomolecules in a pool of samples, comprises generating a Design Matrix comprising $m \times n$ measurement values, where m is a number of measurement vectors and n is the number of types of biomolecules; generating a probe library corresponding to the Design Matrix, wherein the probe library comprises a collection of molecules assembled according to the measurement values, where each molecule has a tag for one of the measurement vectors operably linked to a probe for one of the types of biomolecules; contacting the probe library to the pool of samples, resulting in m measurement results for each sample of the pool of samples; generating an Observed Measurement Matrix M comprising the measurements results for each sample in the pool of samples; and applying a sparse coding solving process to the Observed Measurement Matrix to learn system matrix S as indicative of relative abundance of the biomolecules in each of the samples.

[00455] In another aspect of the present invention, the measurement values in the Design Matrix are independent. In another aspect, the measurement values in the Design Matrix are random.

[00456] In an aspect of the present invention, the biomolecule is a transcript, protein, DNA, non-naturally occurring nucleic acid, peptide. In an aspect of the present invention, the samples include cells, blood, hair, nails, mucus, tissue, feces or urine. In an aspect of the present invention, the probe is a molecule that binds to the biomolecule. In an aspect, the probe is a

complex of molecules. In another aspect, the probe comprises an antibody or binding fragment thereof.

[00457] In another aspect of the present invention, the molecules of the probe library are single-stranded DNA. In an aspect, each molecule in the probe library further comprises a tag for one of the samples.

[00458] In another aspect of the present invention, the contacting includes binding and the contacting may, in some aspects, include hybridization.

[00459] In another aspect of the present invention, the tag uniquely corresponds to one of the measurement vectors. In an aspect, the tag is a barcode. In an aspect, each molecule in the probe library further comprises a tag for one of the samples.

[00460] In an aspect, the number of measurement vectors is approximately 1-30,000. In an aspect, the number of measurement vectors is based on an estimate of system sparsity. In an aspect, the number of measurement vectors is based on the log of the number of types of biomolecules. In an aspect the number of measurement vectors is approximately $k \log(n)$, where k is an estimate of sparsity. In an aspect, the Design Matrix may be adjusted according to a basis. In an aspect, k is approximately equal to 10. In another aspect, m is less than n . In another aspect, n is greater than 10.

[00461] In another aspect of the present invention, generating the Observed Measurement Matrix include tag quantification. In an aspect, tag quantification includes sequencing.

[00462] The term "device" as used herein, refers to any composition capable of measuring expression levels of transcripts. For example, a device may comprise a solid planar substrate capable of attaching nucleic acids (i.e., an oligonucleotide microarray). Alternatively, a device may comprise a solution-based bead array, wherein nucleic acids are attached to beads and detected using a flow cytometer. Alternatively, a device may comprise a nucleic-acid sequencer.

[00463] The term "probe" as used herein, refers to any molecule capable of attaching and/or binding to a nucleic acid (i.e., for example, a barcode nucleic acid). For example, a capture probe may be an oligonucleotide attached to a bead, wherein the oligonucleotide is at least partially complementary to another oligonucleotide. Alternatively, a capture probe may comprise a polyethylene glycol linker, an antibody, a polyclonal antibody, a monoclonal antibody, an Fab fragment, a biological receptor complex, an enzyme, a hormone, an antigen, and/or a fragment or portion thereof. A probe may be a nucleic acid sequence, the nucleic acid

being for example deoxyribonucleic acid (DNA), ribonucleic acid (RNA), peptide nucleic acid (PNA) or other non-naturally occurring nucleic acid.

[00464] The term "Design Matrix" as used herein, refers to a collection of the relative number of times up to which a type of transcript can be counted for a specific type of measurement. The collection may be described as a table and will be so in the following for the sake of the description. However, it is understood that the collection need not be generated as a table but can also be generated in any other form suitable for the measuring relative abundances of transcripts, e.g. a string of data. Each entry in the table is intended to be a randomly generated number, at least initially. In general, each row in the matrix corresponds to a specific type of measurement, while each column in the matrix corresponds to a type of transcripts in the sample pool. Thus, the number of rows in the Design Matrix corresponds to the number of specific types of measurements, while the number of columns in the Design Matrix corresponds to the number of types of transcripts in the sample pool. Of course, one of skill may choose to adjust the Design Matrix to reflect fewer than all available types of measurement and fewer than all types of transcripts in the sample pool. Further, one of skill in the art will appreciate that the Design Matrix itself may be transposed such that the number of rows indicates the number of types of transcripts, while the number of columns indicates the number of types of measurement.

[00465] The term "relative number of times" as used herein, means that if for a type of measurements, the relative number is a_1 for type of transcripts 1, a_2 for type of transcripts 2, ... a_n for type of transcripts n, then type of transcripts 1 can be counted up to $k.a_1$ times, type of transcripts 2 up to $k.a_2$ times, ... type of transcripts n up to $k.a_n$ times; k being an integer.

[00466] The term "count" and its derivatives as used in the definitions, encompasses any method that yield a measured value indicative of the count.

[00467] The term "Observed Measurement Matrix" as used herein, refers to a collection of measured values of a specific type of measurement for a specific type of sample in the sample pool or portion of sample (such as a cell, or a type of cell). The collection may be described as a table and will be so in the following for the sake of the description. However, it is understood that the collection need not be generated as a table but can also be generated in any other form suitable for the measuring relative abundances of transcripts, e.g. a string of data. In general, each row in the matrix corresponds to a specific type of measurement, while each column in the matrix corresponds to a type of samples in the sample pool or portion of sample. Thus, the

number of rows in the Observed Measurement Matrix corresponds to the number of specific types of measurements, while the number of columns in the Observed Measurement Matrix corresponds to the number of types of samples in the sample pool. Of course, one of skill may choose to adjust the Observed Measurement Matrix to reflect fewer than all available types of measurement and fewer than all types of samples in the sample pool. Further, one of skill in the art will appreciate that the Observed Measurement Matrix itself may be transposed such that the number of rows indicates the number of types of sample, while the number of columns indicates the number of types of measurement.

[00468] Correspondingly, any reference to a Matrix used herein, refers to a collection of values, which may be described as a table and will be so in the following for the sake of the description. However, it is understood that the collection need not be generated as a table but can also be generated in any other form suitable for the measuring relative abundances of transcripts, e.g. a string of data. Any reference to a vector (in the mathematical sense) used herein, refers to a collection of values which may be described as row or a column, in some instances row or column of a table, and will be so in the following for the sake of the description. However, it is understood that the collection need not be generated as a row or a column but can also be generated in any other form suitable for the measuring relative abundances of transcripts.

[00469] The term "Connectivity Map" as used herein, refers to a public database of transcriptome-wide gene-expression profiles derived from cultured human cells treated with a plurality of perturbagens, and pattern-matching algorithms for the scoring and identification of significant similarities between those profiles and external gene-expression data, as described by Lamb et al., "The Connectivity Map: using gene-expression signatures to connect small molecules, genes and disease". Science 313:1929 (2006). Build02 of the Connectivity Map contains 7,056 full-transcriptome gene-expression profiles generated with Affymetrix high-density oligonucleotide microarrays representing the biological effects of 1,309 small-molecule perturbagens, and is available at broadinstitute.org/cmap.

[00470] The term "query signature" as used herein, refers to any set of up- and down-regulated genes between two cellular states (e.g., cells treated with a small molecule versus cells treated with the vehicle in which the small molecule is dissolved) derived from a gene-expression profile that is suitable to query Connectivity Map. For example, a 'query signature' may comprise a list of genes differentially expressed in a distinction of interest; (e.g., disease

versus normal), as opposed to an 'expression profile' that illustrates all genes with their respective expression levels.

[00471] The term "connectivity score" as used herein, refers to a relative measure of the similarity of the biological effects of a perturbagen used to generate a query signature with those of a perturbagen represented in the Connectivity Map based upon the gene-expression profile of a single treatment with that perturbagen. For example, one would expect every treatment instances with vorinostat, a known histone deacetylase (HDAC) inhibitor, to have a high connectivity score with a query signature generated from the effects of treatments with a panel of HDAC inhibitors.

[00472] The term "enrichment score" as used herein, refers to a measure of the similarity of the biological effects of a perturbagen used to generate a query signature with those of a perturbagen represented in the Connectivity Map based upon the gene-expression profiles of multiple independent treatments with that perturbagen.

[00473] The term "small organic molecule" as used herein, refers to any molecule of a size comparable to those organic molecules generally used in pharmaceuticals. The term excludes biological macromolecules (e.g., proteins, nucleic acids, etc.). Preferred small organic molecules range in size from approximately 10 Da up to about 5000 Da, more preferably up to 2000 Da, and most preferably up to about 1000 Da.

[00474] The sample may be a biological sample, for example a blood, buccal, cell, cerebrospinal fluid, mucus, saliva, semen, tissue, tumor, feces, urine, and vaginal sample. It may be obtained from an animal, a plant or a fungus. The animal may be a mammal. The mammal may be a primate. The primate may be a human. In other embodiments, the sample may be an environmental sample, such as water or soil.

[00475] As used herein, the term "sample template" refers to nucleic acid originating from a sample which is analyzed for the presence of a target sequence of interest. In contrast, "background template" is used in reference to nucleic acid other than sample template which may or may not be present in a sample. Background template is most often inadvertent. It may be the result of carryover, or it may be due to the presence of nucleic acid contaminants sought to be purified away from the sample. For example, nucleic acids from organisms other than those to be detected may be present as background in a test sample.

[00476] “Target sequence” is intended to designate either one target sequence or more than one target sequence, *i.e.* any sequence of interest at which the analysis is aimed. Thus, the sample may comprise more than one target sequence and preferably a plurality of target sequences, the number of which may be 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20 and above.

[00477] The present invention is related to the field of genomic informatics and gene-expression profiling. Gene-expression profiles provide complex molecular fingerprints regarding the relative state of a cell or tissue. Similarities in gene-expression profiles between organic states (*i.e.*, for example, normal and diseased cells and/or tissues) provide molecular taxonomies, classification, and diagnostics. Similarities in gene-expression profiles resulting from various external perturbations (*i.e.*, for example, ablation or enforced expression of specific genes, and/or small molecules, and/or environmental changes) reveal functional similarities between these perturbagens, of value in pathway and mechanism-of-action elucidation. Similarities in gene-expression profiles between organic (*e.g.* disease) and induced (*e.g.* by small molecule) states may identify clinically-effective therapies. Improvements described herein allow for the efficient and economical generation of full-transcriptome gene-expression profiles by identifying cluster centroid landmark transcripts that predict the expression levels of other transcripts within the same cluster.

[00478] Some embodiments of the present invention contemplate measuring relative gene abundances of transcripts in a pool of samples to allow genome-wide transcriptional profiling for applications including, but not limited to, disease classification and diagnosis without resort to expensive and laborious microarray technology (*i.e.*, for example, Affymetrix GeneChip microarrays). Other uses include, but are not limited to, generating gene-expression data for use in and with information databases (*i.e.*, for example, connectivity maps). A connectivity map typically may comprise a collection of a large number of gene-expression profiles together with allied pattern-matching software. The collection of profiles is searched with the pattern-matching algorithm for profiles that are similar to gene-expression data derived from a biological state of interest. The utility of this searching and pattern-matching exercise resides in the belief that similar biological states may be identified through the transitory feature of common gene-expression changes. The gene-expression profiles in a connectivity map may be derived from known cellular states, or cells or tissues treated with known chemical or genetic perturbagens. In

this mode, the connectivity map is a tool for the functional annotation of the biological state of interest. Alternatively, the connectivity map is populated with gene-expression profiles from cells or tissues treated with previously uncharacterized or novel perturbagens. In this mode, the connectivity map functions as a screening tool. Most often, a connectivity map is populated with profiles of both types. Connectivity maps, in general, establish biologically-relevant connections between disease states, gene-product function, and small-molecule action. In particular, connectivity maps have wide-ranging applications including, but not limited to, functional annotation of unknown genes and biological states, identification of the mode of action or functional class of a small molecule, and the identification of perturbagens that modulate or reverse a disease state towards therapeutic advantage as potential drugs. See Lamb et al, "The Connectivity Map: using gene-expression signatures to connect small molecules, genes and disease" *Science* 313: 1929-1935 (2006), and Lamb, "The Connectivity Map: a new tool for biomedical research" *Nature Reviews Cancer* 7: 54-60 (2007). However, the high cost of generating gene-expression profiles severely limits the size and scope of connectivity maps. A connectivity map populated with gene-expression profiles derived from every member of an industrial small-molecule drug-screening library, a saturated combinatorial or diversity-orientated chemical library, a comprehensive collection of crude or purified plant or animal extracts, or from the genetic ablation or forced expression of every gene in a mammalian genome, for example, would be expected to facilitate more, and more profound, biological discoveries than those of existing connectivity maps.

[00479] The present invention contemplates compositions and methods for making and using a transcriptome-wide gene-expression profiling platform that "under samples" the total number of transcripts. Because gene expression is believed to be highly correlated, direct measurement of a small number allows the expression levels of the remainder to be inferred. The present invention, therefore, has the potential to reduce the cost and increase the throughput of full-transcriptome gene-expression profiling relative to the well-known conventional approaches that require all transcripts to be measured.

[00480] Gene expression data are highly structured, such that the expression level of some genes is predictive of the expression level of others. Knowledge that gene expression data are highly structured allows for the assumption that the number of degrees of freedom in the system

are small, which allows for assuming that the basis for computation of the relative gene abundances is sparse.

[00481] It is possible to make several biologically motivated assumptions that allow Applicants to recover the nonlinear interaction terms while under-sampling without having any specific knowledge of which genes are likely to interact. In particular, if Applicants assume that genetic interactions are low rank, sparse, or a combination of these, then the true number of degrees of freedom is small relative to the complete combinatorial expansion, which enables Applicants to infer the full nonlinear landscape with a relatively small number of perturbations. Working around these assumptions, analytical theories of matrix completion and compressed sensing may be used to design under-sampled combinatorial perturbation experiments. In addition, a kernel-learning framework may be used to employ under-sampling by building predictive functions of combinatorial perturbations without directly learning any individual interaction coefficient.

[00482] In the framework of matrix completion, or more generally tensor completion, the goal is to fill in all the values of a matrix using a small collection of sampled entries. Applicants hypothesize that the rank of a matrix of second-order interactions, or tensor of higher-order interactions, is a fraction of the number of genes being tested. Applicants can test this idea by explicitly calculating the rank from a dense sampling of second or third order knockouts from a small collection of genes. If the rank of interactions is indeed limited, then Applicants can randomly sample pairs of genes to knockout from a larger collection (in the second order case), and then fill in the remaining values via nuclear norm regularized least-squares optimization [Matrix Completion with Noise. E. Candes and Y. Plan]. Provable guarantees for this recovery suggest that if the rank, r , is small relative to the number of genes, n , then m samples such that

$$\mathbf{[00483]} \quad m \geq O\left(\frac{6}{5}nr \log(n)\right)$$

[00484] are sufficient. However, these guarantees assume rough uniformity in the loadings of interaction singular vectors, and this assumption is unlikely to hold up if the interaction matrix itself is very sparse. In this case, Applicants can perform the same random sampling scheme, and then simultaneously regularize over both the nuclear norm, and L^1 norm of the interaction matrix [Estimation of Simultaneously Sparse and Low Rank Matrices. Richard et. al.]. As may be appreciated by one of skill in the art, the “sparsity” of a matrix is a measure of the non-zero

elements of a matrix relative to the total number of elements of a matrix. A sparse matrix is a matrix in which most of the elements are zero, which is indicative of loose correlation between systems.

[00485] As may be appreciated by one of skill in the art, the “rank” of a matrix is the maximum number of linearly independent row vectors in the matrix or the maximum number of linearly independent column vectors in the matrix. The rank of a matrix denotes the “information content of the matrix. The lower the rank, the lower is the information content. A basis for a vector space is a collection of vectors that form a set that is linearly independent and that spans the space.

[00486] In the related setting of compressed sensing Applicants can directly leverage the sparse nature of interactions. In this case, instead of working with a tensor of interaction terms, Applicants work with a basis that spans all higher order interactions. This can be, for example, the Fourier basis of higher-order polynomials, in which case Applicants assume that there is a small number, s , of nonzero Fourier coefficients, which are not known a priori, but can accurately capture the effects of any combinatorial perturbation. Then, if Applicants make m random combinatorial knockouts, where $m \geq O(s \log f)$ and f is the magnitude of combinatorial expansion (e.g. $f = \binom{100}{3} = \text{“100 choose 3”} = 161700$), Applicants can recover the nonzero coefficients, and, thus, the set of nonlinear interaction terms, through L^1 regularized regression. Compressed sensing tells Applicants that if the perturbations are de-coherent under the given basis, then exact recovery is possible with dramatic under-sampling (in the noiseless case) [Compressive sampling. E. Candes.]. Applicants use this measure of de-coherence, then, as a guide in designing combinatorial perturbations. In particular, if each of m experiments includes a Poisson random sampling of knockouts, then Applicants expect the measurements to have good de-coherence under the Fourier basis, provided the mean number of knockouts is not too low (e.g. > 1).

[00487] The formulations for matrix completion and compressed sensing each rely on a degree of sparsity – either in the rank or in the coefficients – to achieve recovery with under-sampling. For each of these Applicants can estimate the level of sparsity by performing dense combinatorial perturbations on a small collection of genes, or, in the case of compressed sensing, simply randomly sampling higher-order knockouts. However, Applicants also consider the

alternative formulation of kernel learning, which does not have such strict demands on sparsity. In this case Applicants build predictive functions of the effects of combinatorial perturbations using a kernel of experimental similarity. More specifically, assume m experiments with Poisson random sampling of knockouts for each measurement. Then, Applicants can define an m -by- m polynomial kernel, for example, based on the overlap in knockouts between any pair of experiments. If Applicants build such a kernel, and learn the weighted combination of kernel vectors that fits a collection of training data, Applicants can then use these coefficients to predict the outcome of new experiments. In this case the density of nonlinear interaction terms can be much greater, since Applicants do not directly learn any particular interaction coefficient, but rather a kernelized version of the entire polynomial. In fact, if the interaction terms are too sparse, the kernel learning framework is unlikely to be successful with any significant under-sampling. However, together, kernel learning, matrix completion, and compressed sensing represent a complimentary range of approaches for inferring the effects of higher-order, combinatorial perturbations, with different assumptions on the underlying structure of the data for each framework.

[00488] The overall sparsity level of biological networks has been estimated in several ways. By comparing colony sizes between single and double mutants, the genetic landscape of *S. cerevisiae* has been mapped revealing that <5% of 5.4 million gene pairs interact significantly with respect to fitness. [Costanzo*, Baryshnikova* et. al., 2011] In mammalian cells, a systematic screen for ricin toxicity using a hierarchical approach revealed that among 60 genes with significant individual effects, interactions broadly divided genes into functional complexes and were of a relatively higher density [Bassik et. al., 2013]. A higher dimensional phenotype using imaging between 282 target genes, 20 query genes, and 11 phenotypes found approximately 5,000 significant interactions [Laufer et. al., 2013]. Sparsity with respect to gene expression phenotypes has not been well studied.

[00489] The empirical limits on the number of DSB that might be tenable for an individual cell to survive vary, but rough estimates using imaging of phosphorylation of H2AX and 53BP1 suggest this number could be in the 10-50 range.

[00490] According to embodiments of the present invention, each measurement integrates signal across many genes (for example, all 29k genes). Thus, the measurements are not sensitive to the stochastic capture of any single gene, but rather to the average signal across a broad range

of genes. In this sense the present scheme is robust to technical noise. In addition, probe sets according to the present invention can be designed independently from any existing data. The probes are, in this sense, universally appropriate for any system desired to be measured. This method recognizes that there may be many different signature gene networks across many cell types / tissues, but that these are used to a sparse degree in any single cell. That is, in a single cell the number of active gene networks is small relative to the total number of all existing gene networks. This assumption can be used to under sample the transcriptome, such that a computational analysis is designed to recover signal from sparse systems. It is assumed that the gene networks (or their abstracted analogs) are sparsely used in any single cell, and specific knowledge of when the networks are active, or even the definitions of the networks themselves, are not required. Every gene is represented in every measurement, even if exact correlations between measurements, genes or expressions are not known. By analyzing transcriptomic measurements across many cells or tissues, it is possible to identify structures in the data that reflect the underlying cell circuitry. These same structures can be exploited to recover gene abundances, while dramatically under-sampling the full signal. In particular, if a given cell or tissue has a sparse representation in some basis, then Applicants can use the theories of compressive sensing and sparse coding to guide the design of RNA probe libraries that Applicants believe reduces sequencing requirements nearly 1000-fold.

[00491] To design compressive measurements, Applicants first consider the canonical basis described by Principle Component Analysis (PCA) as an example of an embedding for RNA-Seq data. Suppose, for the sake of argument, that Applicants knew of a universal set of 100 PCs, and that a linear combination of these could be used to describe the gene abundances in any system. By building molecular probes for each PC, Applicants could measure the response along each component in a given sample, and then take the linear combination of these 100 measurements to recover the abundance of every gene. Generally speaking, a basis that is universally appropriate for all systems (in a sense, this is an ultimate goal of biology) is not known. However, there is a surprising result which tells Applicants that making random measurements is universally appropriate for any basis that Applicants are likely to encounter [Compressive sampling. In Proc. Int. Congress of Math., Madrid, Spain, Aug. 2006.]. After making random measurements, and using methods of sparse coding [Efficient sparse coding algorithms. In NIPS, 2006.], Applicants

can simultaneously infer the basis and its sparse contributions to the samples in question, and recover the desired signal.

[00492] Measurements may be made via hybridization with targeted gene probes, which are barcoded and sequenced to count hybridization events. In such case, each measurement uses a given collection of probes, and there may be multiple probes per gene. For one of m measurements, a random library may be constructed from a barcode pool, and a pool of gene-specific oligos. Each oligo in the barcode pool may contain a universal adapter, a cell/experiment ID, a molecular ID, a measurement ID (one of m total), and a linker region. Gene specific oligos may contain a complimentary linker region, and a target sequence for the given gene. The final probe pool may be created by annealing and extension reactions.

[00493] The choice of probe sets may be a random design. Random probe sets have several advantages. Two random measurements are orthogonal (not correlated) with high probability, which means that each random measurement is effectively measuring something “new.” This is difficult, if not impossible, to ensure when each measurement consists of a single gene, and, thus, random probes do a better job of maximizing the information content of a small number of measurements. Existence of a gene correlation structure is assumed without knowing what it is. The probe sets are substantially universally appropriate. By randomly scattering measurements throughout the transcriptome, it can be ensured that Applicants sample from the relevant structures with high probability. Applicants then employ computational methods of Sparse Coding to learn the basis, or gene expression structure which is appropriate for the sample at hand. Importantly, the learning is integrated into the same process for inferring the full gene expression profile, and does not require a separate set of measurements. The number of types of measurements can be approximated by $k \log(n)$, where k is an estimated value of sparsity.

[00494] At least two methods for randomizing the amount of probe for each gene in each library are considered. First, the molecular IDs are synthesized randomly. Each of these IDs are assigned to $+1/-1$, such that when counting hybridization events for a given measurement Applicants either increment or decrement the total depending on the molecular ID. Thus, Applicants could have the same number of probes for every gene in each measurement, but the sum across molecular IDs for a given gene is a binomial random variable. Alternatively, a Liquid Handling Robot may be used to mix the barcode and gene pools according to a randomly

generated matrix. In either case, the probe library is sequenced after construction to determine its composition.

[00495] In the present system, measurements are read out via sequencing, after hybrid selection and amplification. The number of types of measurements is contemplated to be between approximately 100 and 30,000 different types of measurements. The number of types of measurements can be based on an estimate of system sparsity. The number of types of measurements per sample can be, for example, 100. To the extent that there is any knowledge of assumptions about certain measurements or between genes, the number and type of measurements may be adjusted. For each measurement barcode within a cell, Applicants can have gene specific probes for each of ~29k genes. However, the number of probes for a given gene will vary across each of the 100 measurements, according to random design. Once types of measurements are established, and the number or transcripts to be measured are determined, a Design Matrix having a number of rows corresponding to the number of types of measurements (measurement vectors) and the number of column corresponding to the number of transcripts can be generated. The Design Matrix is populated by “measurement values,” which are random numbers of measurements such that for each type of measurement each transcript can be counted up to a relative number of times that is random and independent of other elements/entries in the Design Matrix. In other words, the relative number of times up to which a certain measurement type can count each respective transcript is random and independent of the relative number of times up to which that certain measurement type can count other transcripts in the sample and the relative number of times up to which a certain transcript can be counted for a certain type of measurement is random and independent of the relative number of times up to which that certain transcript can be counted for other types of measurement. While the measurements may be random and independent, such is not required. Knowledge of some relationships between measurements may allow manipulation of the measurement values to improve calculations. For example, steps according to the present application may be iterated multiple times to refine results by incorporating knowledge gained through a first round of measurements. Similarly, knowledge of relationships between genes, measurements and transcripts may be used to inform and select some measurement values for the Design Matrix, although this is not required. For example, the Design Matrix can be adjusted according to a known basis.

[00496] A probe library corresponding to the Design Matrix may be generated. For example, a probe library comprises 100 or more molecules assembled according to a set of random measurement values forming at least one measurement vector, where each molecule comprises a tag for the at least one measurement vector operably linked to a probe for one type of transcript of a plurality of types of transcripts. The probe library includes a collection of molecules assembled according to the measurement values, where each molecule has a tag, e.g., a barcode, for one of the measurement types operably linked to a probe for one of the types of transcripts. The tag may uniquely correspond to a measurement type/vector. Probes contemplated according to the present invention may be single-stranded DNA, transcript, protein, DNA, non-naturally occurring nucleic acid, peptide, or the like. That is, each amount of probe can thus be adjusted based on the random measurement values in the Design Matrix such that the amount of a specific probe is known relative to the amounts of the other probes based on the Design Matrix.

[00497] While the present example describes cells as the sample for measurement, the sample may include cells, blood, hair, nails, mucus, tissue, feces, urine, body secretion or the like.

[00498] For one application, measurements may be made in single cells, and probe sets may be held as oligos on beads. For example, in one aspect each oligo may hold at least one of a sequencing adapter, a cellular barcode, a measurement barcode, a molecular barcode, and a primary gene-specific probe. cDNA targets may be captured on the primary probes, and then secondary probes, which target sequences immediately adjacent to the primary probes, may be introduced. Ligation of the primary and secondary probes produces a full-length fragment that can be amplified via universal sequences on the 5' end of the bead-attached oligo, and 3' end of the free-floating secondary probe. After amplification, Applicants sequence the barcode-containing regions to count the number of molecular events for each cell / measurement barcode combination. Contacting the probe library to the pool of samples thus results in a number of measurement results for each cell of the pool of cells.

[00499] The sample types are not limited to cells and may be any of numerous types of biomolecules. A biomolecule is any molecule that is present in living organisms, such as large macromolecules such as proteins, polysaccharides, lipids, and nucleic acids, as well as small molecules such as primary metabolites, secondary metabolites, and natural products. In an advantage embodiment, the biomolecule is a nucleic acid, such as but not limited to, deoxyribonucleic acid (DNA), ribonucleic acid (RNA), peptide nucleic acid (PNA) or other non-

naturally occurring nucleic acid. In another advantageous embodiment, the biomolecule is a protein, such as but not limited to, peptides, antibodies, immunogenic molecules or enzymes.

[00500] An Observed Measurement Matrix may thus be constructed using the measurement results for each cell/sample to populate the elements/entries of the Observed Measurement Matrix. Construction of the Observed Measurement Matrix can include tag quantification, and tag quantification can include sequencing. For example, a relative count is conducted of the tags such that the relative number of times a certain tag is counted per cell is entered into an appropriate element of the Observed Measurement Matrix M . A sparse coding solving process may then be applied to the Observed Measurement Matrix M to learn system matrix S as indicative of relative abundance of the transcripts in each of the samples.

[00501] An example of a sparse coding process can be exemplified by

$$\begin{array}{ccccc} \mathbf{M} & = & \mathbf{DM} & \cdot & \mathbf{S} \\ m \times c & & m \times n & & n \times c \end{array}$$

[00502] where m is the number of measurement vectors, c is the number of cell vectors, and n is the number of gene vectors. S is a System Matrix, which can be determined once M and DM are populated (M with actual measurements/counts and DM with random measurement values). S is thus indicative of the relative gene abundances.

[00503] Moreover, once S is known, a Basis Matrix B , can be determined based on knowledge of X , a matrix populated by known contributions of genes n and cells c to the sample.

$$\begin{array}{ccccc} \mathbf{S} & = & \mathbf{B} & \cdot & \mathbf{X} \\ n \times c & & n \times n & & n \times c \end{array}$$

[00504] Knowledge of the Basis can then be used to generate new measurement values of the Design Matrix to refine calculation of a new System Matrix S' .

[00505] Testing for linearity of hybrid selection is contemplated, along with dynamic range, by, for example, splitting bulk cDNA prepared from a population of K562 cells. In one split, qRT-PCR is performed for a collection of 25 genes that range from the highest abundance, to being absent. In a second split, hybrid selection is tested on a collection of 24 random probe

libraries (3 experimental IDs, each with 8 measurement indices). The goal is to optimize the hybridization conditions for maximal dynamic range, while ensuring that the readout from sequencing after hybridization is linear with regards to qPCR controls, and the amount of probe for a given gene in a given library. The size of the probe library is then be successively increased to 100 genes, then 1000, and finally to the set of all genes. The measurements are more robust with larger panels of genes, since Applicants' readout is effectively the sum of hybridization events across all probes for a given measurement index. In other words, since only barcodes are sequenced the barcodes and not the gene specific regions, the process is not sensitive to the stochastic capture of transcripts for any single gene, but rather to the average capture across all genes. Finally, this design is compatible with single cell technologies, and has indeed been designed with this application in mind. In particular, random probe libraries held on beads in place of barcoded poly-dT oligos are contemplated, but this is just one of many possible constructions.

[00506] In a first round of experiments, Applicants have constructed the random probe sets as follows. One pool of oligos is synthesized to contain a universal sequence, a cellular barcode, a measurement barcode, a molecular barcode, and a linker region. A second pool is synthesized to contain a complimentary linker region, and a gene-specific target. Oligos from the two pools can be annealed, and then extended in the 3' direction on the barcoded fragment with the Klenow fragment, or T4 DNA polymerase. These can be treated with lambda exonuclease, and purified by size selection for the full-length fragment. This can be done for each gene, and then the fragments can be mixed and pooled according to random design. Alternatively, fragments for each gene in each measurement can be pooled in equal amounts. In this case, Applicants achieve random measurements by assigning each molecular identifier (which was synthesized randomly) to a +1 / -1. When integrating the counts for a given measurement, Applicants then increment or decrement the count according to this assignment. Since the molecular barcodes and assignments are made randomly, the sum of +1 / -1 assignments for a given gene behave as a random variable across the different measurements. Secondary gene probes can be synthesized independently.

[00507] It is anticipated that probe libraries may be constructed and provided and present methods applied to provide cost-effective gene-expression profiling. Accordingly a product comprising completed probe libraries, or a kit for making the libraries according to principles of

the present invention are possible. For some applications it may make sense to focus on a subset of genes, in which case the kit would be particularly appropriate.

[00508] In addition, random probes and sparse coding could be applied as a service. That is, a customer may provide a sample or sample pools to an entity for constructing a probe library according to customer objectives and sample, applying the steps of the methods according to principles of the present invention and providing as its result a report of the relative gene-abundances or gene-expression profile or the like. This could be particularly useful for applications that require the analysis of tens to hundreds of thousands of transcriptional profiles, since Applicants estimate that methods according to principles of the present invention reduce sequencing requirements by 100-1000 fold. For example, when screening for combinatorial effects of small molecules, or when testing for the presence of very rare cell types (such as cancer stem cells), a large number of profiles need to be generated.

[00509] Rapidly advancing experimental methods are opening new ground for the exploration of higher-order genetic interactions. These methods include the CRISPR-Cas system for knockouts, which enables greater degrees of combinatorial perturbations, as well as high-throughput readout techniques such as single-cell RNA-Seq. The ability to sequence the transcriptomes of tens of thousands of cells using droplet microfluidics has recently been demonstrated [Macosko et. al., 2015]. By constructing a polyadenylated RNA barcode that is associated with a CRISPR sgRNA, Applicants was able to couple the experimental scale of microfluidic technologies with perturbation experiments in mammalian cells.

[00510] However, even with these advances, Applicants are limited to full combinatorial samplings of second- or third-order interactions for a small number of genes. For example, for a collection of 30 genes, performing all third-order knockouts is at the boundary of experimental feasibility ($30 \text{ choose } 3 = 4060$ experiments), while a dense sampling of fourth-order terms for the same size collection ($30 \text{ choose } 4 = 27405$), or third-order terms for a larger collection ($100 \text{ choose } 3 = 161700$) is intractable. Therefore, Applicants are motivated to identify schemes that leverage Applicants' ability to perform combinatorial perturbations, while also dramatically under-sampling the full combinatorial space.

[00511] **Compressive Sensing the Transcriptome**

[00512] **1st Trial Run – Analysis**

[00513] **Probe library creation**

[00514] The new protocols generated a much higher percentage of fragments of the correct size (eg 70-99% stitched) as compared with the original L1000 protocol (19-42% stitched). Having larger cDNA input also resulted in a greater percentage of stitched fragments.

[00515] UMI counting

[00516] Correlations of individual gene counts with expectation based on RPKM from ENCODE were generally positive, but not stunning. Results with gene counts from qPCR were largely unchanged. For the new ligation protocols, correlations were generally positive, except for row1 which generally has low (negative) values. For the L1000 ligation protocol correlations were more consistently positive, except for condition 21, which shows negative correlations for many rows/conditions.

[00517] Applicants also considered UMIs aggregated across genes, with randomly chosen rows. This mimics the intended setting. In this case correlations were generally much higher (range 47-79%). The best condition by this measure was 26_27 (new protocol with low ligation, low L/R ratio with undiluted cDNA)

[00518] The normalization strategy used in initial trials is based on sequenced barcode counts from direct ligation of probes, without cDNA targets. Quantifying barcode counts in this way gives an overall estimate of the abundance of each member of the library. In future experiments, applicants will compare the results of this normalization strategy to a normalization based on ligation to a pool of genomic DNA.

[00519] **X18 samples (3 FBC vs 6 RBC)**

[00520] <P5 barcoding>

[00521] Original L1000: P5-30

[00522] New protocol (high ligation 5ul mix input for R1, see notes below): P5-25

[00523] New protocol (low ligation 1ul mix input for R1, see notes below): P5-26

[00524] <P7 barcoding>

[00525] 1. High L/R ratio probe mix + 5µL cDNA input (~50ng/µL): P7-14

[00526] 2. High L/R ratio probe mix + 5µL (1/8 diluted) cDNA input (~6ng/µL): P7-21

[00527] 3. Low L/R ratio probe mix + 5µL cDNA input (~50ng/µL): P7-27

[00528] 4. Low L/R ratio probe mix + 5µL (1/8 diluted) cDNA input (~6ng/µL): P7-28

[00529] 5. High L/R ratio probe mix direct ligation with T7: P7-15

[00530] 6. Low L/R ratio probe mix direct ligation with T7: P7-20

[00531] Read Counts

		Total		
F	R	Fragments	Stitched	% Stitched
26	27	799365	795574	99.5%
25	27	890612	881037	98.9%
26	20	839036	825118	98.3%
26	14	401287	383828	95.6%
26	28	628901	589907	93.8%
26	21	237079	217109	91.6%
25	14	124932	109626	87.7%
26	15	1067299	906534	84.9%
25	15	1033888	838764	81.1%
25	28	726668	570421	78.5%
25	20	849166	628951	74.1%
25	21	65325	46374	71.0%
30	20	437882	186679	42.6%
30	14	774238	257647	33.3%
30	15	938377	268906	28.7%
30	21	745495	144001	19.3%
30	28	1066525	189338	17.8%
30	27	966163	149402	15.5%

[00532] Correlation with RPKM

[00533] For each read Applicants parse barcodes into structures such as:

[00534] cell1 AGGTCGTA row1 ACTB ACTB

[00535] Stitched reads that do not fit this format are discarded. For direct ligation (no cDNA), Applicants allow the left and right gene barcodes to differ, but for all others Applicants require that the gene barcodes are identical.

[00536] Once parsed, Applicants count the number of UMIs for every gene in each condition and row (8 rows per condition). After normalizing by the number of UMIs in the corresponding

direct ligation condition, Applicants can compare the normalized UMI count with actual RPKM for each gene.

[00537] For a given experimental condition, probes were included at 8 dilution levels (rows 1-8). In row 1, probes for all targets were included at approximately the same concentration. In each successive dilution level (rows 2-8) the overall concentration of probes was serially diluted so that the probe concentration in row 8 was approximately 100-fold lower than the concentration in row 1. Thus, correlations of observed versus expected counts within each of the 8 dilution levels of a given experiment reflect the sensitivity of each assay across a 100-fold change in probe concentration. Correlations from these results are included at the end of this document. In Figure 3 the results for one example condition (24_26) are depicted in a scatter plot.

[00538] Applicants can also ask how each condition would perform in the random setting where Applicants aggregate counts across all genes. To do this Applicants begin with a single gene, and choose a random row from the given condition. Applicants increment the “observed” count by the number of UMIs for that gene in that row, and increment the “expected” count by the product of the direct ligation UMI count for that gene/row with the RPKM of the gene. This is repeated for every gene such that Applicants have “observed” and “expected” sums across all genes with randomly chosen rows. This was repeated 8 times to generate a vector of random observations and expectations. Below Applicants report the correlation of these vectors for each condition, averaged across 50 random trials. In Figure 4 the simulation results for one example condition (24_26) are depicted in a scatter plot.

condition	correlation
25_21	46.6%
25_14	49.0%
26_21	52.9%
30_14	57.8%
30_27	63.4%
30_21	65.0%
30_28	68.9%
25_27	69.7%

26_14	70.2%
25_28	74.1%
26_28	76.6%
26_27	79.3%

[00539] Results for individual genes in each condition (condition IDs offset by 1):

[00540] norm.24_13.row1.txt:correlation: 0.101860
[00541] norm.24_13.row2.txt:correlation: 0.198498
[00542] norm.24_13.row3.txt:correlation: 0.125971
[00543] norm.24_13.row5.txt:correlation: 0.184587
[00544] norm.24_13.row6.txt:correlation: -0.065662
[00545] norm.24_13.row7.txt:correlation: 0.548366
[00546] norm.24_13.row8.txt:correlation: 0.188612
[00547] norm.24_20.row1.txt:correlation: 0.018865
[00548] norm.24_20.row2.txt:correlation: 0.011796
[00549] norm.24_20.row3.txt:correlation: -0.067194
[00550] norm.24_20.row5.txt:correlation: -0.220015
[00551] norm.24_20.row6.txt:correlation: -0.157439
[00552] norm.24_20.row7.txt:correlation: 0.364206
[00553] norm.24_20.row8.txt:correlation: 0.080670
[00554] norm.24_26.row1.txt:correlation: 0.372970
[00555] norm.24_26.row2.txt:correlation: 0.778385
[00556] norm.24_26.row3.txt:correlation: 0.794532
[00557] norm.24_26.row5.txt:correlation: 0.501830
[00558] norm.24_26.row6.txt:correlation: 0.340869
[00559] norm.24_26.row7.txt:correlation: 0.562235
[00560] norm.24_26.row8.txt:correlation: 0.576566
[00561] norm.24_27.row1.txt:correlation: 0.381881
[00562] norm.24_27.row2.txt:correlation: 0.771273
[00563] norm.24_27.row3.txt:correlation: 0.684149
[00564] norm.24_27.row5.txt:correlation: 0.445928

[00565] norm.24_27.row6.txt:correlation: 0.170683
[00566] norm.24_27.row7.txt:correlation: 0.543163
[00567] norm.24_27.row8.txt:correlation: 0.530076
[00568] norm.25_13.row1.txt:correlation: 0.309518
[00569] norm.25_13.row2.txt:correlation: 0.490934
[00570] norm.25_13.row3.txt:correlation: 0.349481
[00571] norm.25_13.row5.txt:correlation: 0.397947
[00572] norm.25_13.row6.txt:correlation: 0.272807
[00573] norm.25_13.row7.txt:correlation: -0.158994
[00574] norm.25_13.row8.txt:correlation: 0.256215
[00575] norm.25_20.row1.txt:correlation: 0.141412
[00576] norm.25_20.row2.txt:correlation: 0.275601
[00577] norm.25_20.row3.txt:correlation: -0.170215
[00578] norm.25_20.row5.txt:correlation: -0.115663
[00579] norm.25_20.row6.txt:correlation: -0.306711
[00580] norm.25_20.row7.txt:correlation: -0.590045
[00581] norm.25_20.row8.txt:correlation: -0.327281
[00582] norm.25_26.row1.txt:correlation: 0.513396
[00583] norm.25_26.row2.txt:correlation: 0.135664
[00584] norm.25_26.row3.txt:correlation: 0.493248
[00585] norm.25_26.row5.txt:correlation: 0.592195
[00586] norm.25_26.row6.txt:correlation: 0.649929
[00587] norm.25_26.row7.txt:correlation: 0.638488
[00588] norm.25_26.row8.txt:correlation: 0.592366
[00589] norm.25_27.row1.txt:correlation: 0.435686
[00590] norm.25_27.row2.txt:correlation: -0.043748
[00591] norm.25_27.row3.txt:correlation: 0.249897
[00592] norm.25_27.row5.txt:correlation: 0.615696
[00593] norm.25_27.row6.txt:correlation: 0.591814
[00594] norm.25_27.row7.txt:correlation: 0.550319
[00595] norm.25_27.row8.txt:correlation: 0.433558

[00596] norm.29_13.row1.txt:correlation: -0.351870
[00597] norm.29_13.row2.txt:correlation: 0.144392
[00598] norm.29_13.row3.txt:correlation: 0.057765
[00599] norm.29_13.row5.txt:correlation: 0.393736
[00600] norm.29_13.row6.txt:correlation: 0.031146
[00601] norm.29_13.row7.txt:correlation: 0.043738
[00602] norm.29_13.row8.txt:correlation: 0.382760
[00603] norm.29_20.row1.txt:correlation: -0.092086
[00604] norm.29_20.row2.txt:correlation: 0.040947
[00605] norm.29_20.row3.txt:correlation: 0.265565
[00606] norm.29_20.row5.txt:correlation: 0.589432
[00607] norm.29_20.row6.txt:correlation: 0.117522
[00608] norm.29_20.row7.txt:correlation: 0.386742
[00609] norm.29_20.row8.txt:correlation: 0.557595
[00610] norm.29_26.row1.txt:correlation: -0.218187
[00611] norm.29_26.row2.txt:correlation: 0.578199
[00612] norm.29_26.row3.txt:correlation: 0.549576
[00613] norm.29_26.row5.txt:correlation: 0.131541
[00614] norm.29_26.row6.txt:correlation: 0.271451
[00615] norm.29_26.row7.txt:correlation: 0.429337
[00616] norm.29_26.row8.txt:correlation: 0.189393
[00617] norm.29_27.row1.txt:correlation: -0.246045
[00618] norm.29_27.row2.txt:correlation: 0.449658
[00619] norm.29_27.row3.txt:correlation: 0.425474
[00620] norm.29_27.row5.txt:correlation: 0.376062
[00621] norm.29_27.row6.txt:correlation: -0.444794
[00622] norm.29_27.row7.txt:correlation: -0.125790
[00623] norm.29_27.row8.txt:correlation: 0.399004
[00624] In a second round of optimizations, the probe library was expanded to target 24 genes that range from the most abundant, to non-abundant in the sample. As illustrated in the plots exemplified in Figures 5-8, correlation values for individual dilution levels (rows 1-8) were

significantly higher than initial experiments. Correlation levels may be determined for a probe library and various dilution conditions of that probe library. For example, Figure 5 is plot of expected counts versus observed counts for a predetermined set of probes under a range of dilution criteria for the protocol Purified cDNA High. In the experiment exemplified, 24 different genes from most abundant to not present were measured and correlation values determined within a dilution level, where a correlation of closer to 1 indicates a higher degree of correlation. In the first dilution level (row 1), probes for each of 24 genes were included at approximately equal abundance. In the second dilution level (row 2) the overall concentration of probes was reduced by a factor of 2. Each row down to row 8 was successively diluted, so that the concentration of probes in row 8 was approximately 100-fold lower than in row 1. Thus, the 8 correlation values in Figure 5 represent the sensitivity of the assay across a 100-fold change in the concentration of gene probes.

[00625] Similarly, Figure 27 is plot of expected counts versus observed counts for a predetermined set of probes under a range of dilution criteria for unPurified cDNA High.

[00626] Figure 28 is plot of expected counts versus observed counts for a predetermined set of probes under a range of dilution criteria for Purified cDNA Low.

[00627] Figure 29 is plot of expected counts versus observed counts for a predetermined set of probes under a range of dilution criteria for unPurified cDNA Low.

[00628] The same data can be used to computationally simulate the random setting, where counts are aggregated across all genes from randomly chosen rows, as described above for initial optimizations. Results for these simulations are presented in Figures 30-33.

[00629] The population of cells with a plurality of genomic sequence or perturbation conditions involves a plurality of cells and perturbations to be tested and measurements sampled to obtain meaningful data and to infer appropriate circuits. The number of genes perturbed, and how many are perturbed simultaneously (the order of the perturbation, pairs, triplets, etc.) varies. In a tissue with n cell types, the rarest present in $m\%$, how many cells X do you need to sequence so that you have at least Y of the rarest subtype.

[00630] For example, ~ 500 cells ensures $\geq 95\%$ chance of including ≥ 10 of each type, based on the following calculation (Figure 34). Assume the most conservative scenario that of M cell subtypes (for example, 12), all but one having the lowest predicted proportion (for example, $p_{min}=5\%$). Assuming that the Central Limit Theorem holds (a reasonable assumption when

solving to detect at least 10 cells of each type) the number of cells of each type i , termed T_i , will distribute as $E[T_i] = N * p_{min}$, $STDV[T_i] = \sqrt{(N * p_{min} * (1 - p_{min}))}$. The minimal N (total number of cells to profile) can be solved such that all $(m-1)$ subtypes have at least n cells (the last, majority, subtype easily clears this threshold since its proportion is much higher). Applicants confirmed with simulation that the strategy conservatively holds in practice even for $n < 10$, and take a margin of additional (conservative) error, to allow for subsequent failed RNA-Seq experiments (<20-30%, depending on protocol).

[00631] Amplification may involve thermocycling or isothermal amplification (such as through the methods RPA or LAMP). Cross-linking may involve overlap-extension PCR or use of ligase to associate multiple amplification products with each other.

[00632] For purpose of this invention, amplification means any method employing a primer and a polymerase capable of replicating a target sequence with reasonable fidelity. Amplification may be carried out by natural or recombinant DNA polymerases such as TaqGold™, T7 DNA polymerase, Klenow fragment of E.coli DNA polymerase, and reverse transcriptase. A preferred amplification method is PCR. In particular, the isolated RNA can be subjected to a reverse transcription assay that is coupled with a quantitative polymerase chain reaction (RT-PCR) in order to quantify the expression level of a sequence associated with a signaling biochemical pathway.

[00633] Detection of the gene expression level can be conducted in real time in an amplification assay. In one aspect, the amplified products can be directly visualized with fluorescent DNA-binding agents including but not limited to DNA intercalators and DNA groove binders. Because the amount of the intercalators incorporated into the double-stranded DNA molecules is typically proportional to the amount of the amplified DNA products, one can conveniently determine the amount of the amplified products by quantifying the fluorescence of the intercalated dye using conventional optical systems in the art. DNA-binding dye suitable for this application include SYBR green, SYBR blue, DAPI, propidium iodine, Hoeste, SYBR gold, ethidium bromide, acridines, proflavine, acridine orange, acriflavine, fluorcoumanin, ellipticine, daunomycin, chloroquine, distamycin D, chromomycin, homidium, mithramycin, ruthenium polypyridyls, anthramycin, and the like.

[00634] In another aspect, other fluorescent labels such as sequence specific probes can be employed in the amplification reaction to facilitate the detection and quantification of the

amplified products. Probe-based quantitative amplification relies on the sequence-specific detection of a desired amplified product. It utilizes fluorescent, target-specific probes (e.g., TaqMan® probes) resulting in increased specificity and sensitivity. Methods for performing probe-based quantitative amplification are well established in the art and are taught in U.S. Patent No. 5,210,015.

[00635] Sequencing may be performed on any high-throughput platform with read-length (either single- or paired-end) sufficient to cover both template and cross-linking event UID's. Methods of sequencing oligonucleotides and nucleic acids are well known in the art (see, e.g., WO93/23564, WO98/28440 and WO98/13523; U.S. Pat. Nos. 5,525,464; 5,202,231; 5,695,940; 4,971,903; 5,902,723; 5,795,782; 5,547,839 and 5,403,708; Sanger et al., Proc. Natl. Acad. Sci. USA 74:5463 (1977); Drmanac et al., Genomics 4:114 (1989); Koster et al., Nature Biotechnology 14:1123 (1996); Hyman, Anal. Biochem. 174:423 (1988); Rosenthal, International Patent Application Publication 761107 (1989); Metzker et al., Nucl. Acids Res. 22:4259 (1994); Jones, Biotechniques 22:938 (1997); Ronaghi et al., Anal. Biochem. 242:84 (1996); Ronaghi et al., Science 281:363 (1998); Nyren et al., Anal. Biochem. 151:504 (1985); Canard and Arzumanov, Gene 11:1 (1994); Dyatkina and Arzumanov, Nucleic Acids Symp Ser 18:117 (1987); Johnson et al., Anal. Biochem. 136:192 (1984); and Elgen and Rigler, Proc. Natl. Acad. Sci. USA 91(13):5740 (1994), all of which are expressly incorporated by reference).

[00636] The present invention may be applied to (1) single-cell transcriptomics: cDNA synthesized from mRNA is barcoded and cross-linked during in situ amplification, (2) single-cell proteomics: cDNA or DNA synthesized from RNA- or DNA-tagged antibodies of one or multiple specificities maps the abundance and distributions of different protein-antigens and (3) whole-tissue transcriptomic/proteomic mapping (molecular microscopy or VIPUR microscopy): using the frequency of cross-contamination between cells to determine their physical proximity, and via applications (1) single-cell transcriptomics and (2) single-cell proteomics, determining the global spatial distribution of mRNA, protein, or other biomolecules in a biological sample. This may be used, for example, to screen for anti-cancer/pathogen immunoglobulins (by analyzing co-localization of B-cells and T-cells within affected tissue) for immunotherapy.

[00637] The terms “polynucleotide”, “nucleotide”, “nucleotide sequence”, “nucleic acid” and “oligonucleotide” are used interchangeably. They refer to a polymeric form of nucleotides of any length, either deoxyribonucleotides or ribonucleotides, or analogs thereof. Polynucleotides may

have any three dimensional structure, and may perform any function, known or unknown. The following are non-limiting examples of polynucleotides: coding or non-coding regions of a gene or gene fragment, loci (locus) defined from linkage analysis, exons, introns, messenger RNA (mRNA), transfer RNA, ribosomal RNA, short interfering RNA (siRNA), short-hairpin RNA (shRNA), micro-RNA (miRNA), ribozymes, cDNA, recombinant polynucleotides, branched polynucleotides, plasmids, vectors, isolated DNA of any sequence, isolated RNA of any sequence, nucleic acid probes, and primers. The term also encompasses nucleic-acid-like structures with synthetic backbones, see, e.g., Eckstein, 1991; Baserga et al., 1992; Milligan, 1993; WO 97/03211; WO 96/39154; Mata, 1997; Strauss-Soukup, 1997; and Samstag, 1996. A polynucleotide may comprise one or more modified nucleotides, such as methylated nucleotides and nucleotide analogs. If present, modifications to the nucleotide structure may be imparted before or after assembly of the polymer. The sequence of nucleotides may be interrupted by non-nucleotide components. A polynucleotide may be further modified after polymerization, such as by conjugation with a labeling component.

[00638] As used herein the term “variant” should be taken to mean the exhibition of qualities that differ, such as, but not limited to, genetic variations including SNPs, insertion deletion events, and the like.

[00639] The terms “non-naturally occurring” or “engineered” are used interchangeably and indicate the involvement of the hand of man. The terms, when referring to nucleic acid molecules or polypeptides mean that the nucleic acid molecule or the polypeptide is at least substantially free from at least one other component with which they are naturally associated in nature and as found in nature.

[00640] “Complementarity” refers to the ability of a nucleic acid to form hydrogen bond(s) with another nucleic acid sequence by either traditional Watson-Crick or other non-traditional types. A percent complementarity indicates the percentage of residues in a nucleic acid molecule which can form hydrogen bonds (e.g., Watson-Crick base pairing) with a second nucleic acid sequence (e.g., 5, 6, 7, 8, 9, 10 out of 10 being 50%, 60%, 70%, 80%, 90%, and 100% complementary). “Perfectly complementary” means that all the contiguous residues of a nucleic acid sequence hydrogen bond with the same number of contiguous residues in a second nucleic acid sequence. “Substantially complementary” as used herein refers to a degree of complementarity that is at least 60%, 65%, 70%, 75%, 80%, 85%, 90%, 95%, 97%, 98%, 99%,

or 100% over a region of 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 30, 35, 40, 45, 50, or more nucleotides, or refers to two nucleic acids that hybridize under stringent conditions.

[00641] As used herein, “stringent conditions” for hybridization refer to conditions under which a nucleic acid having complementarity to a target sequence predominantly hybridizes with the target sequence, and substantially does not hybridize to non-target sequences. Stringent conditions are generally sequence-dependent, and vary depending on a number of factors. In general, the longer the sequence, the higher the temperature at which the sequence specifically hybridizes to its target sequence. Non-limiting examples of stringent conditions are described in detail in Tijssen (1993), *Laboratory Techniques In Biochemistry And Molecular Biology-Hybridization With Nucleic Acid Probes Part I, Second Chapter “Overview of principles of hybridization and the strategy of nucleic acid probe assay”*, Elsevier, N.Y.

[00642] “Hybridization” refers to a reaction in which one or more polynucleotides react to form a complex that is stabilized via hydrogen bonding between the bases of the nucleotide residues. The hydrogen bonding may occur by Watson Crick base pairing, Hoogstein binding, or in any other sequence specific manner. The complex may comprise two strands forming a duplex structure, three or more strands forming a multi stranded complex, a single self-hybridizing strand, or any combination of these. A hybridization reaction may constitute a step in a more extensive process, such as the initiation of PCR, or the cleavage of a polynucleotide by an enzyme. A sequence capable of hybridizing with a given sequence is referred to as the “complement” of the given sequence.

[00643] As used herein, the term “genomic locus” or “locus” (plural loci) is the specific location of a gene or DNA sequence on a chromosome. A “gene” refers to stretches of DNA or RNA that encode a polypeptide or an RNA chain that has functional role to play in an organism and hence is the molecular unit of heredity in living organisms. For the purpose of this invention it may be considered that genes include regions which regulate the production of the gene product, whether or not such regulatory sequences are adjacent to coding and/or transcribed sequences. Accordingly, a gene includes, but is not necessarily limited to, promoter sequences, terminators, translational regulatory sequences such as ribosome binding sites and internal ribosome entry sites, enhancers, silencers, insulators, boundary elements, replication origins, matrix attachment sites and locus control regions.

[00644] As used herein, “expression of a genomic locus” or “gene expression” is the process by which information from a gene is used in the synthesis of a functional gene product. The products of gene expression are often proteins, but in non-protein coding genes such as rRNA genes or tRNA genes, the product is functional RNA. The process of gene expression is used by all known life - eukaryotes (including multicellular organisms), prokaryotes (bacteria and archaea) and viruses to generate functional products to survive. As used herein “expression” of a gene or nucleic acid encompasses not only cellular gene expression, but also the transcription and translation of nucleic acid(s) in cloning systems and in any other context. As used herein, “expression” also refers to the process by which a polynucleotide is transcribed from a DNA template (such as into mRNA or other RNA transcript) and/or the process by which a transcribed mRNA is subsequently translated into peptides, polypeptides, or proteins. Transcripts and encoded polypeptides may be collectively referred to as “gene product.” If the polynucleotide is derived from genomic DNA, expression may include splicing of the mRNA in a eukaryotic cell.

[00645] As described in aspects of the invention, sequence identity is related to sequence homology. Homology comparisons may be conducted by eye, or more usually, with the aid of readily available sequence comparison programs. These commercially available computer programs may calculate percent (%) homology between two or more sequences and may also calculate the sequence identity shared by two or more amino acid or nucleic acid sequences.

[00646] Sequence homologies may be generated by any of a number of computer programs known in the art, for example BLAST or FASTA, etc. A suitable computer program for carrying out such an alignment is the GCG Wisconsin Bestfit package (University of Wisconsin, U.S.A; Devereux et al., 1984, *Nucleic Acids Research* 12:387). Examples of other software that may perform sequence comparisons include, but are not limited to, the BLAST package (see Ausubel et al., 1999 *ibid* – Chapter 18), FASTA (Atschul et al., 1990, *J. Mol. Biol.*, 403-410) and the GENWORKS suite of comparison tools. Both BLAST and FASTA are available for offline and online searching (see Ausubel et al., 1999 *ibid*, pages 7-58 to 7-60). However it is preferred to use the GCG Bestfit program. % homology may be calculated over contiguous sequences, i.e., one sequence is aligned with the other sequence and each amino acid or nucleotide in one sequence is directly compared with the corresponding amino acid or nucleotide in the other sequence, one residue at a time. This is called an “ungapped” alignment. Typically, such

ungapped alignments are performed only over a relatively short number of residues. Although this is a very simple and consistent method, it fails to take into consideration that, for example, in an otherwise identical pair of sequences, one insertion or deletion may cause the following amino acid residues to be put out of alignment, thus potentially resulting in a large reduction in % homology when a global alignment is performed. Consequently, most sequence comparison methods are designed to produce optimal alignments that take into consideration possible insertions and deletions without unduly penalizing the overall homology or identity score. This is achieved by inserting “gaps” in the sequence alignment to try to maximize local homology or identity. However, these more complex methods assign “gap penalties” to each gap that occurs in the alignment so that, for the same number of identical amino acids, a sequence alignment with as few gaps as possible - reflecting higher relatedness between the two compared sequences - may achieve a higher score than one with many gaps. “Affinity gap costs” are typically used that charge a relatively high cost for the existence of a gap and a smaller penalty for each subsequent residue in the gap. This is the most commonly used gap scoring system. High gap penalties may, of course, produce optimized alignments with fewer gaps. Most alignment programs allow the gap penalties to be modified. However, it is preferred to use the default values when using such software for sequence comparisons. For example, when using the GCG Wisconsin Bestfit package the default gap penalty for amino acid sequences is -12 for a gap and -4 for each extension. Calculation of maximum % homology therefore first requires the production of an optimal alignment, taking into consideration gap penalties. A suitable computer program for carrying out such an alignment is the GCG Wisconsin Bestfit package (Devereux et al., 1984 *Nuc. Acids Research* 12 p387). Examples of other software that may perform sequence comparisons include, but are not limited to, the BLAST package (see Ausubel et al., 1999 *Short Protocols in Molecular Biology*, 4th Ed. – Chapter 18), FASTA (Altschul et al., 1990 *J. Mol. Biol.* 403-410) and the GENWORKS suite of comparison tools. Both BLAST and FASTA are available for offline and online searching (see Ausubel et al., 1999, *Short Protocols in Molecular Biology*, pages 7-58 to 7-60). However, for some applications, it is preferred to use the GCG Bestfit program. A new tool, called BLAST 2 Sequences is also available for comparing protein and nucleotide sequences (see *FEMS Microbiol Lett.* 1999 174(2): 247-50; *FEMS Microbiol Lett.* 1999 177(1): 187-8 and the website of the National Center for Biotechnology information at the website of the National Institutes for Health). Although the final % homology may be

measured in terms of identity, the alignment process itself is typically not based on an all-or-nothing pair comparison. Instead, a scaled similarity score matrix is generally used that assigns scores to each pair-wise comparison based on chemical similarity or evolutionary distance. An example of such a matrix commonly used is the BLOSUM62 matrix - the default matrix for the BLAST suite of programs. GCG Wisconsin programs generally use either the public default values or a custom symbol comparison table, if supplied (see user manual for further details). For some applications, it is preferred to use the public default values for the GCG package, or in the case of other software, the default matrix, such as BLOSUM62.

[00647] Alternatively, percentage homologies may be calculated using the multiple alignment feature in DNASISTM (Hitachi Software), based on an algorithm, analogous to CLUSTAL (Higgins DG & Sharp PM (1988), *Gene* 73(1), 237-244). Once the software has produced an optimal alignment, it is possible to calculate % homology, preferably % sequence identity. The software typically does this as part of the sequence comparison and generates a numerical result.

[00648] Embodiments of the invention include sequences (both polynucleotide or polypeptide) which may comprise homologous substitution (substitution and replacement are both used herein to mean the interchange of an existing amino acid residue or nucleotide, with an alternative residue or nucleotide) that may occur i.e., like-for-like substitution in the case of amino acids such as basic for basic, acidic for acidic, polar for polar, etc. Non-homologous substitution may also occur i.e., from one class of residue to another or alternatively involving the inclusion of unnatural amino acids such as ornithine (hereinafter referred to as Z), diaminobutyric acid ornithine (hereinafter referred to as B), norleucine ornithine (hereinafter referred to as O), pyriylalanine, thienylalanine, naphthylalanine and phenylglycine.

[00649] Hybridization can be performed under conditions of various stringency. Suitable hybridization conditions for the practice of the present invention are such that the recognition interaction between the probe and sequences associated with a signaling biochemical pathway is both sufficiently specific and sufficiently stable. Conditions that increase the stringency of a hybridization reaction are widely known and published in the art. See, for example, (Sambrook, et al., (1989); *Nonradioactive In Situ Hybridization Application Manual*, Boehringer Mannheim, second edition). The hybridization assay can be formed using probes immobilized on any solid support, including but are not limited to nitrocellulose, glass, silicon, and a variety of gene

arrays. A preferred hybridization assay is conducted on high-density gene chips as described in U.S. Patent No. 5,445,934.

[00650] For a convenient detection of the probe-target complexes formed during the hybridization assay, the nucleotide probes are conjugated to a detectable label. Detectable labels suitable for use in the present invention include any composition detectable by photochemical, biochemical, spectroscopic, immunochemical, electrical, optical or chemical means. A wide variety of appropriate detectable labels are known in the art, which include fluorescent or chemiluminescent labels, radioactive isotope labels, enzymatic or other ligands. In preferred embodiments, one will likely desire to employ a fluorescent label or an enzyme tag, such as digoxigenin, β -galactosidase, urease, alkaline phosphatase or peroxidase, avidin/biotin complex.

[00651] The detection methods used to detect or quantify the hybridization intensity will typically depend upon the label selected above. For example, radiolabels may be detected using photographic film or a phosphorimager. Fluorescent markers may be detected and quantified using a photodetector to detect emitted light. Enzymatic labels are typically detected by providing the enzyme with a substrate and measuring the reaction product produced by the action of the enzyme on the substrate; and finally colorimetric labels are detected by simply visualizing the colored label.

[00652] Examples of the labeling substance which may be employed include labeling substances known to those skilled in the art, such as fluorescent dyes, enzymes, coenzymes, chemiluminescent substances, and radioactive substances. Specific examples include radioisotopes (e.g., ^{32}P , ^{14}C , ^{125}I , ^3H , and ^{131}I), fluorescein, rhodamine, dansyl chloride, umbelliferone, luciferase, peroxidase, alkaline phosphatase, β -galactosidase, β -glucosidase, horseradish peroxidase, glucoamylase, lysozyme, saccharide oxidase, microperoxidase, biotin, and ruthenium. In the case where biotin is employed as a labeling substance, preferably, after addition of a biotin-labeled antibody, streptavidin bound to an enzyme (e.g., peroxidase) is further added.

[00653] Advantageously, the label is a fluorescent label. Examples of fluorescent labels include, but are not limited to, Atto dyes, 4-acetamido-4'-isothiocyanatostilbene-2,2'-disulfonic acid; acridine and derivatives: acridine, acridine isothiocyanate; 5-(2'-aminoethyl)aminonaphthalene-1-sulfonic acid (EDANS); 4-amino-N-[3-vinylsulfonyl]phenyl]naphthalimide-3,5 disulfonate; N-(4-anilino-1-naphthyl)maleimide;

anthranilamide; BODIPY; Brilliant Yellow; coumarin and derivatives; coumarin, 7-amino-4-methylcoumarin (AMC, Coumarin 120), 7-amino-4-trifluoromethylcoumarin (Coumarin 151); cyanine dyes; cyanosine; 4',6-diaminidino-2-phenylindole (DAPI); 5''-dibromopyrogallol-sulfonaphthalein (Bromopyrogallol Red); 7-diethylamino-3-(4'-isothiocyanatophenyl)-4-methylcoumarin; diethylenetriamine pentaacetate; 4,4'-diisothiocyanatodihydro-stilbene-2,2'-disulfonic acid; 4,4'-diisothiocyanatostilbene-2,2'-disulfonic acid; 5-[dimethylamino]naphthalene-1-sulfonyl chloride (DNS, dansylchloride); 4-dimethylaminophenylazophenyl-4'-isothiocyanate (DABITC); eosin and derivatives; eosin, eosin isothiocyanate, erythrosin and derivatives; erythrosin B, erythrosin, isothiocyanate; ethidium; fluorescein and derivatives; 5-carboxyfluorescein (FAM), 5-(4,6-dichlorotriazin-2-yl)aminofluorescein (DTAF), 2',7'-dimethoxy-4'5'-dichloro-6-carboxyfluorescein, fluorescein, fluorescein isothiocyanate, QFITC, (XRITC); fluorescamine; IR144; IR1446; Malachite Green isothiocyanate; 4-methylumbelliferoneortho cresolphthalein; nitrotyrosine; pararosaniline; Phenol Red; B-phycoerythrin; o-phthalaldehyde; pyrene and derivatives: pyrene, pyrene butyrate, succinimidyl 1-pyrene; butyrate quantum dots; Reactive Red 4 (Cibacron.TM. Brilliant Red 3B-A) rhodamine and derivatives: 6-carboxy-X-rhodamine (ROX), 6-carboxyrhodamine (R6G), lissamine rhodamine B sulfonyl chloride rhodamine (Rhod), rhodamine B, rhodamine 123, rhodamine X isothiocyanate, sulforhodamine B, sulforhodamine 101, sulfonyl chloride derivative of sulforhodamine 101 (Texas Red); N,N,N',N' tetramethyl-6-carboxyrhodamine (TAMRA); tetramethyl rhodamine; tetramethyl rhodamine isothiocyanate (TRITC); riboflavin; rosolic acid; terbium chelate derivatives; Cy3; Cy5; Cy5.5; Cy7; IRD 700; IRD 800; La Jolla Blue; phthalo cyanine; and naphthalo cyanine

[00654] The fluorescent label may be a fluorescent protein, such as blue fluorescent protein, cyan fluorescent protein, green fluorescent protein, red fluorescent protein, yellow fluorescent protein or any photoconvertible protein. Colorimetric labeling, bioluminescent labeling and/or chemiluminescent labeling may further accomplish labeling. Labeling further may include energy transfer between molecules in the hybridization complex by perturbation analysis, quenching, or electron transport between donor and acceptor molecules, the latter of which may be facilitated by double stranded match hybridization complexes. The fluorescent label may be a perylene or a terrylen. In the alternative, the fluorescent label may be a fluorescent bar code.

[00655] In an advantageous embodiment, the label may be light sensitive, wherein the label is light-activated and/or light cleaves the one or more linkers to release the molecular cargo. The light-activated molecular cargo may be a major light-harvesting complex (LHCII). In another embodiment, the fluorescent label may induce free radical formation.

[00656] In an advantageous embodiment, agents may be uniquely labeled in a dynamic manner (see, e.g., international patent application serial no. PCT/US2013/61182 filed September 23, 2012). The unique labels are, at least in part, nucleic acid in nature, and may be generated by sequentially attaching two or more detectable oligonucleotide tags to each other and each unique label may be associated with a separate agent. A detectable oligonucleotide tag may be an oligonucleotide that may be detected by sequencing of its nucleotide sequence and/or by detecting non-nucleic acid detectable moieties to which it may be attached.

[00657] The oligonucleotide tags may be detectable by virtue of their nucleotide sequence, or by virtue of a non-nucleic acid detectable moiety that is attached to the oligonucleotide such as but not limited to a fluorophore, or by virtue of a combination of their nucleotide sequence and the nonnucleic acid detectable moiety.

[00658] In some embodiments, a detectable oligonucleotide tag may comprise one or more nonoligonucleotide detectable moieties. Examples of detectable moieties may include, but are not limited to, fluorophores, microparticles including quantum dots (Empodocles, et al., Nature 399:126-130, 1999), gold nanoparticles (Reichert et al., Anal. Chem. 72:6025-6029, 2000), biotin, DNP (dinitrophenyl), fucose, digoxigenin, haptens, and other detectable moieties known to those skilled in the art. In some embodiments, the detectable moieties may be quantum dots. Methods for detecting such moieties are described herein and/or are known in the art.

[00659] Thus, detectable oligonucleotide tags may be, but are not limited to, oligonucleotides which may comprise unique nucleotide sequences, oligonucleotides which may comprise detectable moieties, and oligonucleotides which may comprise both unique nucleotide sequences and detectable moieties.

[00660] A unique label may be produced by sequentially attaching two or more detectable oligonucleotide tags to each other. The detectable tags may be present or provided in a plurality of detectable tags. The same or a different plurality of tags may be used as the source of each detectable tag may be part of a unique label. In other words, a plurality of tags may be subdivided into subsets and single subsets may be used as the source for each tag.

[00661] In some embodiments, a detectable oligonucleotide tag may comprise one or more non-oligonucleotide detectable moieties. Examples of detectable moieties include, but are not limited to, fluorophores, microparticles including quantum dots (Empodocles, et al., *Nature* 399:126-130, 1999), gold nanoparticles (Reichert et al., *Anal. Chem.* 72:6025-6029, 2000), biotin, DNP (dinitrophenyl), fucose, digoxigenin, haptens, and other detectable moieties known to those skilled in the art. In some embodiments, the detectable moieties are quantum dots. Methods for detecting such moieties are described herein and/or are known in the art.

[00662] Thus, detectable oligonucleotide tags may be, but are not limited to, oligonucleotides which may comprise unique nucleotide sequences, oligonucleotides which may comprise detectable moieties, and oligonucleotides which may comprise both unique nucleotide sequences and detectable moieties.

[00663] A unique nucleotide sequence may be a nucleotide sequence that is different (and thus distinguishable) from the sequence of each detectable oligonucleotide tag in a plurality of detectable oligonucleotide tags. A unique nucleotide sequence may also be a nucleotide sequence that is different (and thus distinguishable) from the sequence of each detectable oligonucleotide tag in a first plurality of detectable oligonucleotide tags but identical to the sequence of at least one detectable oligonucleotide tag in a second plurality of detectable oligonucleotide tags. A unique sequence may differ from other sequences by multiple bases (or base pairs). The multiple bases may be contiguous or non-contiguous. Methods for obtaining nucleotide sequences (e.g., sequencing methods) are described herein and/or are known in the art.

[00664] In some embodiments, detectable oligonucleotide tags comprise one or more of a ligation sequence, a priming sequence, a capture sequence, and a unique sequence (optionally referred to herein as an index sequence). A ligation sequence is a sequence complementary to a second nucleotide sequence which allows for ligation of the detectable oligonucleotide tag to another entity which may comprise the second nucleotide sequence, e.g., another detectable oligonucleotide tag or an oligonucleotide adapter. A priming sequence is a sequence complementary to a primer, e.g., an oligonucleotide primer used for an amplification reaction such as but not limited to PCR. A capture sequence is a sequence capable of being bound by a capture entity. A capture entity may be an oligonucleotide which may comprise a nucleotide sequence complementary to a capture sequence, e.g. a second detectable oligonucleotide tag. A capture entity may also be any other entity capable of binding to the capture sequence, e.g. an

antibody, hapten or peptide. An index sequence is a sequence which may comprise a unique nucleotide sequence and/or a detectable moiety as described above.

[00665] “Complementary” is a term which is used to indicate a sufficient degree of complementarity between two nucleotide sequences such that stable and specific binding occurs between one and preferably more bases (or nucleotides, as the terms are used interchangeably herein) of the two sequences. For example, if a nucleotide in a first nucleotide sequence is capable of hydrogen bonding with a nucleotide in second nucleotide sequence, then the bases are considered to be complementary to each other. Complete (i.e., 100%) complementarity between a first nucleotide sequence and a second nucleotide is preferable, but not required for ligation, priming, or capture sequences.

[00666] The present invention also relates to a computer system involved in carrying out the methods of the invention relating to both computations and sequencing.

[00667] A computer system (or digital device) may be used to receive, transmit, display and/or store results, analyze the results, and/or produce a report of the results and analysis. A computer system may be understood as a logical apparatus that can read instructions from media (e.g. software) and/or network port (e.g. from the internet), which can optionally be connected to a server having fixed media. A computer system may comprise one or more of a CPU, disk drives, input devices such as keyboard and/or mouse, and a display (e.g. a monitor). Data communication, such as transmission of instructions or reports, can be achieved through a communication medium to a server at a local or a remote location. The communication medium can include any means of transmitting and/or receiving data. For example, the communication medium can be a network connection, a wireless connection, or an internet connection. Such a connection can provide for communication over the World Wide Web. It is envisioned that data relating to the present invention can be transmitted over such networks or connections (or any other suitable means for transmitting information, including but not limited to mailing a physical report, such as a print-out) for reception and/or for review by a receiver. The receiver can be but is not limited to an individual, or electronic system (e.g. one or more computers, and/or one or more servers).

[00668] In some embodiments, the computer system may comprise one or more processors. Processors may be associated with one or more controllers, calculation units, and/or other units of a computer system, or implanted in firmware as desired. If implemented in software, the

routines may be stored in any computer readable memory such as in RAM, ROM, flash memory, a magnetic disk, a laser disk, or other suitable storage medium. Likewise, this software may be delivered to a computing device via any known delivery method including, for example, over a communication channel such as a telephone line, the internet, a wireless connection, etc., or via a transportable medium, such as a computer readable disk, flash drive, etc. The various steps may be implemented as various blocks, operations, tools, modules and techniques which, in turn, may be implemented in hardware, firmware, software, or any combination of hardware, firmware, and/or software. When implemented in hardware, some or all of the blocks, operations, techniques, etc. may be implemented in, for example, a custom integrated circuit (IC), an application specific integrated circuit (ASIC), a field programmable logic array (FPGA), a programmable logic array (PLA), etc.

[00669] A client-server, relational database architecture can be used in embodiments of the invention. A client-server architecture is a network architecture in which each computer or process on the network is either a client or a server. Server computers are typically powerful computers dedicated to managing disk drives (file servers), printers (print servers), or network traffic (network servers). Client computers include PCs (personal computers) or workstations on which users run applications, as well as example output devices as disclosed herein. Client computers rely on server computers for resources, such as files, devices, and even processing power. In some embodiments of the invention, the server computer handles all of the database functionality. The client computer can have software that handles all the front-end data management and can also receive data input from users.

[00670] A machine readable medium which may comprise computer-executable code may take many forms, including but not limited to, a tangible storage medium, a carrier wave medium or physical transmission medium. Non-volatile storage media include, for example, optical or magnetic disks, such as any of the storage devices in any computer(s) or the like, such as may be used to implement the databases, etc. shown in the drawings. Volatile storage media include dynamic memory, such as main memory of such a computer platform. Tangible transmission media include coaxial cables; copper wire and fiber optics, including the wires that comprise a bus within a computer system. Carrier-wave transmission media may take the form of electric or electromagnetic signals, or acoustic or light waves such as those generated during radio frequency (RF) and infrared (IR) data communications. Common forms of computer-readable

media therefore include for example: a floppy disk, a flexible disk, hard disk, magnetic tape, any other magnetic medium, a CD-ROM, DVD or DVD-ROM, any other optical medium, punch cards paper tape, any other physical storage medium with patterns of holes, a RAM, a ROM, a PROM and EPROM, a FLASH-EPROM, any other memory chip or cartridge, a carrier wave transporting data or instructions, cables or links transporting such a carrier wave, or any other medium from which a computer may read programming code and/or data. Many of these forms of computer readable media may be involved in carrying one or more sequences of one or more instructions to a processor for execution.

[00671] The subject computer-executable code can be executed on any suitable device which may comprise a processor, including a server, a PC, or a mobile device such as a smartphone or tablet. Any controller or computer optionally includes a monitor, which can be a cathode ray tube (“CRT”) display, a flat panel display (e.g., active matrix liquid crystal display, liquid crystal display, etc.), or others. Computer circuitry is often placed in a box, which includes numerous integrated circuit chips, such as a microprocessor, memory, interface circuits, and others. The box also optionally includes a hard disk drive, a floppy disk drive, a high capacity removable drive such as a writeable CD-ROM, and other common peripheral elements. Inputting devices such as a keyboard, mouse, or touch-sensitive screen, optionally provide for input from a user. The computer can include appropriate software for receiving user instructions, either in the form of user input into a set of parameter fields, e.g., in a GUI, or in the form of preprogrammed instructions, e.g., preprogrammed for a variety of different specific operations.

[00672] The present invention also contemplates multiplex assays. The present invention is especially well suited for multiplex assays. For example, the invention encompasses use of a SureSelect^{XT}, SureSelect^{XT2} and SureSelect^{QXT} Target Enrichment System for Illumina Multiplexed Sequencing developed by Agilent Technologies (see, e.g., www.agilent.com/genomics/protocolvideos), a SeqCap EZ kit developed by Roche NimbleGen, a TruSeq® Enrichment Kit developed by Illumina and other hybridization-based target enrichment methods and kits that add sample-specific sequence tags either before or after the enrichment step, as well as Illumina HiSeq, MiSeq and NexSeq, Life Technology Ion Torrent, Pacific Biosciences PacBio RSII, Oxford Nanopore MinIon, Promethion and GridIon and other massively parallel Multiplexed Sequencing Platforms.

[00673] Although the present invention and its advantages have been described in detail, it should be understood that various changes, substitutions and alterations can be made herein without departing from the spirit and scope of the invention as defined in the appended claims.

[00674] The present invention will be further illustrated in the following Examples which are given for illustration purposes only and are not intended to limit the invention in any way.

[00675] **PROTOCOL**

[00676] Description of Basic Protocol

[00677] *Probe library generation*

[00678] Each probe library has two groups of probes. One group is called “left probe,” and the other “right probe.” Each of the probes bind a short sequence on the target transcript (or cDNA, or target molecule) and the two binding sites are exactly right next to each other (juxtaposition probe set), so that once bound to target they can be ligated to given a single ligation product.

[00679] Applicants synthesized the two groups probe separately, and one of the probes will have the measurement barcode, the random molecular index (UMI), and optionally the cell barcodes for single-cell experiments, as shown in Figure 13.

[00680] The probe design in Figure 13 is just one group of the probe, in the figure it will be the left side probe, given the 5'-3' direction. Applicants generate the probe through performing Klenow extension synthesis followed by lambda exonuclease digestion. In detail, Oligo 1 is first synthesized, containing the Universal adapter, Cell label, UMI, SMRI, and a LK (linker) region, followed by the synthesis of Oligo 2, the complementary sequence.

[00681] The left probe combined with the right side probe will be able to bind specifically to the designed location in the target transcript. Applicants then ligate the two probes to form amplification-ready product. This is followed by amplification to add on sequencing adapter. The amplified probes are then purified, quantified, and subjected to next-generation sequencing and analysis, as shown in Figure 14.

[00682] *Detailed Experimental Implementation for Bulk RNA*

[00683] 1. Extract RNA from K562 cells (K562 is a human immortalised myelogenous leukemia cell line)

[00684] 2. Prepare cDNA by synthesis, using the RNA as template

[00685] 3. Split part for qPCR and RNA-seq (as a reference control)

[00686] 4. Split part for probe binding/ligation in the compressive sensing experiment

[00687] 5. Synthesize Left and Right Probe and Generate Probe Library Mix

[00688] (1) Synthesize Oligo 1 with the adapter and barcode, and linker region

[00689] (2) Synthesize Oligo 2 with the complementary sequence to the linker region and gene-specific probe

[00690] (3) Anneal Oligo 1 and Oligo 2, and perform Klenow extension to synthesize the complete left side probe. The left side of the probe is either purified by column or used directly without further purification.

[00691] (4) Synthesize Oligo 3 with the right side gene-specific probe sequence, and then phosphorylate the right side probe with T4 Polynucleotide Kinase to phosphorylate the 5' end of the oligo to render it ligation-ready

[00692] (5) Mix probes according to the design matrix to generate a probe library mix with both left and right probe.

[00693] 6. Hybridize the probe library to the cDNA sample using a slow ramping protocol, where the mix is incubated at 95°C for 5 minutes, then gradually cooled down to room temperature (25°C). The annealed mixture containing left and right probes bound to the cDNA molecules are then ligated using Taq ligase to generate the ligation products at 45°C for 2 hours, and the enzymatic activity is heat-inactivated at 65°C for 20 minutes.

[00694] A control reaction where the probe library mixture is directly ligated with T7 DNA ligase without the addition of cDNA is also performed to obtain the abundance of the probe sets within the library. The ligation step is performed in a similar set-up except cDNA input is substituted with water, and ligation reaction is carried out at room temperature (25°C) for 2 hours.

[00695] 7. Amplify the ligation product using the library-amplification primer, which will also add the next-generation sequencer adapter to complete sequencing-ready library construction. The amplified is purified and size-selected between 200-300 bp and quantified.

[00696] *Experimental Design for Implementation for Single Cell Sample*

[00697] Each cell is captured and encapsulated in hydrogel droplet, which contains one side of the probe. These single-cell containing droplets are then lysed to allow initial hybridization, barcoded uniquely. Then, the other side of the probe will be added to the system to allow hybridization, ligation, and amplification. The final product are then purified and sequenced (or quantified/detected in other means), and subject to data analysis.

From genetic interactions to molecular interactions

[00698] Analyzing and determining genetic interactions can be used for dissection cellular circuitry. Results obtained using the method of the invention may be combined with further methods in order to probe molecular pathways. This may be achieved by using further perturbation assays (e.g. for validation purposes), by combining with complementary experiments in animal models (e.g. a Cas9 mouse or cell lines derived therefrom), including further genetic, biochemical or cellular testing. The systematic, sparsity-based approach may also be combined with 'knowledge-based' (informed) approaches, based on previously identified molecular interactions, such as known components of a given metabolic pathway.

[00699] The method of the invention may also be used to identify groups of cells based on their molecular profiling. The method of the invention may be used to identify possible target for therapeutics efforts.

[00700] The method of the invention may also be used to dissect cellular circuitry in relation with a spatial and/or temporal approach. For examples, the cells in the population may be synchronized with respect to their cell cycling.

Screening Methods

[00701] The present invention also envisions screening methods involving the herein described embodiments.

[00702] In one embodiment, the screening involves a a pooled, loss-of-function genetic screening approach suitable for both positive and negative selection that uses a genome-scale lentiviral single-guide RNA (sgRNA) library (see, e.g., Wang et al., Science. 2014 Jan 3;343(6166):80-4. doi: 10.1126/science.1246981. Epub 2013 Dec 12). Briefly, sgRNA expression cassettes were stably integrated into the genome, which enabled a complex mutant pool to be tracked by massively parallel sequencing. A library containing 73,000 sgRNAs was used to generate knockout collections and performed screens in two human cell lines. A screen for resistance to the nucleotide analog 6-thioguanine identified all expected members of the DNA mismatch repair pathway, whereas another for the DNA topoisomerase II (TOP2A) poison etoposide identified TOP2A, as expected, and also cyclin-dependent kinase 6, CDK6. A negative selection screen for essential genes identified numerous gene sets corresponding to fundamental processes. sgRNA efficiency is associated with specific sequence motifs, enabling the prediction

of more effective sgRNAs. See also Chen et al., Genome-wide CRISPR Screen in a Mouse Model of Tumor Growth and Metastasis, *Cell* (2015). DOI: 10.1016/j.cell.2015.02.038.

[00703] The activator screen method of Konermann et al., *Nature* (2014) doi:10.1038/nature14136 may be applied to the present invention. Systematic interrogation of gene function requires the ability to perturb gene expression in a robust and generalizable manner. Konermann et al. describes structure-guided engineering of a CRISPR-Cas9 complex to mediate efficient transcriptional activation at endogenous genomic loci. Konermann et al. used these engineered Cas9 activation complexes to investigate single-guide RNA (sgRNA) targeting rules for effective transcriptional activation, to demonstrate multiplexed activation of ten genes simultaneously, and to upregulate long intergenic non-coding RNA (lincRNA) transcripts. Konermann et al. also synthesized a library consisting of 70,290 guides targeting all human RefSeq coding isoforms to screen for genes that, upon activation, confer resistance to a BRAF inhibitor. The top hits included genes previously shown to be able to confer resistance, and novel candidates were validated using individual sgRNA and complementary DNA overexpression. A gene signature based on the top screening hits correlated with a gene expression signature of BRAF inhibitor resistance in cell lines and patient-derived samples. These results collectively demonstrate the potential of Cas9-based activators as a powerful genetic perturbation technology.

[00704] The mouse of Platt et al., *Cell*. 2014 Oct 9;159(2):440-55. doi: 10.1016/j.cell.2014.09.014. Epub 2014 Sep 25 may also be contemplated in the present invention. Platt et al. established a Cre-dependent Cas9 knockin mouse and demonstrated in vivo as well as ex vivo genome editing using adeno-associated virus (AAV)-, lentivirus-, or particle-mediated delivery of guide RNA in neurons, immune cells, and endothelial cells. Using these mice, Platt et al. simultaneously modeled the dynamics of KRAS, p53, and LKB1, the top three significantly mutated genes in lung adenocarcinoma. Delivery of a single AAV vector in the lung generated loss-of-function mutations in p53 and Lkb1, as well as homology-directed repair-mediated Kras(G12D) mutations, leading to macroscopic tumors of adenocarcinoma pathology.

[00705] The mouse of Platt et al., *Cell*. 2014 Oct 9;159(2):440-55. doi: 10.1016/j.cell.2014.09.014. Epub 2014 Sep 25, may also be contemplated in the present invention. Platt et al. established a Cre-dependent Cas9 knockin mouse. Platt et al. demonstrated in vivo as well as ex vivo genome editing using adeno-associated virus (AAV)-, lentivirus-, or

particle-mediated delivery of guide RNA in neurons, immune cells, and endothelial cells. Using these mice, Platt et al. simultaneously modeled the dynamics of KRAS, p53, and LKB1, the top three significantly mutated genes in lung adenocarcinoma. Delivery of a single AAV vector in the lung generated loss-of-function mutations in p53 and Lkb1, as well as homology-directed repair-mediated Kras(G12D) mutations, leading to macroscopic tumors of adenocarcinoma pathology.

[00706] Genetic screens infer gene function in mammalian cells, but it has remained difficult to assay complex phenotypes – such as transcriptional profiles – in large screens. Applicants have developed Perturb-seq. Perturb-seq advantageously combines single cell RNA-seq and CRISPR based perturbations identified by unique polyadenylated barcodes to perform many such assays in a pooled experiment. Applicants demonstrate Perturb-seq by analyzing 200,000 cells across three screens: transcription factors involved in the response of dendritic cells to LPS, transcription factors bound in K562 cells, and cell cycle regulators in the same cell line. Perturb-seq accurately identifies individual gene targets, gene signatures, and cell states affected by each perturbation and their genetic interactions. The individual gene target predictions were validated by their known functions and ChIP-Seq. Applicants advantageously posit new functions for regulators affecting cell differentiation, the anti-viral response, and mitochondrial function during immune activation. Perturb-seq can be flexibly applied to diverse cell metadata, to customize design and scope of pooled genomic assays.

[00707] Although the present invention and its advantages have been described in detail, it should be understood that various changes, substitutions and alterations can be made herein without departing from the spirit and scope of the invention as defined in the appended claims.

[00708] The present invention will be further illustrated in the following Examples which are given for illustration purposes only and are not intended to limit the invention in any way.

Examples

Example 1: combinatorial genetic perturbations in mouse cells using CRISPR-Cas mediated multiple gene activation

[00709] A pooled perturbation assay is performed in mouse cells, using a CRISPR-Cas system, as described in Shalem O., Sanjana NE., Zhang F. High-throughput functional genomics using CRISPR-Cas9. Nat Rev Genet. May;16(5):299-311. (2015). doi: 10.1038/nrg3899. Epub

2015 Apr 9. The pooled assay is designed to reach MOI values of 3-4, corresponding to orders 3-4 of genetic perturbations, concerning a total of 7,000 target genomic loci.

[00710] An sgRNA library is designed as described in Sanjana, N. E., Shalem, O. & Zhang, F. Improved vectors and genome-wide libraries for CRISPR screening. *Nat Methods*. 11, 783-784, doi:10.1038/nmeth.3047 (2014); Konermann, S., Brigham, M. D., Trevino, A. E., Joung, J., Abudayyeh, O. O., Barcena, C., Hsu, P. D., Habib, N., Gootenberg, J. S., Nishimasu, H., Nureki, O. & Zhang, F. Genome-scale transcriptional activation by an engineered CRISPR-Cas9 complex. *Nature*. 517, 583-588, doi:10.1038/nature14136 (2015). PMID:4420636.

[00711] Effects of the genetic perturbations are determined in terms of single-cell RNA profiling. Said RNA profiling is carried out using droplet technology (Drop-Seq technique, as described in Macosko, E. Z., Basu, A., Satija, R., Nemesh, J., Goldman, M., Tirosh, I., Bialas, A. R., Kamitaki, N., Sanes, J. R., Weitz, D. A., Shalek, A. K., Regev, A. & McCarroll, S. A. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell*. 2015 May 21;161(5):1202-14. doi: 10.1016/j.cell.2015.05.002. PMID:4481139; or Klein, A. M., Mazutis, L., Akartuna, I., Tallapragada, N., Veres, A., Li, V., Peshkin, L., Weitz, D.A., Kirschner, M. W. Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell*. 2015 May 21;161(5):1187-201. doi: 10.1016/j.cell.2015.04.044. PMID: 4441768).

Example 2: combinatorial genetic perturbations in human cells using CRISPR-Cas mediated multiple gene knock-outs

[00712] A pooled perturbation assay is performed in a human using a CRISPR-Cas system, as described in Shalem O., Sanjana NE., Zhang F. High-throughput functional genomics using CRISPR-Cas9. *Nat Rev Genet*. May;16(5):299-311. (2015). doi: 10.1038/nrg3899. Epub 2015 Apr 9. The pooled assay is designed to reach MOI values of 2-3, corresponding to orders 2-3 of genetic perturbations regarding a total of 10,000 target genomic loci.

[00713] An sgRNA library is designed as described in Sanjana, N. E., Shalem, O. & Zhang, F. Improved vectors and genome-wide libraries for CRISPR screening. *Nat Methods*. 11, 783-784, doi:10.1038/nmeth.3047 (2014); Shalem, O., Sanjana, N. E., Hartenian, E., Shi, X., Scott, D. A., Mikkelsen, T. S., Heckl, D., Ebert, B. L., Root, D. E., Doench, J. G. & Zhang, F. Genome-scale CRISPR-Cas9 knockout screening in human cells. *Science*. 343, 84-87, doi:10.1126/science.1247005 (2014). PMID:4089965; Shalem, O., Sanjana, N. E. & Zhang, F.

High-throughput functional genomics using CRISPR-Cas9. *Nat Rev Genet.* 16, 299-311, doi:10.1038/nrg3899 (2015).

[00714] Effects of the genetic perturbations are determined in terms of single-cell RNA profiling. Said RNA profiling is carried out using droplet technology (Drop-Seq technique, as described in Macosko, E. Z., Basu, A., Satija, R., Nemesh, J., Goldman, M., Tirosh, I., Bialas, A. R., Kamitaki, N., Sanes, J. R., Weitz, D. A., Shalek, A. K., Regev, A. & McCarroll, S. A. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell.* 2015 May 21;161(5):1202-14. doi: 10.1016/j.cell.2015.05.002. PMID:4481139; or Klein, A. M., Mazutis, L., Akartuna, I., Tallapragada, N., Veres, A., Li, V., Peshkin, L., Weitz, D.A., Kirschner, M. W. Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell.* 2015 May 21;161(5):1187-201. doi: 10.1016/j.cell.2015.04.044. PMID: 4441768).

Example 3: Perturb-Seq: pooled CRISPR screens with single cell RNA -seq readout

[00715] Applicants devised Perturb-seq by combining the design of a pooled CRISPR screen, where the identity of each perturbation is reported by a unique barcode for each guide that is assayed by sequencing, with massively-parallel, droplet-based scRNA-seq, as a rich genomic phenotype for each cell (**Figure 48A**). In Perturb-seq, the sgRNA barcode is reported by a polyadenylated RNA, which is captured by single cell RNA-seq. Perturb-seq uses a CRISPR lentiviral vector that both delivers an sgRNA to a cell and reports on the identity of the delivered sgRNA by an expressed PolyAdenylated Perturbation Index (PAPI) (**Figure 48B**). As in a typical pooled genetic screen, applicants first infect cells with a pool of lentiviruses, each carrying a construct for a particular guide. Depending on the multiplicity of infection (MOI), applicants expect a different distribution of the number of lentivirus genomic insertions (and hence sgRNAs) in each cell. At low MOI, the vast majority of cells would have one or no guides. In such a case, Perturb-seq can be used to assess the effect of single gene perturbations. By increasing the multiplicity of infection (At higher MOI), a greater fraction of cells possesses more than one guide. Thus, the MOI is a tunable parameter that, for a given number of cells, determines how powered tests will be to study epistatic effects.

[00716] Next, following growth, differentiation, and/or stimulation, Applicants collect the cells and process them through massively parallel scRNA-seq. This process tags the mRNA of each cell, including the synthetic PAPI, with a unique cell barcode (CBC) and a unique

molecular identifier (UMI). The CBC associates the cell's transcriptional profile with the genetic perturbation(s) delivered to this cell, as encoded by the PAPI. Applicants perform Whole Transcriptome Amplification (WTA), library preparation, and sequencing in a single pool. In our current implementation, Applicants rely on a droplet method, which is now commercially available (Zheng et al., 2016) (**Methods and Resources**), but our design is compatible with additional single-cell RNA-seq methods (Fan et al., 2015; Klein et al., 2015; Macosko et al., 2015), and Applicants have tested it successfully with Drop-seq (Macosko et al., 2015) in both K562 cells and BMDCs, albeit with different gene targets than in the rest of this study (AD, OP, BL, and AR, unpublished data).

[00717] Applicants demonstrate Perturb-seq, in two complementary biological systems: K562 cells, a rapidly growing human erythroleukemia cell line, and bone marrow derived mouse dendritic cells (BMDCs), a post-mitotic, short-lived primary cell, which Applicants analyzed pre-stimulation (0h) and three hours following stimulation with LPS. Applicants have previously established both systems for genome-wide CRISPR screens, with a readout of either cell viability (in K562 cells (Wang et al., 2015)), or an expression marker (TNF, in BMDCs (Parnas et al., 2015)). In both systems, the Cas9 protein is stably knocked into the cell either *in vitro* (K562 cells) or in the germline of a transgenic mouse from which the cells are derived (Parnas et al., 2015; Platt et al., 2014).

[00718] Altogether, Applicants performed six Perturb-seq experiments, analyzing 200,000 cells in total (**Figure 48C**). In K562 cells, Applicants targeted either 14 TFs or 10 cell cycle regulators, in separate pooled experiments (**Table S1**). For TFs, Applicants performed experiments in both lower and higher MOI, and at two time points (7 and 13 days post infection). In BMDCs, Applicants targeted 24 TFs (**Table S1**) important for either BMDC cell differentiation or the response to LPS (Amit et al., 2009; Garber et al., 2012), and measured their effects at both pre-stimulation (0h) and three hours (3h) post-stimulation. Finally, in each case, Applicants collected reference single cell RNA-seq data from unperturbed cells in separate experiments: 5,400 K562 cells; 2,700 unstimulated BMDCs; and 1,300 BMDCs stimulated with LPS for 3h (**Table S1**).

Example 4: Perturb-seq efficiently detects expressed RNA barcodes, gene expression profiles, and impact on targeted genes

[00719] To couple massively parallel scRNA-seq with a pooled CRISPR/Cas9 screen, applicants developed a lentiviral vector in which a polyadenylated random barcode is associated with an sgRNA and reports on its identity (**Figure 48C, Methods and Resources**). Applicants first matched each sgRNAs with the identity of its corresponding RNA barcode, by either Sanger or next-generation sequencing (**Methods and Resources**), and then converted it into lentivirus for pooled transduction. The guide barcode follows the ORF for BFP, allowing us to sort only infected cells. Finally, because accurate detection of sgRNA from single cells is critical for the accuracy of the pooled approach, applicants designed a dial-out PCR approach, which applicants applied following WTA, to further amplify and sequence the barcode (**Methods**).

[00720] Applicants also developed an enrichment protocol to optimize our ability to detect the sgRNA(s) delivered to each cell. First, Applicants associated each sgRNA with the identity of its corresponding PAPI, by either Sanger or next-generation sequencing (**Methods and Resources**). Applicants next converted the plasmids into lentiviruses for pooled transduction. The plasmid construct included an ORF encoding a BFP-T2A-Puromycin (**Figure 48B**), allowing us to select for transduced cells by FACS sorting or by antibiotic selection. Finally, to increase the accuracy of detection of sgRNA from single cells, Applicants designed a dial-out PCR approach, to enrich for the PAPI following WTA (**Methods and Resources**). Applicants observed low levels of many PAPIs present in many cells (due to cross talk caused by PCR chimeras and ambient RNA) (**Figure 55F**), and account for it by normalizing the observed number of PAPIs within a cell, and setting a threshold (**Methods and Resources**).

[00721] Applicants estimated the probability of detection of the PAPI by considering two factors: (1) the initial MOI and (2) the transcript capture efficiency of our scRNA-seq protocol. Applicants reasoned that these should be described well by a generative model, where applicants assume a zero-truncated Poisson distribution for infection with a guide-carrying lentivirus (zeros are truncated by BFP+ selection), convolved with a binomial process (for the probability of detection) (**Methods**). Indeed, the predicted fit was indistinguishable from the observed frequencies of number of guides per cell (**Figure 55A-C**, bottom panels, KS-test. Using this approach, applicants estimate a 94% (92-96%) detection probability with an approximate initial MOI of 0.63 in the K562 TF pool, a 98% probability with an MOI of 0.35 in the cell cycle pool, and a 60% detection probability with an approximate initial MOI of 1.4 in BMDCs at 3h (**Figure 55A-E**). The lower detection probability in BMDCs reflects overall lower complexity of these

single cell profiles (**Table S1**), and is an additional source of noise, which Applicants handle in downstream computational analysis. Nevertheless, even in BMDCs, Applicants were able to detect cells with multiple guides per cell (**Figure 48D**).

[00722] Table S1

Pool	Number of Cells	Number of sgRNAs	Cells/sgRNA	Average PAPI detected/cell	Median Genes Detected	Median Transcripts Detected
PBMDC 3hr	37369	57	655.596491 2	1.1	1537	5074
PBMDC unstim	32624	57	572.350877 2	1.1	988	2187
K562 TFs low MOI 7 days	33104	30	1103.466666 7	1.3	3015	11610
K562 TFs low MOI 13 days	19268	30	642.266666 7	1.1	2952	12303
K562 TFs high MOI	51898	30	1729.933333 3	1.9	3080	13974
K562 cell cycle regulators	26020	33	788.484848 5	1.2	3764	11716
Wildtype PBMDC 3hr	1310	NA	NA	NA	728	2162
Wildtype PBMDC unstim	2713	NA	NA	NA	585	1450
Wildtype K562	5409	NA	NA	NA	1284	3285
TOTALS	209715	237	5492.09888 4	7.7	17933	63761
AVERAGES	23301.6666 7	39.5	915.349813 9	1.283333333	1992.55555 6	7084.55555 6

[00723] Single cell RNA-seq profiles typically represent only a fraction of the molecular complexity of the measured cell's transcriptome (estimated at 2-26% for droplet based approaches (Klein et al., 2015; Macosko et al., 2015; Zheng et al., 2016)). In both systems, averages of approximately 100 single cell droplet profiles recovered the population complexity (**Figure 55J,K**). (Applicants do note a bias in capture efficiency of small transcripts (<500 nts; **Figure 55L**), which may be associated with the 3' library preparation protocol.)

[00724] Finally, Applicants showed that a reduction of the targeted gene's expression is typically associated with the corresponding perturbation (**Figure 48E,F**, **Figure 55G-I**). In general, this measure should be a proxy for cutting efficiency, although it is possible that regulatory mechanisms can in some cases cause overexpression of a transcript in the absence of a

functional protein. Indeed, the expression of sgRNA targets was significantly reduced from zero (e.g., in BMDCs: **Figure 48F**,

[00725] $p\text{-value} < 10^{-4}$, one-sample t-test). Notably, our ability to determine a reduction in expression is strongly related to the expression level of the target in the wild type cells (**Figure 55I**).

Example 5: A computational model to predict the effect of perturbations on the expression of individual genes and on global cell states

[00726] Applicants next devised a computational framework, based on a linear model, to estimate the impact of perturbations on gene expression (**Figure 48G**, **Figure 49** and **Figure 56**) (**Figure 49**). In the basic linear model (**Figure 49**) (**Figure 48G and 49A**) the logarithm of, the expression of each gene in a cell is explained as a linear (weighted) combination of the effects of individual sgRNAs on the gene. Thus, the log expression level of each gene in each cell (expression matrix **Y**, **observed**) is a product of two matrices: one describing cell attributes – which guides are present in each cell (design matrix **X**, **observed**), and the other represents the regulatory effect of each attribute on each gene (coefficient matrix **β**). This basic model assumes no interactions: that is, the combined effect of two perturbations is simply their weighted sum. The perturb-seq measurements provide matrix **Y** (expression levels in each cell) and **X** (guides in each cell; based on analysis of dial-out PCR measurements of PAPIs; **Methods and Resources**), with good fidelity (**Figure 48**) (**Figure 48** and **55**). Thus, Applicants can fit matrix **β** to determine the effect of each guide on each gene.

[00727] To fit the coefficient matrix **β** (effect of each guide on each gene) given the observations, Applicants used elastic net regularization, which contains both L_1 penalty, to reduce the number of hypotheses tested, and an L_2 penalty on the magnitude of the coefficients (**Methods and Resources**). The second regularization is especially important in the case of correlated covariates and noisy data. Notably, Applicants did not provide any other prior knowledge when learning the model: in particular, Applicants did not include which gene is targeted by which guide. After the model is fit, Applicants evaluate the significance of the fit coefficients using a permutation-based approach to calculate significance (**Methods and Resources, Figures 56D-G**).

[00728] Next, applicants considered how to handle key confounders that are observed in the data. Applicants account for three key classes of covariates that represent either core biological

phenomena or expected confounders (**Figure 49A-G**): (1) the number of observed UMIs in a cell (“cell quality”); (2) the probability that a perturbation successfully affected the cell; and (3) the presence of distinct subpopulations, of either different cell types or states, in the same sample. Below Applicants describe how each of these aspects manifests in the original data, and how it is handled by our model.

[00729] First, as typical in scRNA-seq experiments (Kharchenko et al., 2014; Shalek et al., 2014; Stegle et al., 2015), scRNA-seq profiles from individual cells in the same experiment can vary substantially in their quality, for example as reflected here by the number of UMIs (RNA barcode) (**Figure 49**), assigned to the cell. Applicants ability to detect the expression of a given transcript in any given cell depends on both technical and biological factors relating to the total number of transcripts Applicants captured in that cell and the transcript’s original expression level in that cell (**Figure 49B**). To address this confounder in the basic linear model, Applicants added a covariate to our model that represents the sum of the total number of UMIs detected per cell as well the square of this feature (**Figure 49C**).

[00730] Next, Applicants considered effect of a perturbation (**Figure 49F**) as well as the possibility that in some cells even when a guide is confidently detected, it may not have had any phenotypic effect in that specific cell. This could happen, for example, because the CRISPR system produced no mutations, a heterozygous cell for a fully or partially haplosufficient gene, an in-frame mutation, or even if the editing event was too recent to have had a phenotypic impact. Indeed, for most guides, while cells carrying the guide are distinguishable on average from those that do not, there are still cells with reliably detected guides, but with profiles that are similar to those of cells without the guide (*e.g.*, **Figure 49H**, **Figure 56I**).

[00731] When profiling RNA from a population of perturbed cells in bulk, this mixture would dilute the true perturbation phenotype. If perturbation effects are strong, it should be possible to distinguish true perturbed cells from unperturbed cells leveraging the single cell nature of our data. The regulatory matrix (β) derived from the observed guide-to-cell assignment, provides an initial assessment of the effects of perturbations. Applicants use this model to revisit each cell and predict whether each of the observed guides had perturbed the cell (**Figure 49H**, right). Applicants can then re-estimate the model with the corrected perturbation-to-cell assignment. This one iteration can be thought of as approximate Expectation-Minimization applied to the linear model.

[00732] Third, Applicants reasoned that the cells in each of our experiments might belong to one of several states or types, and that those global states can be affected by the perturbations (**Figure 49E**). For example, in K562 cells Applicants expected (Macosko et al., 2015) cells in different phases of the cell cycle, as well as apoptotic cells. In BMDCs, unsupervised clustering shows several distinct subpopulations, including overlap with those previously reported (Helft et al., 2015). Indeed, even after accounting for cell complexity, many transcripts still show two distinct modes of expression (*e.g.*, **Figure 49C**). Thus, Applicants may want to distinguish between a global impact of a perturbation on cell proportions and a more specific impact within a cell state. To model sub-populations, Applicants first define them using the profiles of the matched, genetically unperturbed, experiments (**Methods and Resources**), and then incorporate this classification as a set of covariates in our model. Applicants fit the model either with or without these subpopulations covariates. When cell states are fit as covariates, the two modes of expression are no longer observed in the residuals (*e.g.*, **Figure 49D**), and some of the sgRNA's effect are accounted for (compare **Figure 49F and G**), suggesting that those perturbations may have primarily affected the proportions of cell types.

[00733] After the model is fit, Applicants can evaluate the significance of the observed coefficients using a permutation based approach to calculate significance (**Materials and Methods, Figures 56D-G**).

[00734] Applicants framework, MIMOSCA (**Figure 49H**), can be readily extended to incorporate non-linear interactions between covariates. In principle, gene perturbations may interact in nonlinear ways with one another ("genetic interactions"), with cell states (cell state specific gene expression changes), or with our quality measures (*e.g.*, changes in cell size). Applicants describe an extension for interactions between perturbations below.

Example 6: The linear model is robust, reproducible and predictive

[00735] Applicants assessed the quality of the linear model by its robustness, fit, and biological coherence. Applicants determined the contribution of the four core components of the model (cell complexity, batch, phenotypic effects of perturbations, and cell state), by analyzing the proportion of the variance in the data each explains (**Figure 49I, Figure 56A-C**). For example, for stimulated BMDCs, the perturbations themselves explain 5% of the variance in the data, 17% is explained when adding covariates associated with number of transcripts detected,

and up to 20% with cell state covariates (**Figure 49I**). Applicants obtain similar results with the other datasets (**Figure 56A-C**). This highlights the importance of accounting for covariates to ensure sensitive detection of the relevant perturbation signal (**Figure 54**). Furthermore, gene-gene correlations in the residuals were significantly reduced as covariates were added (**Figure 49J**). This suggests that much of the global correlation is mediated through transcripts detected and cell states.

[00736] Considering the multivariate nature of gene expression measurements, Applicants also analyzed the consistency in the effects on genes across cells that all contain a particular perturbation (**Figure 56I**). Overall, guides targeting genes have stronger and more consistent effects than the non-targeting control guides.

Example 7: Perturb-seq dissects the transcriptional program in the LPS response in BMDC

[00737] Applicants next analyzed the effect of each of 24 TFs on the transcriptional response in BMDCs to LPS. Briefly, when DCs are stimulated with LPS they activate the expression of ~2,000 genes (Amit et al., 2009; Parnas et al., 2015), organized into several key response modules (Parnas et al., 2015; Shalek et al., 2014). Previously, Applicants have implicated dozens of factors in the response, and validated some of their roles by genetic perturbations (Amit et al., 2009; Parnas et al., 2015). In addition, ChIP-seq analysis of 29 transcription factors (19 of them among the 24 Applicants perturbed) at 4 time points during the response (Garber et al., 2012) showed that factors can be distinguished into three categories: (1) pioneer factors that commit the cells to their lineage (CEBPB and Spi1/Pu.1); (2) broad binders that prime lineage specific genes, before their stimulus-induced expression (e.g., JunB, Irf4, Atf3); and (3) dynamic factors that activate response-specific genes (e.g., Stat2). Knockdown of some factors, but not all, affect the expression of the genes they bind, suggesting buffering genetic interactions for at least some factors (Garber et al., 2012). The response is not perfectly synchronous, such that at different time points, there is bi-modality across the cell population in the expression of specific modules (Shalek et al., 2013, 2014). Finally, Applicants showed that the cells in a BMDC population may consist of at least two sub-types, and that their relative proportions may shift between experiments (Shalek et al., 2013, 2014). Applicants reasoned that this extensive prior knowledge would help benchmark our results, and that the single cell nature of our screen and computational model would directly address the challenges of heterogeneity.

[00738] Applicants performed Perturb-seq at two time points. Applicants obtained precursor cells from the bone marrow of Cas9 transgenic mice (Platt et al., 2014), and after two days of culturing in GM-CSF, Applicants infected the cells with a pool of lentivirus carrying the Perturb-seq vector with 1-4 guide RNAs targeting each of 24 factors (total of 67 guides) and a non-targeting control. After another seven days, Applicants stimulated the cells with LPS, and collected cells for scRNA-seq at 0 and 3h, profiling 32,624 and 37,369 cells, respectively (**Table S1**).

[00739] Fitting our regulatory model without considering sub-populations or genetic interactions (**Figures 50 and 51, Methods and Resources**), the key predictions – the regulatory effects of guides on genes – were coherent in the sense that guides targeting the same gene typically grouped together by their similar impact on genes in the model (*i.e.*, by the regulatory coefficients, β , in matrix **B**; **Figure 51A**). Thus, the correlation in regulatory profiles (β) between guides targeting the same gene is significantly higher than between those targeting different genes (**Figure 50B**, $P < 10^{-9}$, Wilcoxon signed-rank test), even though the model had no knowledge which guides target the same gene. The model also typically shows the repressive effect of guides on the expression of their direct target, although no information was provided to the model on the identity of the gene that it is designed for (**Figure 50A**, separate column on left).

[00740] Applicants uncovered four key modules of transcription factors (**M1-M4**), based on the overall regulatory effects of the different guides as inferred in the regulatory matrix, consistent with their known function (**Figure 50**) (**Figure 50A yellow squares, and 51A**): (1) the anti-viral TFs (Stat1, 2; **M1**); (2) the pioneer factor Cebpb along with the broad binder JunB and the dynamic binders Rela, Stat3 and Hif1a (**M2**); (3) The dynamic binders Rel and Irf2, along with the broad binder Atf3 (**M3**); and (4) The pioneer factor Spi1/Pu.1, broad binders Runx1 and Irf4, and the dynamic binder Nfkb1 (**M4**). While each TF module consists of guides targeting different TFs, the correlation in the regulatory effects of guides targeting the same gene is especially high (*e.g.*, Stat1, Stat2, Cebpb, Rela, Rel, Irf2, Atf3, Spi1/Pu.1, Irf4, Nfkb1, Runx1 etc.; **Figure 51A, 50B**), supporting the specificity of Perturb-seq and suggesting TF specific effects within the context of shared overall programs.

[00741] The four TF modules regulate five major biological programs (**P1-5**; **Figure 51A and B**): an anti-viral response (**P1**), a program of antigen presentation, cytoskeleton and ribosomal

protein (RP) expression, including some markers of DC maturation and of “cluster disrupted” cells (e.g., SerpinB6) (P2), mitochondrial function and biogenesis (P3), an interferon gamma program, with chemokines and cytokines important in the response to intracellular pathogens (viruses, bacteria or parasites) (P4), and an inflammatory TNF response (P5). For example, perturbation of Stat1 and Stat2 (M1 module) represses both the anti-viral (P1) program and the interferon gamma program (P4). Perturbation of the M2 module (Cebpb, Rela, and other factors) represses both the inflammatory and IFN gamma modules, consistent with their roles of activators of inflammation, and activates the P2 program of antigen presentation, cytoskeleton and RPs. The M4 module (Spi1, Nfkb1, Runx1, etc.) has an opposite effect to that of M2. The Rel module (M3) has a partially opposing effect to that of the Stat (M1) module, for example inhibiting the antiviral program.

[00742] These results highlight both known functional associations and novel predictions. For example, Stat1 and Stat2 are well known activators of the anti-viral program (Shalek et al., 2013). The predicted repression of the antiviral program by the Rel module, including Irf2 and Atf3, is supported by a recent report in macrophages (Labzin et al., 2015) that Atf3 can act as a direct transcriptional repressor of the interferon beta (IFN- β) gene, that it further modulates the expression of other genes downstream of IFN signaling, and that Atf3-null macrophages have an enhanced clearing of virus.

[00743] Our model also predicted that Stat1 and 2 are repressors of a program of mitochondrial function and biogenesis in BMDCs (Figure 51A and B). BMDCs undergo a shift from mitochondrial to glycolytic metabolism during the LPS response, including mitophagy (Cohen and Sheva, 1998) (Jovanovic et al., 2015). Recent studies showed that Stat1 inhibits mitochondrial biogenesis in mouse liver (Sisler et al., 2015), and reported Stat2 deficiency in two children as a novel disorder of mitochondrial fission (Shahni et al., 2015). Stat1 and 2 were also reported as physically present in the mitochondria, where they may provide both sensing and signaling functions (Meier and Larner, 2014).

[00744] The comprehensive model, predicting the effects of each factor on each gene, also provides a detailed view into the way in which DCs control different processes. For example, Gbp2,2b,3,4,5 and 7, a set of IFN γ regulated genes important for the response for intracellular infections by bacteria or parasites, are all positively regulated by Stat1 and 2, and negatively regulated by Rel and Irf2. Stat1 has been reported as an activator of GBPs (Ramsauer et al.,

2007) and Rel is reported to bind the Gbp2 promoter (Wei et al., 2008). However, Gbp3,4,5, and 7 are also repressed by **M2** (CEBPB, Rela and others), and hence are members of **P4**, whereas Gbp2 and 2b are not (and are members of **P1**). This may reflect the distinct roles of different GBPs in infections with different parasites [REF]. Moreover, perturbation of the **M1** and **M3** TF modules affects not only these effector genes, but also Irf8 and Batf, key transcription factors that control the GBPs (Tussiwand et al., 2012). Irf8 expression is repressed under Stat2 perturbation, suggesting that the impact on the effector genes may be at least partially indirectly mediated through Irf8 perturbation. Moreover, the TF Batf is induced during Stat2 perturbation, which may reflect the cell's attempt at compensation, since Batf can act in a compensatory capacity in the response to intracellular parasites and *Leishmania* (Tussiwand et al., 2012). Although both IRF8 and BATF regulate differentiation, this is unlikely to be the case here, as the regulatory effect of Stat2 on GBPs, IRF8 and is maintained – and even strengthened – in a model that accounts for cell states and controls for other differentiation effects (below).

Example 8: Opposing programs of BMDC differentiation controlled by two modules each with a distinct pioneer factor

[00745] Two of the TF modules – **M2** (Cebpb, JunB, Rela, Stat3 and Hif1a) and **M4** (Spi1, Runx1, Irf4, and Nfkb1) had diametrically opposed effects on broad programs: **P2** was repressed by factors in **M2** and induced by those in **M4**, whereas **P4** and **P5** were both induced by **M2** factors and repressed by **M4** factors. (**P4** is regulated by Stats whereas **P5** is not, as Applicants discuss below.) The **P4** and **P5** programs each reflect key aspects of the DC response to LPS and pathogens. Conversely, the **P2** program is less expected: in addition to genes for antigen presentation (MHC II system) and cytoskeleton proteins (important for DC migration after maturation), it is also enriched for ribosomal proteins. **P2** also includes SerpinB6 and CD86: Applicants have previously associated SerpinB6 with “cluster disrupted cells” (Shalek et al., 2014), a sub-population in BMDCs that appears to express some maturation genes even prior to stimulation (Shalek et al., 2014) and refs within]; CD86 is known to be associated with DC maturation (Helft et al., 2015). The expression of ribosomal, cytoskeletal and MHC II proteins has been previously reported to be induced in pre-DCs along DC differentiation from progenitors [REF: PMID 26054720 and Biorxiv: biorxiv.org/content/biorxiv/early/2016/10/07/079509.full.pdf]. Several genes expressed by these pre-DCs, but not earlier progenitors are members of **P2**, including genes associated with early

pre-DCs (Iglas3, Itgax), actin polymerization (Crip1), and late-pre-DCs (Cd74, H2-Ab1, H2-AA, H2-Eb1). Furthermore, a recent study showed that protein synthesis is reduced upon full differentiation [Biorxiv: [biorxiv.org/content/biorxiv/early/2016/10/07/079509.full.pdf](https://www.biorxiv.org/content/biorxiv/early/2016/10/07/079509.full.pdf)]. Thus, **P2** may reflect a distinct cell state or type, for example, maybe cells that are either less differentiated or abortive or arrested in their differentiation *ex vivo*.

[00746] Applicants reasoned these effects may thus be best interpreted by a global impact of some of the perturbed TFs on set of possible BMDC-like cell types or states, such that an increase at one cell type under a given perturbation, would come at the proportional decrease of others. Indeed, the BMDCs population Applicants study do have substantial heterogeneity, including the presence of distinct sub-types, such as “cluster-disrupted” cells among the Cd11c⁺ cells (Shalek et al., 2013, 2014) as well as ~10% of the cells that are Cd11c⁺ (Amit 2009, Garber 2012)]. For example, if a perturbation is detrimental to proper derivation of BMDCs *ex vivo*, it would lead to an increase in the number of cells with the active **P2** program, over those that assume **P4** and/or **P5**.

[00747] To test for this possibility, Applicants first used unsupervised analysis to identify a set of seven cell clusters present in 1,310 wild type cells stimulated with LPS (**Figure 50C, Methods and Resources**). The clusters are significantly associated with the five programs identified by genetic perturbations, with distinct enrichments for **P2** (cluster 2, 5 and 6), **P3** (cluster 1) and **P4/5** (clusters 0, 1, 3, 4) (**Figure 58A**). This supports the hypothesis that **P2** and **P4/5** are states that exist in different cells in the somewhat heterogeneous BMDC population, absent perturbation, but the probability that a cell assumes any of these states over the other options may be impacted by at least some of the TF perturbations. Applicants next tested the association of each guide or targeted gene with each state by either predicting the probability that a guide associates with a state, or assigning cells to states and testing for their enrichment in the state (**Figure 50D, Methods and Resources**).

[00748] Both methods show that cells perturbed for specific TFs are significantly associated with one or more of the cell state clusters, compared to the expected proportion of cells in that state when considering their frequencies in the WT cells?/non-targeted cells? in the same pool. In particular, cells perturbed for Cebpb, and to a lesser extent Rela or JunB, all from **M2**, are enriched in one or more cell clusters matching **P2** and concomitantly depleted in cell clusters

matching **P5**, whereas those perturbed for *Spi1*, *Irf4*, *Nfkb1*, or (to a lesser extent) *Runx1*, have the opposite pattern (**Figure 50D**).

[00749] This suggests a model where factors in **M2** and **M4** – *each involving a different pioneer factor* (Garber et al., 2012) – CEBPB and *Spi1*, respectively – have mutually opposing effects on DC differentiation programs. In this model, **M2** promotes differentiation, leading to LPS-responsive programs (**P4** and **P5**), whereas **M4** promotes a mutually exclusive cell state that is either less differentiated, or less productive (**P2**). Importantly, the two states are present in different cells in a BMDC population absent genetic perturbation. Furthermore, consistent with this model, Applicants observed the two TF modules and their effects with the same pattern prior to stimulation with LPS, when the other programs are largely not detected (**Figure 57A,B and 58B**).

Example 9: Diametrically opposed programs are generated by a regulatory circuit with self-reinforcing and mutually inhibitory transcriptional modules

[00750] Applicants next asked how is the transcriptional circuit wired to control these mutually exclusive states or types. Examining the genes in the key transcriptional programs **P2** and **P5** Applicants noted that they include as members their key positive and negative regulators (**Figure 51B**, bottom, blue and red gene names, respectively). This would lead to both positive (reinforcing) and negative (inhibiting) feedback loops that can support the diametrically opposed, mutually exclusive behavior of **M2** and **M4** (**Figure 51B**) and could explain how disrupting one of the factors in a TF module would ‘flip’ the cell to the other state. Specifically, *IRF4*, a TF in **M4** is a target gene in **P2**, a program that **M2** positively controls, generating a positive feedback loop. Conversely, *Stat3*, another target genes in **P2**, is a member of **M2**, generating a negative feedback loop. Similarly, *Cebpb* and *Hif1a* of **M2** are members of its positive target program **M5**, but so is *Spi1*, a member of its negative regulator module, **M4**. The same organizational principle is present in **P1**, the antiviral program controlled by **M1** (positively) and **M3** (negatively): *Stat1* is a member of the program (based on its positive regulation by *Stat2*; by definition, our genetic experiment cannot ascertain its auto-regulation), but so are *Irf2* and *Atf3*, members of **M3**. This architecture – observed in other immune cell systems (Yosef et al., 2013); reviewed in (Yosef and Regev, 2016) may be a general principle in the DC network, resulting in balanced regulation of the cell’s diverse functions, for example based on functional demands (Okabe and Medzhitov, 2014).

[00751] Applicants further sharpened the details of this balanced architecture by constructing a gene level genetic circuit based on the individual significant effects of the key transcription factors on each other (**Figure 51C and D**). Most of the TF modules (**Figure 51D**, shaded areas), have internally reinforcing positive regulation of one module TF by another (*e.g.*, Hif1a and Cebpb by each other and by JunB, Stat3, Rela (M2); Stat1 and Stat2 (M1); Spi1 by Irf4 and Irf4 by NfKB1 and Runx1 (M4)), and negative regulation between modules (*e.g.*, Cebpb and HIF1a in M2 negatively regulated by Runx1 and NFKB1 in M4; Rela in M4 negatively regulates Irf2 in M2; Stat3 and Rela in M4 repress Rel in M3). Importantly, there is a regulatory asymmetry in the network: some factors (*e.g.*, Nfkb1, Rela, Runx1) regulate many others, but are not regulated themselves, whereas a few factors – Cebpb, Hif1A, Spi1, Irf4, Rel, Atf3 – are particularly highly regulated by many others (~6 incoming edges compared to ~1 for the other factors). These latter factors are typically themselves members of key transcriptional programs, as noted above, and thus would account for broad-impact feedback loops.

Example 10: The genetic circuit is strongly supported by the binding patterns of transcription factors

[00752] To support the genetic model and relate it to underlying mechanisms, Applicants compared the effects of the individual factors to their binding patterns, which Applicants previously measured for 22 of the 24 TFs by ChIP-Seq at bulk populations of BMDCs at 0, 30, 60, and 120 min following stimulation with LPS (Garber et al., 2012). Applicants tested whether the bound targets of each individual TF (by ChIP-Seq) are more likely to be negatively or positively regulated by the TF in our model than genes that are not bound by the factor (**Figure 51E, Methods and Resources**).

[00753] The negative effects of a TF perturbation on individual genes are strongly supported by corresponding binding of the TF to these genes' regulatory regions, and consistent with the TF's known role as an activator (**Figure 51F,G**). In particular, genes bound by Rela and Cebpb, both prior to LPS stimulation and at each subsequent time point, are enriched for negative regulatory effects of the guides perturbing these factors. This effect on targets bound across the time course is consistent with the role of these factors as transcriptional activators and with the strong effect perturbing Cebpb and Rela on cell sub-populations at both unstimulated cells (**Figure 57B**) and at 3h post LPS (**Figure 50D**). Furthermore, genes bound by Stat1 and Stat2 at 2h post LPS, but not earlier, are enriched for negative regulatory effects of the guides perturbing

these factors. This is consistent with the role of Stat1 and 2 as activators and with the fact that they only bind most of their targets later in the response (Garber et al., 2012).

[00754] Similarly, the analysis strongly supports the role of factors as repressors, as the genes that they bind are enriched for positive regulatory effects of the perturbing guides (**Figure 51F,G**). These include Irf2, Atf3, and Nfkb1. Irf2 is bound to its targets both pre- and post-LPS stimulation, and its perturbation indeed induces their expression even at time 0 (**Figure 58B**), both in our current study and in an earlier one using RNAi (Garber et al., 2012). The targets bound by Atf3 and preferentially induced by its perturbation are enriched for anti-viral genes (e.g., $P < 1.32 \times 10^{-6}$), directly supporting our model's prediction that Atf3 is a repressor of the antiviral response, and related observations in macrophages (Labzin et al., 2015), as discussed above. Finally, Nfkb1 (encoding p50) acts as a repressor for its known bound targets, suggesting that its major impact might be as a p50-p50 homodimer, a known transcriptional repressor that recruits HDAC-1 (Cartwright et al., 2016; Elsharkawy et al., 2010), rather than as part of the p65 (Rela)-p50 heterodimer, which is a transcriptional activator (Rela perturbation represses its bound target, consistent with p65's role in activating NFkB complexes.) Applicants return to the relationship between Nfkb1 and Rela when Applicants address genetic interactions below.

[00755] The distinct binding patterns of many of the TFs support their specific and direct associations with three of the regulatory programs: P1, P4 and P5 (**Figure 51A,B**). Compared to a stringent background of only those genes in any of the programs (rather than a genome-wide background), genes bound by the **M1** factors, Stat1 and 2, are enriched in **P1** and **P4** (both positively regulated by **M1**); genes bound by Atf3 and Irf2 (from **M3**) are enriched in **P1** (negatively regulated by **M3**, and consistent with the role of the two factors as repressors; above); genes bound by Atf3 are also enriched in **P4** (including key targets such as IL-6, where Atf3 is known to act as a repressor in the TLR4/LPS response in macrophages [PMID: 16688168]); Stat3 and Rela targets are enriched in **P4**, and Cebpb targets are enriched in **P5**, both programs positively regulated by their module **M2**. **M4** factors Irf4 and Runx1, and to a lesser extent NFKB1, are enriched for bound targets in **P4**, and NFKB1 targets are also weakly enriched in its repressed target program P5.

[00756] The remaining two programs do not show such enrichments for bound TF targets. The mitochondrial **P3** program is not enriched for targets bound by any of the factors, even when comparing to a more lenient, genome-wide background (**Figure 58C**). This supports a model

where Stat1 and 2 do not regulate these genes directly, but possibly through another function, such as mitochondrial signals (Meier and Larner, 2014). This is further consistent with the fact that Stat1 and 2 perturbation affects the mitochondrial module even pre-stimulation, when Stat1 and 2 are not bound to promoters (Garber et al., 2012) and that they repress the module (while act as transcriptional activators). The P2 genes, however, are enriched for binding by many key factors when comparing to a genome-wide background genes (**Figure 58C**). These include Cebpb and Spi1, but also many other factors that do not regulate the program. One possibility is that binding by some of these factors (*e.g.*, Spi1 in **M4**) does directly activate this module. Another is that the **P2** program and associated cell state/type are the outcome of the competing positive and negative effects of **M2** and **M4** on P4/P5, such that it is an alternative default or arrested state taken by the cell instead of the differentiated or responsive state **P4/5**. Finally, in some cases the bound targets of TFs (*e.g.*, Rela, SPI1) are enriched (relative to the stringent background) in programs (*e.g.*, **P1**) where their perturbation does not appear to have a substantial effect on the target genes (**Figure 51A and 51G**).

Example 11: TF-specific programs are revealed once accounting for the global effects on sub-populations

[00757] The presence of a small number of TF modules that share global effects on the one hand, and of factor-specific effects on the other hand, motivated us to distinguish global effects on the relative proportions of sub-populations, from effects of TFs that are independent of the context of a cell type or its global state. Applicants also reasoned that absent such distinction, the unique contribution of individual TFs may be masked by entire TF modules with similar global effects. Applicants therefore introduced the assignment of the perturbed cells into the seven states (subpopulations, **Figure 50C**) as covariates in the model. Once Applicants explicitly modeled these global cell states, guides targeting the same TF grouped particularly tightly together (Stat2, Stat1; Irf2, NFkB and Irf4; Rel; Rela; Irf1; and Irf2 and Egr2; **Figure 50E,F** and **Figure 58D**), compared to their membership in larger modules of guides with similar global profiles observed when global states were not explicitly modeled (**Figure 50A,B**).

[00758] Rather than emphasizing the regulatory system controlling the two major cell types/states (which were now reflected in the covariates), the resulting model (**Figure 58D**) highlighted subtler effects unique to one or a few factors. For example, in this model, Applicants distinguished a positive regulatory effect of Runx1 and Ctcf, and a negative effect of Rel, on the

expression of the antigen presentation (MHCII) system, which was not detectable before. Similarly, Applicants found a strong repressive effect of Irf2 (consistent with its role as a repressor) on the IFN gamma response. Importantly, the two models are complementary: the first emphasized global phenomena, and uncovered a densely inter-connected transcriptional circuit that controls the mutually exclusive choice between two global states or types and allows to maintain both in a mixed population. The second controlled for this major phenomenon, and thus uncovered additional, more specific effects that are independent of (and may be masked by) global cell states.

Example 12: Genetic interactions between TFs affect gene expression and global cell states

[00759] The network Applicants infer from observing the impact of single gene perturbations is densely interconnected. Multiple TFs often impact the same gene and program (**Figure 50** and **51**), but can either have the same, opposite, or no effect on transcription when perturbed individually. These effects can be direct when the two transcription factors bind the target gene. Applicants leveraged the fact that in each pooled Perturb-Seq experiment Applicants obtained cells containing more than one guide (e.g., **Figure 48D**), to analyze the extent to which the observed transcriptional defect (on individual gene targets or on global cell states) in cells where two genes are perturbed deviates from what Applicants would expect based on a linear supposition (**Figure 52**).

[00760] First, Applicants found several significant (non-additive) interaction terms for the extent to which perturbing more than one gene affected the probability of being in one of our seven cell states: that is the proportion of cells observed with the combination of guides was higher or lower than expected given the proportion of cells observed with the individual guides (**Figure 52B**). For example (**Figure 52C**), Runx1-NFKB1 are depleted from cells in state 4 (related to **P4/5**), even though each alone is either enriched (Runx1) or unchanged (NFKB1) in proportion in this state. This may be consistent with a recent report that co-expression of NFKB1 (p50) and Runx1 in macrophages enhances the expression of inflammatory genes (such as Il1b from **P5**) [PMID 27573239]. Applicants even uncovered three-way genetic interactions, such as one between Nfkb1, Rel, and Stat2 (**Figure 52C**), where cells with a 3-way perturbation have a higher probability of assuming cell state 2 (partially corresponding to the **P2** program), even though neither cells carrying single nor double perturbations show a related effect. Two other 3-

way combinations involving NFKB1 and Stat2 are each enriched in a different cell state: NFKB1-Stat2-Spi1 is enriched in cell state 1 (corresponding to **P3**, mitochondrial biogenesis), and NFKB1-Stat2-Stat3 is enriched in cell state 3 (corresponding to **P5**, the inflammatory response).

[00761] Next, Applicants assessed the effect of genetic interactions on the expression of each gene, by learning a new model, expanding the original one to include interaction terms with two genes per term (**Figure 52A, Methods and Resources**). Then, Applicants stratified the genes by each of the 27 possible categories when considering the negative, positive, or no effect of each member of the pair and their interaction (**Figure 52D,E**). The 27 categories in turn belong to four major classes: additive (no interaction), synergistic, buffering, or dominant (when the two factors have opposing effects, and the interaction term enhances one of them). Finally, Applicants assessed for each pair of TFs the relative proportion of target genes that are in each of the 27 categories (**Figure 52E**).

[00762] While some TFs participate in genetic interactions that impact many target of genes, others mostly have additive effects. In particular, most TF pairs involving one of Runx1, Irf1, Irf2, or Irf4 has mostly additive effects. In contrast, interactions that include combinations of Stat1, Stat2, Stat3, Rela, Nfkb1 and Spi1 tend to be highly enriched for interactions: they had non-additive effects on a large proportion of their targets. Surprisingly, only a subset of these, all involving Nfkb1 (Stat1-Nfkb1, Stat3-Nfkb1, Rel-Nfkb1, Spi1-Nfkb1) contain buffering interactions (**Figure 52E**).

[00763] The gene-level analysis of genetic interactions allows us to relate the presence (or absence) of interactions to the molecular mechanisms by which TFs act on their targets, as is illustrated in the case of NFKB1 and Rela (**Figure 52F**). As noted herein, NFKB1 and Rela individually have opposing effects on the genes in both P4 and P5 (the programs reflecting the response to intracellular parasites and the inflammatory response): NFKB1 acts as a repressor, and Rela as an activator, ChIP-Seq profiles indicates that these affected genes are directly bound by the respective factors. The model with interaction terms shows that this group of genes (**Figure 52F**, hatched boxes) splits in two: in some genes, the joint perturbation of Rela and NFKB1 leads to an additive effect (no interaction: final gene expression is a sum of the two opposing effects), in the others, there is a dominant interaction, such that the observed expression reflects the effect of NFKB1 over that of Rela. While both sets belong to the same broad

programs (**Figure 52F**, hatched boxes), and both are enriched for CHIP targets of both NFKB1 and Rela (**Figure 52F**, right panel), only the set with the dominant interaction are *co-bound* by both factors (**Figure 52F**, right panel). Applicants observed a very similar phenomenon – where genetic interactions between two TFs is present only in those genes that they co-bind (but not in other genes from the same program, despite their similar individual effects) – for additional pairs of TFs (*e.g.*, buffering interactions for Runx1-Rel and Irf4-Nfkb1 for **P4** and **P5** targets) (**Figure 59**).

Example 13: Global transcriptional modules and specific effects of TFs in K562 cells

[00764] To explore the generality of our approach, Applicants also performed Perturb-seq targeting TFs in K562 cells, a rapidly proliferating cell line that is not synchronized by a stimulus. The TFs were chosen to be not essential for cell viability, but are known to bind targets in K562 cells. Fitting a linear model without cell state co-variates, the TFs partition into two global modules by the correlations in their coefficients: one including EGR1, CREB1, and ETS1, and the other including GABPA, NR2C2, and ELF1, as well as a distinct response to YY1 (**Figure 53A** and **Figure 60**). Although some guides partitioned inconsistently between the two modules, there was an overall significant consistency between guides to the same gene, both within an experiment (**Figure 60**), and across experiments conducted at different time points (**Figure 53E**), as tested by growing a subset of the cells used for the promoter experiments an additional six days (**Figure 53E**). (Such a later time point is not possible to study in BMDCs, a primary short-lived cell culture.)

[00765] As in BMDCs, Applicants explored the role of global cell states. Initial unsupervised clustering partitioned the cells to 9 “states”, reflecting biological processes such as G2/M transition, mitotic cell cycle, anaphase promoting complex, the ribosome, cholesterol biosynthesis, and hemoglobin complex (**Figure 53C**). Applicants tested each perturbation for enrichment in each state (**Figure 53D**) and found significant associations. For example, cells perturbed in EGR1 are depleted from State 6, enriched for hemoglobin, consistent with EGR1’s known role in this process (PMID 19597182). In another example, cells perturbed in YY1 are enriched from State 5 (enriched for cholesterol biosynthesis genes), consistent with the role of YY1 as a repressor of key genes in this process (PMID: 17925399, 26310456, 11145566).

[00766] Next, Applicants fitted a model that accounted for cell states. Applicants reasoned that, unlike BMDCs, many of the states are in fact continuous (*e.g.*, phase in the cell cycle) rather than discrete (*e.g.*, cell types). Applicants therefore performed PCA on the unperturbed dataset, scored the cells from the Perturb-Seq experiment against those principle component (PC) scores, and introduced these state PC scores as covariates to the linear model. In the resulting model (**Figure 53B**) the two global modules are far less correlated, but individual guides to the same gene are much more consistent in their effects (**Figures 53B,F**), suggesting Applicants have controlled for the shared program. This higher consistency is also observed when comparing across experiments and time points (**Figure 53F**), suggesting that TF-specific effects can be reproducibly identified even if global states and cell proportions change over time.

[00767] Applicants used the model that accounts for cell states to explore the individual functions of TFs. For example, Applicants find a distinctive repression of mitochondrial functions upon perturbation of GABPA (**Figure 60Y**) consistent with its role as a required regulator of mitochondrial biogenesis ($P < 2.05 \times 10^{-9}$) (PMID 24958105) and weakly supported by ChIP-Seq binding [REF]. YY1 perturbation represses oxidative phosphorylation ($P < 1.11 \times 10^{-11}$) (known to be regulated by YY1; PMID 12525853, and enriched for its binding by ChIP-Seq [REF]) and induces an NFkB innate immune response ($P < 2.82 \times 10^{-11}$) (PMID 22065573).

Example 14: Perturbations of cell cycle regulators reveal distinct profiles associated with similar fitness effects and mitotic arrest

[00768] Individual cells in a rapidly dividing cell line, such as K562 cells, naturally vary in their cell cycle state. This continuous spectrum is readily detected in scRNA-seq profiles (Buettner et al., 2015; Kowalczyk et al., 2015; Macosko et al., 2015). Factors that impact the cell cycle can be discovered by their impact on morphological features or key marker genes, but the fact that perturbation of two genes has the same overall phenotypic endpoint (*e.g.*, arrest at a particular phase) does not necessarily reflect identical underlying mechanisms.

[00769] To address this question, Applicants targeted in K562 cells 13 genes (with 33 guides) that were previously identified by a mitotic arrest phenotype in a genome-wide arrayed imaging screen in HeLa cells (Neumann et al., 2010) (Table S1). The genes were associated with diverse phenotypes, ranging from a reduction in migration speed to “grape-like” cells with clusters of micronuclei (**Figure 60**, reproduced from (Neumann et al., 2010)).

[00770] First, Applicants used the preponderance of cells with each PAPI to determine the fitness effects of each gene. Unlike BMDCs, where fitness effects are mild, K562 cells are rapidly proliferating and our cell cycle targets are known to affect proliferation (Neumann et al., 2010). Applicants assessed the number of cells carrying each PAPI relative to their expected frequency given the input pool with either of two complementary Bayesian probabilistic models that compute the expected probability of each guide in the resulting cell population (**Methods and Resources**). Applicants identified a strong proliferative advantage conferred by targeting PTGER2, CABP7 and CIT, while a subset of guides targeting genes such as AURKA, TOR1AIP1, and RACGAP1 show proliferative disadvantage (**Figure 53I**). (Similar analysis of the TFs in K562 cells showed a proliferative disadvantage of ELK1 and EGR1, whereas – as expected – there were few fitness effects in DCs). Surviving cells may be enriched for in-frame mutations or off target (yet compensatory) effects.

[00771] Applicants next asked whether these fitness effects can be explained by the genes and processes regulated by the perturbation of each gene. Applicants collapsed all guides for a given gene into a single covariate to increase our ability to discern the true impact of a gene. Applicants found distinct processes affected by those factors with positive and negative effects on fitness, but also two different routes to increased fitness. First, supervised analysis using signature gene sets for the cell cycle phases (Macosko et al., 2015), showed that perturbation of the three Aurora Kinase genes (A, B, and C) and of TOR1AIP1 (which decrease fitness) are all associated with an increase in G2/M and M signature expression (**Figure 53H**). Conversely, perturbation of CABP7 and PTGER2 – which increase fitness – are associated with a strong opposite effect: decrease in G2/M and M signature expression and increase in the M/G1 signature (**Figure 53H**). Notably, perturbation of CIT, which also increases fitness, has an increase in G1/S and S cell cycle states, highlighting a different route to an overall phenotype of increased fitness.

[00772] Furthermore, signature-free analysis of our perturbation by our gene-level model (**Figure 60**) shows that perturbation of CABP7 strongly induces a program of mitochondrial respiration and biogenesis ($P < 1.59 \times 10^{-11}$), NFkB signaling ($P < 2.45 \times 10^{-6}$), and mitotic division ($P < 2.45 \times 10^{-9}$), consistent with the fitness advantage. Perturbation of CIT and PTGER2 repressed these global programs and instead induced the expression of genes of other phases, especially 11 histone genes induced by CIT.

[00773] Notably, the overall partitioning of factors based on their global effects on gene expression also mostly followed their previous reported groupings based on cell biological features in HeLa cells (**Figure 60**, reproduced from (Neumann et al., 2010)): (1) CIT, PTGER2 and RACGAP1 (binuclear phenotype); (2) Cenpe and Arhgef17 (grape-like phenotype and mitotic delay); and (3) Aurorak kinases A, B, and C (proliferation and migration defects). A notable exception is CABP7, which despite a similar binuclear phenotype to that of CIT, PTGER2 and RACGAP1 has a distinct transcriptional phenotype, as described above (**Figure 60**).

[00774] Finally, when examining the nonlinear effects associated with multiple perturbations, Applicants found several significant interactions (**Figure 60**), and noted a significant nonlinear decrease in genes detected when more than three cell cycle genes are perturbed (**Figure 60D**).

[00775] Fitness considers individual guides and synthetic lethality (Livnat model). The effect on G1/S, G2/M cell states is shown (Figure 51C), “virtual FACS” plus Livant model. Figure 51D shows the underlying transcriptional patterns, such that there is a correlation between imaging phenotype and transcription phenotype, such that there is a different expression signature with same imaging phenotype.

Example 15: A brief guide to the miserly: the effects of perturbations on gene signatures and PC scores can be recovered with substantially lower numbers of cells and reads

[00776] The regulatory coefficients associated with the perturbations are highly structured in each of the studied systems (**Figure 50-53**), consistent with the modular organization of genes regulatory systems (Heimberg et al., 2016; Kemmeren et al., 2014; Sokolov et al., 2016). Thus, sets of target genes are affected in coherent ways across sets of perturbations. Satisfyingly, these gene sets are enriched for members of coherent biological pathways (**Figure 50-53**).

[00777] Applicants thus reasoned that recovery of effects at a pathway level should be more robust – and could be detected even if a substantially smaller number of cells is analyzed and/or a smaller number of reads is sequenced per cell. Such a scenario would be highly favorable for researchers who may often desire to employ Perturb-seq to determine key phenotypes (as assayed by their signatures) over fine regulatory connections (such as the targets of each TF, or the wiring and logic of a TF circuit). Applicants thus quantified our ability to detect gene level regulation (**Figure 54A**) vs. signature or state level regulation (**Figure 54B**), when Applicants down-sample cells and reads/cell across our dataset (Applicants had sequenced each library to

near-saturation to permit such analysis; **Methods and Resources**). In particular, Applicants considered signatures quantified in two ways: (1) cell state changes, as defined by our clusters on unperturbed cells; and (2) previously described pathways, such as the antiviral response in BMDCs (Shalek et al., 2014) (**Methods and Resources**).

[00778] Applicants found that the number of cells and reads needed to detect signature level effects (**Figure 54B, 61B-H**) is indeed substantially lower than that needed to detect individual effects on target genes (**Figure 54A, and Figure 61I-L**). In fact, there is almost no change in the observed effects on signatures as Applicants down-sample to approximately *half* the transcripts and cells in our full dataset (148 cells/perturbation (32-671); 780 transcripts); and a Spearman correlation of 0.5 between the effects is still discernable with just 10% of the cells and 25% of the transcripts (30 cells/perturbation (6-134), 390 transcripts). In contrast, the scale required to reach confident estimates in transcriptional effects on individual target genes is higher (~100-200 cells, 500+ transcripts). These estimates – approximately consistent across our three systems (primary cell and cell lines; TFs and cell cycle regulators) – should provide helpful guidelines for future researchers using Perturb-seq.

Example 16:

[00779] Turning to **Figure 66**, GO terms associated with the following sgRNAs: c_sgCIT_1, c_sgRACGAP1_9, c_sgPTGER2_3, c_sgPTGER2_4, c_sgOGG1_4, c_sgECT2_3, c_sgCIT_7, c_sgAURKB_4, c_sgRACGAP1_3, c_sgAURKB_6, c_sgTOR1AIP1_1, c_sgAURKA_3, c_sgOGG1_2, c_sgCENPE_4, c_sgPTGER2_2, c_INTERGENIC216151, c_INTERGENIC345439, c_sgARHGEF17_4, c_sgECT2_2, c_sgCENPE_1, c_sgOGG1_3, c_sgCENPE_2, c_sgPTGER2_1, c_sgCEP55_1, c_INTERGENIC393453, c_sgAURKC_1, c_sgAURKC_7, c_sgCEP55_4, c_sgCABP7_1, c_sgCABP7_4, c_sgCABP7_2, c_sgARHGEF17_1 are shown. The Table below describes the labeling of the figure.

X-axis left to right	Y-axis bottom to top
c_sgCIT_1	chemical_synaptic_transmission
c_sgRACGAP1_9	double-strand_break_repair_via_homologous_recombination
c_sgPTGER2_3	ribosomal_small_subunit_biogenesis
c_sgPTGER2_4	ribosomal_large_subunit_assembly
c_sgOGG1_4	anchored_component_of_membrane
c_sgECT2_3	chromosome
c_sgCIT_7	error-prone_translesion_synthesis
c_sgAURKB_4	positive_regulation_of_viral_transcription
c_sgRACGAP1_3	transmembrane_signaling_receptor_activity

c_sgAURKB_6	oxidative_phosphorylation
c_sgTOR1AIP1_1	ribosomal_large_subunit_biogenesis
c_sgAURKA_3	exon-exon_junction_complex
c_sgOGG1_2	serine-type_endopeptidase_activity
c_sgCENPE_4	heparin_binding
c_sgPTGER2_2	hemoglobin_complex
c_INTERGENIC216151	regulation_of_ion_transmembrane_transport
c_INTERGENIC345439	visual_perception
c_sgARHGEF17_4	potassium_ion_transmembrane_transport
c_sgECT2_2	growth_factor_activity
c_sgCENPE_1	7-methylguanosine_mRNA_capping
c_sgOGG1_3	snRNA_transcription_from_RNA_polymerase_II_promoter
c_sgCENPE_2	chromosome_centromeric_region
c_sgPTGER2_1	chromosome_segregation
c_sgCEP55_1	endosomal_transport
c_INTERGENIC393453	mitotic_metaphase_plate_congression
c_sgAURKC_1	aerobic_respiration
c_sgAURKC_7	viral_life_cycle
c_sgCEP55_4	mitochondrion_organization
c_sgCABP7_1	DNA_duplex_unwinding
c_sgCABP7_4	postsynaptic_membrane
c_sgCABP7_2	small_nuclear_ribonucleoprotein_complex
c_sgARHGEF17_1	telomere_maintenance_via_recombination
	DNA_replication_initiation
	cell_proliferation
	cell_cycle
	proteasome_regulatory_particle
	proteasome_regulatory_particle_base_subcomplex
	COPII_vesicle_coating
	formation_of_translation_preinitiation_complex
	thioredoxin_peroxidase_activity
	double-stranded_RNA_binding
	intracellular_protein_transport
	toxin_transport
	nuclear_matrix
	vesicle
	nucleotide-excision_repair_DNA_incision_5'-to_lesion
	nucleotide-excision_repair_DNA_incision
	error-free_translesion_synthesis
	signal_transduction
	mitochondrial_ribosome
	nucleosomal_DNA_binding
	DNA_damage_response_detection_of_DNA_damage
	damaged_DNA_binding
	G1_S_transition_of_mitotic_cell_cycle
	mitochondrial_proton-transporting_ATP_synthase_complex
	ATP_biosynthetic_process
	mitochondrial_nucleoid

	mitochondrial_small_ribosomal_subunit
	mitochondrial_translation
	mitochondrial_large_ribosomal_subunit
	mitochondrial_intermembrane_space
	gene_expression
	nuclear_import
	SMN-Sm_protein_complex
	spliceosomal_complex_assembly
	protein_peptidyl-prolyl_isomerization
	peptidyl-prolyl_cis-trans_isomerase_activity
	U5_snRNP
	ribonucleoprotein_complex_binding
	RNA_splicing_via_transesterification_reactions
	U4_U6_x_U5_tri-snRNP_complex
	preribosome_large_subunit_precursor
	maturation_of_SSU-rRNA_from_tricistronic_rRNA_transcript_(SSU-rRNA__5.8S_small-subunit_processor)
	positive_regulation_of_GTPase_activity
	response_to_unfolded_protein
	tRNA_binding
	GTPase_activity
	GTP_binding
	GDP_binding
	receptor_complex
	tricarboxylic_acid_cycle
	cell_redox_homeostasis
	NAD_binding
	electron_carrier_activity
	regulation_of_translational_initiation
	nucleotide-excision_repair_DNA_damage_recognition
	spliceosomal_snRNP_assembly
	U2_snRNP
	precatalytic_spliceosome
	chromatin
	methylosome
	U1_snRNP
	proteasome_assembly
	DNA_recombination
	telomere_maintenance
	proteasome_core_complex_alpha-subunit_complex
	positive_regulation_of_protein_localization_to_Cajal_body
	chaperonin-containing_T-complex
	zona_pellucida_receptor_complex
	positive_regulation_of_telomere_maintenance_via_telomerase
	positive_regulation_of_establishment_of_protein_localization_to_telomere
	positive_regulation_of_telomerase_RNA_localization_to_Cajal_body
	tRNA_aminacylation_for_protein_translation

	negative_regulation_of_apoptotic_process
	proteasome_accessory_complex
	ATPase_activity
	ubiquitin-dependent_protein_catabolic_process
	identical_protein_binding
	endoplasmic_reticulum-Golgi_intermediate_compartment
	protein_stabilization
	sister_chromatid_cohesion
	condensed_chromosome_kinetochore
	protein_transport
	spindle_pole
	kinetochore
	DNA_damage_response_signal_transduction_by_p53_class_mediator_resulting_retrograde_vesicle-mediated_transport_Golgi_to_ER
	ER_to_Golgi_vesicle-mediated_transport
	transcription_coactivator_activity
	protein_complex_binding
	intracellular_membrane-bounded_organelle
	regulation_of_signal_transduction_by_p53_class_mediator
	regulation_of_cellular_response_to_heat
	mitotic_nuclear_envelope_disassembly
	intracellular_transport_of_virus
	protein_sumoylation
	RNA_processing
	ER-associated_ubiquitin-dependent_protein_catabolic_process
	nuclear_envelope
	covalent_chromatin_modification
	nuclear_membrane
	protein_complex
	nuclear_chromatin
	spindle
	G2_M_transition_of_mitotic_cell_cycle
	ubiquitin_protein_ligase_binding
	midbody
	endoplasmic_reticulum_membrane
	endoplasmic_reticulum
	ATP_binding
	chromatin_binding
	protein_kinase_binding
	regulation_of_transcription_from_RNA_polymerase_II_promoter
	ligase_activity
	enzyme_binding
	Golgi_apparatus
	ubiquitin_protein_ligase_activity
	ubiquitin-protein_transferase_activity
	transcription_initiation_from_RNA_polymerase_II_promoter
	cellular_response_to_DNA_damage_stimulus

	DNA_binding
	transcription_DNA-templated
	negative_regulation_of_transcription_from_RNA_polymerase_II_promoter
	transcription_factor_binding
	zinc_ion_binding
	focal_adhesion
	cytokine_activity
	cell_surface
	ribosome
	small_ribosomal_subunit
	cytoplasmic_translation
	nucleus
	protein_binding
	nucleoplasm
	extracellular_region
	poly(A)_RNA_binding
	cytosol
	cytoplasm
	membrane
	nucleolus
	integral_component_of_plasma_membrane
	G-protein_coupled_receptor_activity
	G-protein_coupled_receptor_signaling_pathway
	extracellular_space
	mitochondrion
	plasma_membrane
	SRP-dependent_cotranslational_protein_targeting_to_membrane
	nuclear-transcribed_mRNA_catabolic_process_nonsense-mediated_decay
	viral_transcription
	cytosolic_small_ribosomal_subunit
	cytosolic_large_ribosomal_subunit
	translation
	translational_initiation
	rRNA_processing
	structural_constituent_of_ribosome
	homophilic_cell_adhesion_via_plasma_membrane_adhesion_molecules
	cell_adhesion
	calcium_ion_binding
	proteinaceous_extracellular_matrix
	cell_division
	mitotic_nuclear_division
	centrosome
	cell-cell_adherens_junction
	cell-cell_adhesion
	cadherin_binding_involved_in_cell-cell_adhesion

	viral_process
	nucleotide_binding
	mRNA_processing
	nuclear_speck
	termination_of_RNA_polymerase_II_transcription
	mRNA_3'-end_processing
	nuclear_chromosome_telomeric_region
	mRNA_binding
	single-stranded_DNA_binding
	RNA_export_from_nucleus
	transcription-coupled_nucleotide-excision_repair
	external_side_of_plasma_membrane
	U12-type_spliceosomal_complex
	extracellular_matrix_organization
	intracellular_ribonucleoprotein_complex
	transcription_elongation_from_RNA_polymerase_II_promoter
	mitochondrial_respiratory_chain_complex_I
	mitochondrial_electron_transport_NADH_to_ubiquinone
	NADH_dehydrogenase_(ubiquinone)_activity
	mitochondrial_respiratory_chain_complex_I_assembly
	cell-cell_signaling
	DNA_repair
	DNA_replication
	immune_response
	inflammatory_response
	mRNA_export_from_nucleus
	mRNA_splicing_via_spliceosome
	RNA_splicing
	mitochondrial_inner_membrane
	mitochondrial_matrix
	spliceosomal_complex
	mitochondrial_translational_elongation
	mitochondrial_translational_termination
	RNA_binding
	integral_component_of_membrane
	negative_regulation_of_ubiquitin-protein_ligase_activity_involved_in_mitotic_cell_anaphase-promoting_complex-dependent_catabolic_process
	positive_regulation_of_ubiquitin-protein_ligase_activity_involved_in_regulation_of_proteasome_complex
	NIK_NF-kappaB_signaling
	regulation_of_cellular_amino_acid_metabolic_process
	catalytic_step_2_spliceosome
	regulation_of_mRNA_stability
	proteasome-mediated_ubiquitin-dependent_protein_catabolic_process
	protein_polyubiquitination
	extracellular_exosome

	stimulatory_C-type_lectin_receptor_signaling_pathway
	Fc-epsilon_receptor_signaling_pathway
	Wnt_signaling_pathway__planar_cell_polarity_pathway
	antigen_processing_and_presentation_of_exogenous_peptide_antigen_via_MHC_negative_regulation_of__canonical_Wnt_signaling_pathway
	T_cell_receptor_signaling_pathway
	positive_regulation_of_canonical_Wnt_signaling_pathway
	tumor_necrosis_factor-mediated_signaling_pathway
	threonine-type_endopeptidase_activity
	proteasome_core_complex
	translation_initiation_factor_activity
	protein_folding
	unfolded_protein_binding
	myelin_sheath
	melanosome

[00780] Whereas the nature of the experimental system minimizes batch effects when comparing perturbations, as all cells are grown in one pool, applicants sought to determine to what extent the effects observed are reproducible across when compared to a later time point. A subset of the cells used for the promoter experiments were grown for seven more days, and then assayed for expression in the same manner.

[00781] Reports suggest that CRISPRi uniformly reduces expression of a gene across all cells, while the CRISPR-KO system creates an assortment of heterogeneous indel mutations. In order to compare these two systems, applicants perturbed three of the same genes using four sgRNAs per gene with the CRISPRi system.

[00782] Example of cost for conditions and cells:

Cells	Conditions – low (800 cells/condition)	Conditions – high (80 cells/condition)	Cost (\$) – low (1000 reads/cell)	Cost –high (5,000 reads/cell)
10,000	12.5	XXX	250	1,250
50,000	62.5	625*	1,250	6,250

100,000	125	1,250*	2,500	12,000
500,000	625*	6,250*	12,500	62,500
1,000,000	1,250*	12,500*	25,000	125,000
10,000,000	12,500*	125,000*	250,000	1.25 million

*significantly more than previously feasible (100's of conditions)

Example 17: Vectors for capturing guide RNAs in single cells and RT primers for labelling transcripts

[00783] Assays typically measure either a great number of cells, or a great number of genes. RNA-seq allows for the measuring of the expression of every gene in the genome, but can only get average behavior. High content microscopy/flow cytometry allows for the measuring of many individual cells, but only a handful of reporters. Single-cell resolution exposes new phenotypes and the present invention allows for determining rich phenotypes for each perturbation, thus enabling mechanistic understanding and functional clustering (**Figure 68**). Before the present invention guide RNAs could not be easily captured. Applicants developed a novel vector for measuring the guide RNA in a single cell that may be distinguished from guide RNAs present in the pool of cells of the present invention. The vector relies on the capture of a reporter transcript having a unique guide RNA specific barcode (**Figure 70**).

[00784] To build a highly parallel platform for conducting functional genetic studies with complex phenotypes (defined by single cell expression profiling), Applicants developed a simple method for encoding the identity of genetic perturbations into the DNA of individual cells (**Figure 81A-B**). These DNA “barcodes” can be expressed as mRNAs and, therefore, captured during single-cell RNA-sequencing library preparation (**Figure 81A-B**). This strategy is broadly applicable because our barcode design allows the encoding of different pieces of information; however, here Applicants focus on the case where the barcode encodes the identity of a perturbation made using CRISPRi/a (23849981, 25307932). In this system, CRISPR guide RNAs (sgRNAs) are delivered to cells using a lentiviral delivery vector, called the perturb-seq vector, that carries tandem expression cassettes: one encoding an sgRNA under control of a polymerase III-driven promoter and the other encoding a polymerase II-driven mRNA (**Figure**

81B). This mRNA carries a guide-mapping barcode, or “guide barcode,” that uniquely identifies the sgRNA in the former expression cassette, as well as two cell selection markers, blue fluorescent protein (BFP) and puromycin (PURO). Because the majority of single cell RNA sequencing platforms rely on polyadenylation-mediated priming for reverse transcription and capture only the 3' end of message RNAs in sequencing amplicons, Applicants engineered our barcode expression transcript to terminate with a strong polyadenylation signal sequence (BGH polyadenylation signal cloned from pcDNA3.1(+)) and placed the sgRNA-mapping barcode adjacent to this sequence. To incorporate this design into a functional lentiviral delivery system and prevent the BGH polyadenylation sequence from disrupting provirus expression, both the guide and barcode expression cassettes were encoded in reverse orientation with respect to the lentiviral genome. After incorporation into a target cell genome, the expressed barcode transcript can be captured by any oligo-dT-based protocol for cDNA construction. Specific amplification of the barcode transcript can then be used to enrich this important mapping information within RNA-sequencing libraries.

[00785] In the Drop-seq method beads are used to deliver RT primers (**Figure 70**), which are used to tag the 3' ends of transcripts with: (1) Cell barcode, a sequence common to all primers on that bead, used to group transcripts after sequencing; and (2) UMI (Unique Molecular Identifier), a random sequence used to identify duplicates introduced during PCR amplification. Therefore, all of the transcripts in a cell, including the reporter transcript are tagged. The method allows one to introduce diverse perturbations using CRISPRi/a into single cells, pool, and profile responses by droplet RNAseq. Perturb-seq in droplets relies on 3 barcodes to map single cell transcriptomes and analyze expression (**Figure 72**). Vectors can also be used having multiple guide RNAs. In one embodiment, three guides are expressed from a single vector (**Figure 78, 79**).

[00786] Applicants performed pilot experiments with eight cell lines containing transcription factor knockdowns by CRISPRi, ~6,000 cells total. The targets had minimal growth defects, but still showed phenotypes in some screens. One control guide at 3-fold higher representation was used in addition to guides for seven transcription factors: EP300, BHLHE40, CREB1, DDIT3, ZNF326, SNAI1, SPI1. Two sequencing runs were performed on the same library. A HiSeq 2500 rapid run resulted in 132 million reads total, ~22,900 reads per cell, ~10,000 UMIs per cell, and had 2.29X coverage. A HiSeq 4000 lane run resulted in 300 million reads, ~52,000 reads per

cell, ~17,000 UMIs per cell, and had 3X coverage. The merged data set resulted in 75,000 reads per cell, ~20,700 UMIs per cell, and 3.6X coverage.

[00787] A dial-out PCR approach was designed for guide identification (Figure 81). Transcripts containing guide barcodes are specifically amplified by PCR and prepared as separate sequencing libraries. Sequencing coverage is high (16.5 reads/UMI), so almost every molecule is detected (Figure 74). Applicants also capture the guide barcode transcripts during RNAseq and the identities agree, but are less certain (Figure 75).

[00788] Guide identification is done by consensus because most cells have a clearly dominant guide. When two guides show similar abundance, that cell barcode is ignored as a cell doublet.

cell BC	guide identity	read count	UMI count	coverage	good coverage
AGCCTCACAGGCCGA-1	*	132	86	1.534884	False
	EP300_pDS268	3888	243	16.000000	True
	ZNF326_pDS262	829	37	17.000000	True
ATAGCGTGCAACCA-1	*	87	87	1.288507	False
	62(mod)_pBA581	3906	213	18.338028	True

[00789] Guide identification is illustrated the following in test experiment.

rank	guide identity	number of cells	percentage
1	62(mod)_pBA581	1552	28.22
2	SPI1_pDS255	600	10.91
3	EP300_pDS268	525	9.55
4	ZNF326_pDS262	492	8.95
5	SNAI1_pDS266	477	8.67
6	BHLHE40_pDS258	472	8.58
7	CREB1_pDS269	436	7.93
8	DDIT3_pDS263	410	7.45
	*	4	0.07
	Multiplets	353	6.42
	Total uniquely identified	4964	90.25
	Total unidentifiable	183	3.33
	Total number of cells	5500	100.00

[00790] The guide barcodes also permit overloading (Figure 76). If barcodes are sufficiently diverse and recoverable with high confidence, then one can overload and discard doubles, triples, etc. The peak loading possible is 37.8%, although there is a “sweet spot” with sequencing costs.

[00791] Single cell sequencing may be performed by any method known in the art and each has advantages and disadvantages (Figure 80).

Example 18: Discussion

[00792] Here, Applicants developed Perturb-seq a method to analyze the transcriptional defect associated with genetic manipulations as a readout of the cell's phenotype and the architecture of its molecular circuitry – from individual genes to specific processes to global states or types. Perturb-seq dramatically decreases the time and cost limitations associated with assaying the effects of large numbers of perturbations. By decomposing expression into the effects due to technical vs. biological covariates, Applicants recover expected expression changes and uncover new effects, from the effect on individual genes to cell states. Because of the single cell nature of the experiment, Perturb-seq can analyze heterogeneity in cell populations – distinguishing transcriptional effects associated with changes in cell state proportions from transcriptional responses within states.

[00793] The present invention shows novel methods to improve scalability, single cell phenotypes associated with perturbations, and interactions. As both the cost of both sequencing and high throughput library preparation methods decrease, the experimental and computational methods applicants have outlined should become versatile tools in armamentarium of the systems biologist. By sequencing shallowly across large numbers of cells, major transcriptional defects can be readily identified and, by sequencing deeper, smaller changes become apparent. In this invention, applicants exceeded the necessary scale to provide the community with a necessary resource to design future experiments. Excitingly, the signature analysis applicants describe, with as few cells as the number of cells that were analyzed in this work, could be used to gain a broad survey of transcriptional phenotypes across thousands of genes.

[00794] In certain examples, a diverse set of RNA barcodes can be used as an error correcting mechanism when more than one guide is present in a cell and sufficient clonality exists in the population. This feature can be leveraged in future studies to more intensively study the effects of combinations of perturbations and condition our regression framework on clonal similarity. Greater than 1 guide per cell can inform measurements of 1 guide per cell. Greater than 1 cell/drop can reduce costs by an order of magnitude.

Example 19: Perturb-seq and MIMSOCA correctly determine broad phenotypic and gene specific effects of individual perturbations and genetic interactions

[00795] Perturb-Seq analysis of a set of TFs in BMDCs demonstrated the accuracy and resolution of our approach, as well as made specific new discoveries in this cell type. The grouping of individual guides and sets of TFs into TF modules, and many of the genes and processes that perturb-seq and our MIMOSCA framework predicted they regulate were coherent with respect to the extensive knowledge in the system. The specific regulatory connections were strongly supported by ChIP-Seq data. In addition, Applicants uncovered new and emerging regulatory connections, such as between Stat1 and 2 and mitochondrial biogenesis and Atf3, Irf2 and Rel and anti-viral responses, and showed that most programs induce both positive and negative feedback loops.

[00796] Because this novel approach supports such detailed dissection without pre-sorting the cells based on specific marker or state, Applicants could uncover two diametrically opposed programs that govern two global cell states, which precede stimulation with LPS, and may reflect distinct or mutually exclusive differentiation programs. The specific functional nature of “cluster disrupted cells” has been debated, including a suggestion that the two states represent dendritic cells (“cluster disrupted”; here, **P2**) and macrophages (“non cluster disrupted”, here **P4/5**). Our result suggest that both states not only co-exist in unperturbed cells, but are coupled through a circuit with a specific set of positive and negative feedback loops, such that genetic perturbations (in single TFs or their combinations) that enhance one state inhibit the other and vice versa. The specific genes regulated in P2 include many recently highlighted in pre-DC cells , and may thus represent a less differentiated, arrested or abortive state.

Example 20: Alternative embodiments to enhance the precision and facility of Perturb-seq

[00797] Several enhancements can further improve Perturb-Seq. First, including a veritable cornucopia of sgRNAs that are non-targeting or target intergenic regions along with nontransduced controls in the same experiment would help create a stronger reference from which to evaluate the effect of perturbations.

[00798] Second, an important advantage of using Perturb-seq is that higher order interactions can be reached without further need to generate complex reagents (such as new vectors). The same experiment used for a single perturbation, due to the Poisson loading of perturbation per cell, can also uncover the genetic interactions between the perturbed genes. However, the number of sgRNAs present in each cell can vary over a large range. Applicants can leverage the

cassette structure of the endogenous bacterial Cas9 crRNAs (Makarova et al., 2011) and the ability of Cpf1 to autonomously process an entire array, and deliver several sgRNAs (or an unprocessed array) on one construct (Kabadi et al., 2014; Wong et al., 2016; Zetsche et al., 2015, 2016). Such a system has advantages, such as having all members of the cassette be jointly detected or not, but depending on implementation disadvantages as well, such as cloning difficulty, and oligonucleotide synthesis scale.

[00799] At its current scale, given the required number of cells and depth, and with reasonable estimates for the current cost of single cell assays (**Figure 54C**), Perturb-seq can be readily applied for targeted screens of a subset of genes of interest and their interactions, as Applicants have done here for TFs and cell cycle regulators. In some systems, growth or marker based screens may first be performed to identify this subset upstream of Perturb-seq (as done for the unfolded protein response described herein as Applicants have done earlier in BMDCs (Parnas et al., 2015)). The scale of Perturb-seq, will increase as both the cost of both sequencing and of high throughput scRNA-seq preparation methods continues to decrease (**Figure 54C**).

[00800] A researcher can balance the key parameters of the screens based on their wish to focus on major transcriptional phenotypes (cell states or signatures) vs. effects on individual genes' expression, by varying the number of surveyed cells and the sequencing depth. Major transcriptional defects can be readily identified when sequencing shallowly across large numbers of cells; by sequencing deeper (and more cells per perturbation), finer effects become readily apparent. Applicants purposefully exceeded in this study the sample size and sequencing depth required to reach saturation, thus providing the community with a necessary resource to design future experiments. Surprisingly, the signature analysis Applicants describe showed that with as few as ~10 cells per perturbation, a broad survey of transcriptional phenotypes across thousands of perturbations can be performed (**Figure 54B**).

[00801] As cell numbers increase towards genome-wide or large combinatorial screens, Applicants will also need a concomitant increase in the bandwidth of computational methods. In our analysis framework, MIMOSCA (**Figure 56H**), Applicants used a range of efficient algorithms, such as Facebook's Fast PCA implementation (research.facebook.com/blog/fast-randomized-svd/), Infomap clustering (Rosvall and Bergstrom, 2008), parallelizable regression frameworks, and binned permutation significance testing, foreseeing a need for scalability towards millions of cells. In future extensions, when the total number of perturbations is larger,

expression matrices can be stored in chunked form, split by perturbation identity, and the relevant subsets can be read, combined with control cells, and analyzed in parallel.

Example 21: Challenges and opportunities for understanding the vast space of possible genetic interactions

[00802] By considering interactions between the covariates, Applicants showed that Perturb-seq can dissect molecular effects at unprecedented resolution. A key extension is to screen the combinatorial interaction space between perturbations (Hartl, 2014; Phillips, 2008; Sackton and Hartl, 2016; Shao et al., 2008; Zuk et al., 2012), to understand genetic interactions not only between pairs of genes but even at higher order, such as the three-way interaction between NFKB1-Rel-Stat2 Applicants uncovered in BMDC (**Figure 52C**). This can also help uncover global principles of the genetic architecture of a system, such as how robust it is and what is the characteristic number of permissible perturbations (as in **Figure 60**).

[00803] However, there are several major hurdles to achieve this ambitious goal. Both the probability of detecting all perturbations and the probability of all perturbations resulting in an effect scale exponentially with the order of perturbation (**Figure 61A**). For example, when examining cells containing three sgRNAs, assuming half of the cells contain homozygous mutations, and all sgRNAs are independent, the probability of all targets being fully perturbed is one-eighth. Overcoming this source of noise will be a major challenge. However, the alternating-descent framework Applicants described (**Figure 82**) and future improvements to it can help to deconvolve mixtures of knockouts, initialized by the observed assignments between cells and perturbations, and is potentially scalable to higher order interaction terms.

[00804] More critically, while Perturb-seq significantly reduces the time and cost required for transcriptional assays of many perturbations, these limiting factors still scale linearly with the number of perturbations assayed. The size of the combinatorially complete parameter space, p , scales as 2^g , where g is the number of perturbations. For any sizeable g , p is enormous; far more enormous, in fact, that the number of *cells available* in any imaginable setting, regardless of costs, efforts or other technological improvements.

[00805] Applicants hypothesize however, that an alternative strategy should combine a substantial under-sampling of this vast space, with appropriate analytical approaches—inspired by both classical and recent methods (Beerenwinkel et al., 2007; Carter et al., 2007; Jin et al., 2006; Weinberger, 1991), and should aid greatly in our endeavor to infer the impact of

interactions even when p – the number of possible combinations – is much larger than n , the number of samples. Applicants are motivated in pursuing this framing by two assumptions on biomolecular systems, such as the circuits that control gene expression: (1) modularity, as Applicants have shown in this study for both the perturbations and the gene targets, and has been well-established in many genomic and genetic studies (e.g., PMID 14973197, 12134151, 15448693, 27708008); and (2) sparsity, for example, such that Applicants expect most genes not to have genetic interactions (this is supported by those studies – almost entirely in yeast – of pairwise genetic interactions; PMID 27708008). Under these two assumptions, it may be possible, for example, to perform a subset of experiments, when most cells receive a relatively large number of perturbations (e.g., 5) and infer both partially observed and unobserved interactions at a lower order (e.g., 2 or 3).

[00806] Perturb-seq was experimentally designed with such future studies in mind, and our initial framework already started to rely on such assumptions. One example is the L_1 penalty in the elastic net, in our regression framework. A similar L_1 penalty as Applicants apply to the coefficients of \mathbf{B} could in principle be also applied to its rank (Candès and Recht, 2009); considering the group structure in our pathway analysis of perturbations, this may represent a reasonable constraint. More technically, when attempting to fit the combinatorial interaction space, our approach is consistent with the general machine learning framework of finding a function $f(x)$ that minimizes a loss function to best explain our observations y . Applicants effectively search for a sparse polynomial of the form $y = x_1 + x_2 + x_3 + \dots + x_1 x_2 + x_1 x_2 x_3 \dots$ to best explain our observations.

Example 22: Perturb-seq is a first example of a general framework to combine rich readout with cellular metadata

[00807] The current implementation of Perturb-Seq can be readily extended to other traces of the cell's history and to additional rich readouts. First, other CRISPR-based perturbations are readily compatible with Perturb-Seq, including CRISPRi, CRISPRa, and alternative editors such as CRISPR/Cpf1. In particular, Modeling SNPs associated with common diseases or variants mutated in cancers requires efficient homologous recombination followed by cloning of edited cells, but this obstacle may be reduced by the ability to deconvolve wildtype and perturbed cells based on expression signature, as Applicants have done here, or by assays that will directly

detect the edited variant (through DNA+RNA assay or through its encoded transcript) in each cell in the pool.

[00808] The RNA barcodes (PAPIs) used in this study served to identify perturbations, but can also be used to mark cells derived from a common ancestor for the purposes of lineage tracing (**Figure 54D**), thus relating cell lineage and cell fate maps. For example the PRISM pool contains a set of barcoded cell lines (Yu et al., 2016). Their expressed RNA barcodes could be combined with Perturb-seq to enable pooled measurements of the transcription effect of different perturbations across cell lines by fitting an interaction term between the cell line and the perturbation.

[00809] Other measurement platforms may be compatible within the experimental and computational framework Applicants describe. For example, multiplex PCR (Fan et al., 2015) can be used to dial-out a subset of transcripts on which to obtain more reliable information (similar to our PAPI protocol). Protein measurements across millions of cells may be feasible by fusing FISH probes to heavy metal and using CyTOF as the readout technique. It may also be possible to devise paracrine assays by co-culturing perturbed cells in droplets and merging the droplets with a massively parallel single cell sequencing technology.

[00810] The experiments and analysis provide a guide the design of future experiments that leverage the scale and tunable resolution possible by combining high throughput single cell sequencing and pooled screens. By rapidly and simultaneously testing causal relations between many genetic manipulations and the resulting molecular phenotypes, Applicants aspired to help bridge the gap between genetic screening and molecular follow up experiments and to provide a valuable tool to the systems biologist. In doing so, Applicants hope to facilitate casual research questions that infer mechanisms by which genotypes result in specific phenotypes.

Example 23: A robust strategy for encoding information in RNA transcripts that allows parallel profiling of separately perturbed cells by droplet-based single-cell RNA-sequencing

[00811] A particularly fascinating example of a complex transcriptional response is the mammalian unfolded protein response (UPR), a branched pathway critical for maintaining homeostasis of the endoplasmic reticulum (ER). In metazoans, this response is coordinated by three distinct ER transmembrane sensor proteins (IRE1 α , ATF6, and PERK), which activate three individual transcription factors (XBP1, the N-terminal cleavage product of ATF6, and ATF4, respectively) to promote survival or, when ER stress cannot be corrected, trigger cell

death pathways (Ron, Walter 2007, Walter, Ron 2011, Hetz 2012). The molecular signaling events downstream of each UPR sensor protein are discrete. The resulting transcriptional responses, however, are overlapping and integrated and, therefore, the response to any given stress is likely to be heterogenous. Among stresses that can activate the UPR are potentially deleterious changes to protein folding, as well as disruption of normal nutrient flux, redox potential, calcium homeostasis, hormone stimulation, lipid biosynthesis, and membrane integrity. Considering the diversity of inputs monitored by the UPR and the necessary plasticity in outcome, comprehensive characterization of this system in mammalian cells requires both unbiased profiling of the physiological stresses that activate the sensors and delineation of the complex transcriptional phenotypes for each input.

[00812] Here, Applicants also apply perturb-seq to the systematic analysis of the mammalian UPR. To independently manipulate the three branches of the UPR, Applicants first developed a programmable strategy for simultaneously depleting up to three genes with high efficacy. Applicants then used perturb-seq with combinatorial depletion of the UPR sensors to delineate the distinct transcriptional programs of the UPR and the epistatic relationships among the three UPR branches. Next, Applicants used a two-tiered approach to interrogate the biological systems monitored by the UPR. Applicants identified hundreds of genes that contribute to ER homeostasis from a genome-wide CRISPRi screen and then applied perturb-seq to interrogate a diverse subset of these genes with single-cell resolution. These experiments allowed us to systematically define the functional relationships between genes with high precision and to dissect the complex, partially overlapping transcriptional responses to ER stress. Furthermore, analysis of the single-cell responses revealed bifurcation of the UPR branches at two levels: among individual cells subject to the same perturbation and at the population level, where differential activation of the three UPR branches occurred across perturbations. The latter indicated a dedicated feedback loop that enables a single arm of the UPR (the IRE1 α /XBP1 branch) to specifically monitor the integrity of the protein translocation machinery. These data, alongside the parallel analysis of dendritic cell activation described herein, demonstrate the ability of perturb-seq to provide rich biological insights and systematically dissect complex biological responses.

[00813] The recent introduction of scalable, droplet-based approaches for single-cell gene expression profiling and the development of commercial devices (10X Genomics) have made

measurement of complex phenotypes highly parallelizable (Macosko et al. 2015, Klein et al. 2015, Zheng et al. 2016). These methods incorporate two barcoding strategies that allow deconvolution of pooled RNA-seq data into single-cell transcriptomes. A cell barcode, affixed to cDNA molecules during reverse transcription of mRNAs from a single cell, is used to computationally link data from individual cells, and a similarly appended unique molecular identifier (UMI) collapses information gathered from individual mRNA molecules to allow molecular counting (Figure 81A). A current limitation of these methods, however, is that they rely on oligo-dT priming during cDNA synthesis to capture mRNA identity and therefore report only on polyadenylated RNA transcripts. To enable the recording of other types of information, Applicants built a platform that allows genetically encoded barcodes to be expressed as polyadenylated transcripts and captured alongside endogenous mRNAs (Figure 81A). These barcodes can be used to mark specific cell perturbations (e.g., the identity of a Cas9-targeting single guide RNA, sgRNA), and thus, allow complex pools of cells to be evaluated by single-cell RNA-sequencing on existing droplet-based platforms.

[00814] Central to this approach is our barcode delivery construct, a third generation lentiviral vector carrying a polymerase II-driven barcode expression cassette that terminates with a strong polyadenylation signal (PAS) (Figure 81B) (Goodwin, Rottman 1992). To ensure that the PAS of this internal cassette does not disrupt transcription of the lentiviral genome, and therefore lentivirus production, the entire barcode expression cassette was placed in reverse orientation with respect to transcription initiated by the upstream hybrid long terminal repeat (LTR) promoter. Within this region, the PAS was encoded directly adjacent to the barcode site (127 nucleotides separate the barcode site from the AAUAAA hexanucleotide), a placement designed to favor faithful transmission of barcode sequences into single-cell RNA-sequencing libraries, which typically capture and count only the 3' ends of polyadenylated transcripts. To test barcode capture, Applicants performed single-cell RNA sequencing using the Chromium™ Single Cell 3' Solution (10X Genomics) on a pool of individually transduced chronic myeloid leukemia cells (K562) carrying 8 distinct barcodes. To ensure robust barcode capture with minimal sequencing depth, Applicants developed a protocol to specifically amplify barcode transcripts (Figure 81A). In a pilot experiment of ~5,700 cells, Applicants observed a median of 45 independent observations (marked by unique molecular identifiers) of the “guide barcodes” per cell using this approach, allowing us to uniquely infer a guide barcode for 92.2% of cells (Figure 81C,D). The

guide barcodes also allowed us to discard droplets that fortuitously received more than one cell (marked by the presence of two or more distinct guide barcodes assigned to one cell barcodes) (Figure 81D).

[00815] A practical way to build complex pools of cells experiencing different perturbations is by variably modulating gene expression, and the advent of CRISPR-based genetic approaches has revolutionized the specificity with which this can be done in mammalian cells (Shalem, Sanjana & Zhang 2015). By including an sgRNA expression cassette in our barcode delivery vector, Applicants tailored our barcoding system to the study of complex CRISPR-based phenotypes (Figure 81B). To begin validating this “perturb-seq” approach, Applicants confirmed that sgRNA expression from our barcoded construct, the “perturb-seq vector,” was capable of generating robust CRISPRi-mediated gene depletion phenotypes (Figure 81E). Specifically, Applicants tested an EGFP-targeting sgRNA (programmed with a previously validated protospacer, EGFP-NT2 (Gilbert et al. 2013)) in the context of our perturb-seq vector for its ability to repress expression of a genomically integrated GFP reporter. Guide activity from this construct was robust and comparable to that from a previously validated sgRNA expression construct (Gilbert et al. 2014) (95.4% and 96.2% repression of GFP fluorescence, respectively). Importantly, this experiment also demonstrated that, consistent with previous findings (Gilbert et al. 2013), CRISPRi mediates highly homogenous protein depletion across a population of perturbed cells. This ability carried over to pooled experiments, in which the guide barcode can be used computationally to identify the subpopulations of cells bearing particular CRISPR perturbations. Within each subpopulation, the RNA-seq profiles always revealed knockdown of the expected target, as seen in example perturbations from our pilot experiment (Figure 81F) demonstrating our ability to genetically encode and reliably capture the identity of CRISPRi-mediated perturbations in pooled format with single-cell resolution.

Example 24: A strategy for multiplexed sgRNA delivery that allows systematic and simultaneous induction of robust and homogeneous genetic perturbations

[00816] Having established the potential for perturb-seq as a high content screening platform, Applicants sought to interrogate the unfolded protein response (UPR). Although well studied at the level of mechanistic signaling and bulk population behavior (Ron, Walter 2007, Walter, Ron 2011, Hetz 2012), much remains unknown about how individual cells execute the UPR. Of particular interest is the potential for cell-to-cell heterogeneities in activation of the metazoan

UPR, because, unlike the linear pathway in yeast, this branched stress response integrates distinct and separable signaling events to coordinate downstream transcriptional programs (Figure 82A). Specifically, there are 3 signaling arms of the metazoan UPR, each controlled by a different ER transmembrane sensor protein. These sensor proteins are IRE1 α , PERK, and ATF6 (Walter, Ron 2011). Briefly, IRE1 α is an endoribonuclease that, upon sensing ER stress, mediates noncanonical splicing of *XBP1* mRNA and facilitates stabilized expression of the active XBP1 transcription factor (XBP1s). PERK is a kinase that, upon activation, phosphorylates the alpha subunit of translation initiation factor eIF2 (eIF2 α). This inhibitory phosphorylation event generally suppresses translation but paradoxically promotes translation of a small group of mRNAs, including one encoding the ATF4 transcription factor. Lastly, under stress conditions, ATF6 is targeted to the Golgi where proteolytic cleavage releases a cytosolic transcription factor domain. Once activated, XBP1, ATF4 and cleaved ATF6 translocate into the nucleus to initiate transcription.

[00817] Applicants sought to systematically interrogate the integrated programs of UPR-mediated transcription using perturb-seq with two motivations: first, to delineate the IRE1 α , PERK, and ATF6 transcriptional regulons, which is inherently challenging because of cross-regulation between branches, and second, to evaluate both cell-to-cell and perturbation-to-perturbation heterogeneity in UPR activation. To accomplish our first goal, Applicants had to expand the perturb-seq platform for the systematic analysis of higher-order genetic interaction. Therefore, Applicants developed an approach that allows combinations of genes to be robustly and homogeneously perturbed in a single cell (Figure 82B).

[00818] Previous efforts to multiplex genetic perturbations by expressing multiple sgRNAs for Cas9 or the alternative programmable nuclease Cpf1 have had limited success achieving consistent and uniform perturbation of multiple genes simultaneously (Kabadi et al. 2014, Nissim et al. 2014, Zetsche et al. 2016). To begin addressing this issue, Applicants adapted a previously reported strategy (Kabadi et al. 2014) of fusing three sgRNA expression cassettes (each composed of an RNA polymerase III promoter, the sgRNA targeting sequence, and the sgRNA constant region) into a lentiviral vector carrying our reverse barcode expression cassette (Figure 82B). To minimize recombination at repetitive nucleotide sequences during lentiviral transduction (Sack et al. 2016, Smyth, Davenport & Mak 2012), Applicants used three different RNA polymerase III promoters (initially Applicants used a modified mouse U6 promoter (mU6),

a modified human U6 promoter (hU6), and a modified human 7SK promoter (h7SK), but replaced h7SK with a more potent modified bovine U6 promoter (bU6) in later work). To evaluate the activity of our initial three-guide expression vectors, Applicants assembled three test constructs that each express EGFP-NT2 from one of the promoters (and negative control sgRNAs from the other two) and monitored GFP depletion in K562 cells using flow cytometry (Figure 88A). Applicants found that the h7SK promoter generally performed poorly in the context of our perturb-seq vector (Figure 88B). More importantly, however, cells transduced with our first generation three-guide vectors expressing EGFP-NT2 from the hU6 promoter or the mU6 promoter partitioned into two distinct subpopulations with either strong depletion (>90%) or no depletion of GFP (Figure 88A). To exclude the possibility that dCas9-KRAB expression levels limited the activity of the three-guide vectors, Applicants engineered a GFP+ K562 cell line with increased levels of dCas9-KRAB expression (Figure 88C, see Methods). When transduced with a three-guide vector expressing EGFP-NT2, these cells still exhibited bimodal GFP expression (Figure 82C). A likely source of this inactive population is intramolecular recombination within the 93-nt constant regions, which would lead to elimination of one or two sgRNA expression cassettes (Sack et al. 2016).

[00819] Consequently, Applicants sought to engineer modified sgRNA constant regions that were functionally equivalent but distinct enough in sequence to suppress recombination. Guided by crystal structures of Cas9 (Nishimasu et al. 2014, Anders et al. 2014, Jiang et al. 2015), literature reports of modified constant regions (Briner et al. 2014, Dang et al. 2015), and sequence conservation of tracrRNAs in *Streptococcus* species (Briner et al. 2014), Applicants identified bases that seemed amenable to mutation and designed a panel of variant constant regions (Figures 82D, 88D). Almost all constant region variants retained strong CRISPRi activity as indicated by a 97-99% reduction in mRFP levels using a bacterial CRISPRi assay suggesting that the introduced mutations do not disrupt sgRNA:Cas9 binding (Figure 82D). By combining sets of neutral mutations, Applicants obtained two new constant regions (cr2 and cr3) that share at most 20 bases of continuous sequence homology with each other and the original constant region. Both cr2 and cr3 exhibit strong CRISPRi activity in *E. coli* (Figure 82D).

[00820] To test the activity of cr2 and cr3 in human cells, Applicants generated sgRNA expression cassettes containing the mU6 promoter paired with cr2 and the hU6 promoter paired with cr3 in the context of our perturb-seq vector. When used to express EGFP-NT2, both

cassettes yielded GFP depletion that was indistinguishable from that of the original perturb-seq vector in K562 cells (Figure 82E). Finally, to identify a third functional RNA polymerase III promoter, Applicants cloned a panel of U6 promoters from different mammalian species into our perturb-seq construct with the original constant region and tested them in our GFP depletion assay. Several U6 promoters mediated strong CRISPRi activity (Figure 88E), with the bovine U6-2 promoter (bU6) (Lambeth et al. 2006) having the largest effect.

[00821] Using our three sgRNA expression cassettes (bU6-2 with original cr (cr1), mU6 with cr2, hU6 with cr3), Applicants then assembled our second iteration of three-guide expression vectors (Figures 82B, Figure 88F) and tested them in our GFP depletion assay. GFP+ K562 cells carrying dCas9-KRAB (higher expression levels) transduced with three-guide vectors expressing EGFP-NT2 from any of the three promoters exhibited near-uniform and strong knockdown of GFP (96-97%), nearly identical to that mediated by our single perturb-seq vector (Figure 82F). These data establish that recombination at the constant regions during lentiviral transduction is an important limitation when fusing multiple sgRNA expression cassettes and that new highly active sgRNAs can be engineered to suppress recombination. Our final three-guide vectors, which use three different RNA polymerase III promoters and three different constant regions, can be faithfully transduced into cells alongside a single perturb-seq barcode to mediate robust depletion of targeted genes.

Example 25: Systematic delineation of the transcriptional responses controlled by the unfolded protein response using perturb-seq

[00822] Applicants next sought to use the perturb-seq approach together with our ability to simultaneously repress expression of multiple genes to systematically explore the three branches of the mammalian UPR (Figure 83A). Using our triple sgRNA expression vector, Applicants transduced sgRNAs targeting each UPR sensor gene (*ERN1* encoding IRE1 α , *EIF2AK3* encoding PERK, or ATF6) in all possible combinations (i.e., single, double, or triple depletions) into K562 cells. Applicants chose K562 because they have a robust UPR and as suspension cells avoid heterogeneity due to difference in their microenvironment (Snijder et al. 2009). Transduced cells were then pooled, sorted for BFP expression (a marker for sgRNA delivery), and after 2 days of growth, cells were treated with thapsigargin, a non-competitive inhibitor of ER calcium pumps, or tunicamycin, an inhibitor of N-linked glycosylation, well-known pharmacological inducers of the UPR. Control cells were treated with DMSO. Applicants

sequenced transcriptomes of ~15,000 cells, of which 5% were called as multiplets and 5% were unidentifiable, to a depth of ~120,000 reads per cell, and observed a median of ~25,000 UMIs per cell arising from ~4,300 unique genes. *ERN1*, *EIF2AK3*, and *ATF6* were depleted >80% when targeted in any combination, and across the three chemical treatments (Figure 82G).

[00823] Applicants then set out to construct an analytical approach for finding statistically robust features within the data (Figure 83B, 89A). Every perturb-seq experiment generates an extremely rich dataset, but single-cell RNA-seq data are both intrinsically noisy (due to efficiency limitations during library preparation and natural stochasticity of transcription at the single-cell level) and of very high-dimension (Grün, van Oudenaarden 2015). The latter problem is usually overcome using unsupervised dimensionality reduction methods like principal components analysis (PCA). As many genes share common regulation, cellular behavior is intrinsically low-dimensional (Heimberg et al. 2016, Pollen et al. 2014), thus cellular behavior can be effectively summarized in terms of tens of aggregated effects instead of thousands of individual genes. A limitation of PCA is that it can be arbitrarily corrupted by outliers, and principal components often do not have clear meanings. In our hands, PCA was strongly influenced by clear outliers, such as apoptotic cells (data not shown).

[00824] To mitigate these problems Applicants developed an analytical pipeline, Applicants term low rank independent component analysis (LRICA), which proceeds in two stages (See Methods). Applicants first leverage recent advances in the mathematics of sparse matrices (Candès et al. 2011, Lin, Chen & Ma 2010) to directly decompose the observed gene expression matrix (X) into a low-rank matrix (L), representing the low-dimensional dynamics shared across the population, and a sparse matrix (S), capturing noise and effects that are highly variable between cells:

$$X = L + S$$

[00825] The low rank matrix contains the major dynamic trends in the population. Applicants identified these trends using independent component analysis (ICA, Methods), which deconstructs the gene expression pattern in each cell into a linear sum of components that are maximally independent of each other (instead of uncorrelated, as in PCA). The components aid interpretation in two ways: by determining which cells are high or low for a particular component, Applicants can identify subpopulations, and by asking which genes maximally

influence a component, Applicants can identify groups of related genes driving a particular behavior.

[00826] Applicants applied LRICA to our thapsigargin-treated population. Four components varied across the different perturbations, including three that tracked the presence of PERK, IRE1 α , and ATF6 (Methods, Figure 89B). When projected to two dimensions using *t*-distributed stochastic neighbor embedding (t-sne) (Van Der Maaten 2014), cells bearing a particular perturbation all grouped together, and biologically reasonable groups of perturbations also clustered in proximity to each other. Thus, the identity of the perturbation contained in each cell was distinguishable, which further validated our triple knockdown strategy. The same analysis applied to the four components that varied across the cell cycle arranged the cells in a circular pattern ordered correctly by cell cycle phase. Thus, LRICA identified and decomposed the two largest effects causing variation in the population in an unbiased way and also computationally decoupled them from each other, enabling them to be studied independently.

[00827] Applicants did observe a small but prominent interaction between the two effects, apparent in a “bulge” seen in the cell cycle analysis (highlighted in Figure 83B) where cells bearing certain perturbations concentrated. Closer examination of a more targeted set of three perturbations suggested that this interaction was caused by PERK-dependent cell cycle arrest in G1 caused by thapsigargin treatment (Figures 83C,D) (Hamanaka et al. 2005). A particular component (right panel of Figure 83C) was bimodal among the cells bearing each perturbation. Defining the cells with that component low as “G1 cells” (cf. middle and right panels of Figure 83C), Applicants looked at the top fifty genes influencing the component (Figure 83E). This list included a set of factors associated with cell cycle progression and the UPR. Intriguingly, it highlighted epistatic interactions between the two effects of PERK-dependent activation of the UPR and progression through G1, showing for example how for some genes the two programs cancel each other out, while for others they may act synergistically. A notable example of the latter was the expression of *MYC*, which has previously been identified as being thapsigargin-induced (Liang et al. 2006). Our data here show that *MYC* induction is most strongly associated with the G1-arrested subpopulation (Figure 83E).

[00828] Given our success in decomposing independent effects, Applicants next turned to decomposing the distinct branches of the UPR itself. Applicants identified a set of genes that were most robustly induced in response to both thapsigargin and tunicamycin treatment and used

hierarchical clustering to order them based on their pattern of co-expression in the 15,000 cells within our experiment (see Methods). Applicants constructed synthetic bulk RNA-seq profiles by averaging together all cells bearing the same perturbation. When these profiles were ordered according to our clustering, clear patterns of regulatory control were apparent (Figure 83F). To estimate regulatory overlap, Applicants decomposed the changes across bulk responses using ICA (bottom of Figure 83F, Methods). PERK/ATF4 had the largest regulon in our experiment, and had many targets uniquely under its control. ATF6 and IRE1 α showed more overlap in their set of regulatory outputs, consistent with their common transcriptional regulatory mechanism. Of the two, IRE1 α had more specific targets, notably components of the translocon (e.g. *SEC61A1*, *SEC61B*), which is consistent with previous reports (Shoulders et al. 2013), but ATF6 had stronger activating effects on common targets (Figure 83F). Many genes showed some sensitivity to all branches, particularly a group of very high abundance core stress response genes (*HSPA5*, *HERPUD1*, *SDFL21*). Our epistasis experiment thus defined and decoupled the three overlapping branches of the mammalian UPR, both at the bulk level and within single cells.

Example 26: Unbiased genome-scale CRISPRi screens identify a library of genetic perturbations that induce the UPR

[00829] Having systematically mapped the transcriptional programs controlled by each of the three UPR sensors, Applicants next sought to leverage this information and our platform to evaluate what role each branch plays in responding to different types of perturbations. A principal advantage of perturb-seq is the ability to evaluate transcriptional phenotypes for many genetic modifications in a single, internally controlled experiment. Thus Applicants can ask fundamental questions about how transcriptional pathways differentially respond to diverse and physiologically relevant inputs, and importantly, Applicants can do this in cells otherwise exposed to identical conditions. Applicants employed a two-tiered approach to use perturb-seq to systematically explore transcriptional pathways, such as the UPR. First, Applicants performed a genome-scale CRISPRi screen to identify genes that are important for ER homeostasis, as evidenced by induction of the UPR upon depletion. Applicants then applied perturb-seq to a diverse sub-library targeting the strongest hits to better understand their role.

[00830] To begin, Applicants constructed a fluorescence reporter of IRE1 α activation. This reporter drives mCherry expression from a minimal *c-fos* promoter controlled by 5 IRE1 α -responsive unfolded protein response elements (UPREs) (Wang et al. 2000) (Figure 84A). The

reporter was transduced into K562 cells stably expressing a dCas9-KRAB fusion protein and a clonally derived line (cBA010), which demonstrated strong reporter activation (Figure 84B), was selected. To control for general effects on gene expression during screening, this cBA010 parental line was subsequently transduced with a constitutive GFP reporter (driven by the EF1a promoter) (Figure 84A). Importantly, when treated with tunicamycin, these cells (cBA011) demonstrated *XBPI*-dependent mCherry induction (maximally 16-fold), which occurred subsequent to endogenous XBP1 splicing (Figures 84B, 90A). As expected, Applicants observed no similar induction of GFP.

[00831] Using our reporter cell line, Applicants separately screened two genome-scale CRISPRi libraries, our first generation library (hCRISPRi-v1) and our recently described second generation library (hCRISPRi-v2) (Figure 84C) (Gilbert et al. 2014, Horlbeck et al. 2016). Briefly, cBA011 K562 cells were transduced with each library. After 8 days of growth, the cells were separated by fluorescence-activated cell sorting into bins according to the ratiometric reporter signal (mCherry/GFP). Cells in the top and bottom thirds of the reporter distribution were collected and processed to measure the frequency of sgRNAs contained within each bin by deep sequencing. Using these data, Applicants calculated the \log_2 sgRNA enrichments in the top mCherry/GFP bin over the bottom bin, and defined gene-level reporter signal phenotypes as the average of the 3 most active sgRNAs targeting each gene. Our hCRISPRi-v2 screen identified 397 hit genes with high mCherry/GFP, indicative of UPR activation (Figure 84D). Importantly, phenotypes were reproducible between replicates of this screen and minimal correlation was observed between hit phenotypes and previously calculated gene growth phenotypes (Spearman $R = -0.2$) (Figure 84E, 90B).

[00832] Our first generation CRISPRi library targets 15,977 genes (20,899 transcriptional start sites, TSSs) with 10 sgRNAs per TSS, and our hCRISPRi-v2 library targets 18,905 genes at 20,526 TSSs. hCRISPRi-v2 incorporates several improvements over hCRISPRi-v1, including a greater frequency of highly active sgRNAs and the flexibility to screen with either 10 sgRNAs per gene or a similarly active half-library containing the top 5 predicted sgRNAs for each gene (Horlbeck et al. 2016). Practically, the compact hCRISPRi-v2 library allowed us to maintain higher screen representation through the flow cytometer with similar sorting times (screen representations: hCRISPRi-v1 ~450, hCRISPRi-v2 ~600). Encouragingly, of the 141 hits from

the hCRISPRi-v1 screen, 103 reproduced from screening the hCRISPRi-v2 library (Fisher's Exact p-value = 8.97e-138) (Figure 90C).

[00833] Among hits from the hCRISPRi-v2 screen are well-characterized regulators of protein folding in the ER, most notably *HSPA5*, which encodes the major ER Hsp70 chaperone BiP (Figure 84F). Consistent with results from a similar screen in yeast (Jonikas et al. 2009), our hits also featured genes functionally involved in *N*-linked glycosylation, including components of the oligosaccharyltransferase (OST) and the dolichol-linked oligosaccharide biosynthesis pathway, ER-associated degradation (ERAD), and protein trafficking. Additionally, genes involved in protein targeting to the ER were enriched among hits (Fisher's Exact p-value = 2.65e-09). This gene set consists of genes encoding components of the signal recognition particle (SRP), the SRP receptor, and the signal peptidase complex. Three out four of the translocon-associated protein complex (TRAP) scored; and strikingly, among the 7 hits with the strongest phenotype were all three genes that encode the ER protein-conducting channel or PCC (*SEC61A1*, *SEC61B*, *SEC61G*) (Figure 84D,F). The phenotypes of SRP targeting factors and PCC were surprising because recent reports have shown that SRP-mediated recruitment to the ER and IRE1 α binding to PCC subunits are required for maximal pathway activation in response to exogenous stress (Plumb et al. 2015, Kanda et al. 2016). Satisfyingly, targeting of both *ERN1* and *XBPI* was observed to decrease reporter signal in the screen (Figure 84D).

[00834] Of note, genes with biological functions unrelated or indirectly related to ER function also scored as UPR hits, some of which are distinct from functional classes seen in the analogous systematic yeast studies (Jonikas et al. 2009). Specifically, sets of genes that control general translation, transcription, and, perhaps most intriguingly, mitochondrial function were enriched among hits (Figure 84F, 90D). While disruption of the functions of these globally important genes may induce the UPR through impairment of ER function, it is also possible that such broad hits, especially the transcriptional components, represent perturbations of our reporter system. During individual testing of 257 sgRNAs targeting 152 select hit genes, Applicants confirmed that a majority induced reporter signaling from the UPRE reporter; however, for some sgRNAs, notably ones targeting the mediator transcriptional complex, Applicants found that GFP levels were also reduced, inflating the mCherry/GFP ratio (Figure 84G). Moreover, from the magnitude of the phenotypes alone, Applicants cannot distinguish among the various mechanisms for disrupting ER homeostasis because perturbation of functionally diverse genes frequently resulted

in similar levels of reporter induction whereas perturbation of genes with shared function, as in the case of ER protein targeting genes, often resulted in a range of phenotypes (Figure 84F). Cumulatively, these data highlight the ability of the pooled screens to robustly identify diverse factors important for ER homeostasis but the challenges of making systematic conclusions about gene function from such simple phenotypic readouts.

Example 27: Perturb-seq analysis of a UPR-inducing CRISPRi sub-library demonstrates functional relationships between genes and differential activation of UPR transcriptional programs in response to distinct forms of stress

[00835] Our next goal was to systematically characterize the role of these different classes of genes. Applicants hypothesized that a mechanistically diverse set of perturbations to ER homeostasis might reveal a richer picture of the underlying structure of its regulation than is possible with a small set of strong chemical perturbants (e.g., thapsigargin and tunicamycin). Applicants thus applied perturb-seq to a CRISPRi sub-library of 91 sgRNAs targeting 83 genes, including many of our strongest hits, and 2 negative controls (Figure 90C). To test platform scalability, sgRNAs were delivered via pooled transduction in this experiment. Applicants then subjected the UPR-inducing CRISPRi sub-library to profiling via perturb-seq, collecting ~65,000 cells and sequencing to an average depth of ~67,000 reads per cell, which collapsed to ~15,300 UMIs arising from ~3,700 unique genes per cell. Approximately 50,500 cells were called a containing unique sgRNA, while ~9,000 contained multiple sgRNAs (either due to double infection during our pooled transduction or encapsulation of multiple cells during emulsion generation). ~6,000 cells were unidentifiable. All of the expected sgRNAs were detected, with nearly even representation across sgRNAs that were pooled equally during viral titering (457 ± 108 cells per sgRNA, mean \pm standard deviation).

[00836] Applicants then constructed synthetic bulk expression profiles by averaging normalized expression across cells containing a given sgRNA. Hierarchical clustering of these average profiles revealed that sgRNAs targeting the same gene always clustered together (Figure 91A), and the level of knockdown of guides targeting the same gene was similar (Figure 85C). Knockdown was robust, with a median 90% depletion of the guide target (Figure 85C). Knockdown appeared as a shift of the entire distribution of expression, rather than a bifurcation into perturbed and unperturbed subpopulations (Figure 91B). To assess homogeneity, Applicants computationally assigned each cell a score measuring the deviation of its expression profile from

the distribution of expression seen in control cells. Comparing the most perturbed cells (top third) to the least perturbed (bottom third) showed a median difference in knockdown of 8% (Figure 91C). These findings establish the general robustness of the CRISPRi reagents in producing uniform knockdowns as well as the ability of the barcoding scheme to accurately assign sgRNAs to the appropriate cells. Given the similarity in phenotypes between sgRNAs targeting the same gene, in subsequent analyses Applicants grouped cells by sgRNA target, rather than by sgRNA.

[00837] At the level of bulk analysis, perturb-seq offers a simple and robust platform for highly parallel RNA-seq. Since sgRNAs can be introduced in pooled format, variability from exogenous factors is minimized, and the loss of information from shallow capture and sequencing is reduced once many cells are averaged together. These bulk profiles are a rich phenotypic fingerprint that can be used to understand how different perturbations relate to each other (Parnas et al. 2015). For example, hierarchical clustering of bulk profiles for each targeted gene revealed high precision functional clusters (apparent as lines on the diagonal in Figure 85A), consistent with known protein-protein interactions such as the oligosaccharyltransferase complex (*DDOST / DAD1 / OST4*) and Triple T complex (*TELO2 / TTI1 / TTI2*), as well as factors involved in the ubiquitylation reactions of ERAD (*SYVN1 / SEL1L*), SRP-mediated protein targeting (*SRP68 / SRPRB / SRPR*), protein trafficking (*TMED2 / TMED10*), and UFMylation (*UFML1 / UFM1 / DDRGK1*). Genes also clustered by role, including ER-Golgi trafficking (the YIPF5 / IER3IP1 complex clustering with *SCYL1* and *TMEM167A*). Interestingly, genes encoding particular translation initiation factors, the EIF2B complex (*EIF2B2 / EIF2B3 / EIF2B4*) and eIF1 α (*EIF21*), were part of a large, highly correlated cluster. Similar to inhibition of eIF1 α , EIF2B loss-of-function is known to induce ATF4 signaling (Fogli, Boespflug-Tanguy 2006). This presents the possibility that the repression of other genes within this cluster also induce ATF4. Indeed, within this large cluster Applicants observed smaller clusters of genes whose common roles are monitored by ATF4, including a set of aminoacyl tRNA synthetases (*JARS2 / SARS / FARSB / HARS / QARS / CARS / AARS / MARS / DARS / TARS*) and mitochondrial genes (*TIMM23 / SAMM50 / ATP5B / MRPL39* and *DNAJC19 / HSPA9 / IARS2*) (Bao et al. 2016, Tyynismaa et al. 2010).

[00838] An advantage of perturb-seq is that these bulk properties can then be examined at the single-cell level to better understand the phenotypes. For example, decomposing the populations

by cell cycle position (Figure 85B) revealed that some perturbations, such as many of the aminoacyl tRNA synthetases, elicited an accumulation of cells in G2. In contrast, CAD, a gene essential for pyrimidine synthesis, elicited a clear accumulation of cells in S phase.

[00839] Applicants next sought to analyze how individual hits effect activation of the different branches of the UPR. As single-cell data are prone to dropouts, scoring systems based on defined linear combinations of gene expression are prone to wide variation. Applicants thus adopted a data-driven strategy, and trained random forest regressors to score branch activation using the cells in our epistasis experiment, in which the branches are definitively separated, as training data (Methods). This scoring method had superior accuracy compared to defined gene lists and superpositions derived from ICA (Methods, Figure 91F).

[00840] Branch activation scores (Figure 85D) showed that hits from the screen elicited activation of all three branches of the UPR with clear correlations in activation among functionally related groups of genes. Perhaps most interestingly, different groups of perturbations elicited differential activation of the three branches. For example, repression of *HSPA5*, which encodes the core component of the UPR and central ER chaperone BiP, robustly activated all three branches. The amino acid synthetases activated both IRE1 α and ATF4 branches of the UPR, with the latter possibly influenced by GCN2 signaling, which is responsive to amino acid stress and, like PERK, inactivates eIF2 α . (Because our training data was experimentally obtained using PERK depletion in the presence of ER stress, Applicants score activation of the corresponding signature with a “PERK score.” However, the activity of four separate kinases that upregulate ATF4 may formally contribute to these effects). Finally, all three components of the PCC (*SEC61A1/SEC61G/SEC61B*) appeared to selectively activate only the IRE1 α branch. Comparison with alternate scoring methods and expression of branch-specific genes Applicants identified in the epistasis experiment showed good agreement with these calls (Figure 91D,E). Thus at the bulk level, our data reveal how different genetic perturbations can selectively activate the different branches of the UPR.

Example 28: Single-cell expression profiling of cells perturbed by HSPA5 depletion uncovers a bifurcated response

[00841] The observation of differential activation of the different UPR branches in bulk raises an immediate question: do single cells coordinately carry out an identical response to a given perturbation, or might the branches also operate independently at the single-cell level? To

explore this issue, Applicants examined cells depleted of BiP, where all three branches of the UPR are active.

[00842] When compared to cells containing control sgRNAs, cells transduced with *HSPA5*-targeting sgRNAs were readily distinguishable as a distinct population (Figure 86A), and had markedly different patterns of gene expression (~2,100 genes differentially expressed at $P < 0.01$). Using LRICA, Applicants decomposed these differences into sixteen independent components. Several of these reflected cell cycle effects, reflecting the altered composition seen previously (Figure 85B).

[00843] Two of the independent components (IC1 and IC2) varied substantially between control and *HSPA5*-perturbed cells (Figure 86B), and examination of the genes that most strongly affected these components showed they were driven by genes Applicants had previously identified as under UPR control. Comparing these hypothesis-free results to the branch activation scores given by our random forest regressors (Figure 86C) showed that our analysis pipeline had independently discovered a subpopulation structure within the *HSPA5*-perturbed cells in which ATF6 alone was activated (IC1 and IC2 both low), ATF6 and IRE1 α were activated (IC1 low, IC2 high), or ATF6, IRE1 α , and ATF4 were simultaneously activated (IC1 high, IC2 low). Indeed, when Applicants ordered the cells by the value of IC1 and examined the expression of UPR-induced genes (as defined in Figure 83F), the trends defining these subpopulations were apparent (Figure 86D).

[00844] Of particular note was the switch-like induction of the ATF4 branch, revealing that these cells represented a discrete subpopulation. Importantly, these differences did not reflect differences in BiP depletion, as the subpopulations with IC1 low and high (Figures 86B,D) had indistinguishable expression of *HSPA5* (Figure 86E). They were distinguishable, however, in that the ATF4-induced subpopulation had a substantially altered cell cycle distribution, with many cells accumulating in G2 compared to control cells or *HSPA5*-perturbed cells without ATF4 induced (Figure 86F). These results thus reveal that the phenotype observed at bulk level does not accurately capture the biology of single cells, and that the UPR can be executed in markedly different ways within an apparently homogeneous population.

Example 29: Gene-gene covariance analysis of perturb-seq data reveals transcriptional regulons

[00845] Perturb-seq offers substantial insights into the biology of individual cells, but also reveals the genes driving those behaviors. Figure 86D underscores a key point: correlated up- or downregulation of genes can be a signature of shared regulation (Klein et al. 2015). As a single perturb-seq experiment provides tens of thousands of data points for each gene, these coordinated motions provide a rich phenotype that can be used to identify related genes by looking for correlated patterns of expression (Figure 86G).

[00846] For example, Applicants identified 200 genes that were perturbation-induced and highly variable across the 50,500 cells in our perturb-seq experiment. When these genes were clustered based on their patterns of co-expression, obvious functional groupings appeared, including all three branches of the UPR (Figures 86H and 92A). Indeed, when Applicants clustered our defined list of UPR-induced genes (from Figure 83F) using co-expression information from the epistasis experiment or the perturb-seq experiment, Applicants uncovered similar clusterings (cophenetic correlation 0.81, compared to 0.13 when control cells were used to cluster the genes), with differences tending to involve targets with overlapping regulation by ATF6 and IRE1 α (Figure 86I, Figure 92B). The similarity between experiments indicates that much of the regulatory interactions between these genes is implicitly present in the perturb-seq experiment, and also suggests that the organization of the UPR is preserved between the very strong chemical inducers used in the epistasis experiment and the far more varied genetic perturbations used here.

[00847] A key advantage of perturb-seq is that the perturbations tend to elevate correlations above the levels induced by natural variation. These correlations can be further enhanced by a “fishing” strategy, where the population is subsampled to enrich for cells perturbed for a given group of genes. Related genes can then be discovered by how well their expression correlates with the original group (Figure 86J,K, Methods). For example, our preliminary analysis (Figure 86H) identified a five gene cluster of cholesterol biosynthesis genes with correlated expression. When Applicants confined our gene clustering analysis to the roughly ~9,000 cells most perturbed for this cluster, Applicants saw a general strengthening of correlations, and the core group now clustered with a group of genes related to cholesterol biosynthesis (Figure 86K). Similarly, when Applicants selected out the cells most perturbed for the heat shock proteins *HSPA1A* and *HSPA1B*, the genes that clustered nearby were essentially all implicated in the heat shock response based on Reactome annotation, HSF1 binding profiles from Encode, or literature

searches (Figure 86L). When Applicants expanded the analysis to look at nearby gene clusters, Applicants noticed some apparent structure in co-expression (Figure 86L). Intriguingly, the correlation structure seemed somewhat explainable in terms of the most enriched transcription factor binding sites among the gene set, with the middle group of genes for example linked to the first group by a possible regulatory overlap between HSF1 and KAT2A (Figure 86L). Though these demonstrations are not definitive proof, they suggest that gene-gene correlation information from perturb-seq, potentially constrained by other sources of information like transcription factor binding, may enable automated functional clustering of genes of unknown function.

Example 30: Preferential activation of the IRE1 α -controlled arm of the UPR in response to translocon dysfunction initiates homeostatic feedback

[00848] Among the genes targeted in the perturb-seq screen, the PCC genes (*SEC61A1*, *SEC61G*, and *SEC61B*) were perhaps the most intriguing outliers, in that they each displayed a marked preference for activation of the IRE1 branch (Figure 85D), with widespread IRE1 α activation and little or no activation of the other branches (Figure 87A, 87B, 93A). To confirm that our single-cell data were accurately calling IRE1 α activation, Applicants directly probed for XBP1 splicing using RT-PCR (Calfon et al. 2002) (Figure 87C). Depletion of all three PCC coding genes induced XBP1 splicing at levels consistent with the single-cell data and to a degree at or above that provoked by targeting *HSPA5*, a perturbation that induces all three branches of the UPR. Additionally, a time course experiment showed that depletion of PCC components led to sustained XBP1 splicing over the course of days. Expression of SSR2, a translocon-associated protein and strongly selective target of IRE1 α (Figure 87D, and cf. Figure 83F), was also maintained over the same period. These results were in marked contrast to the transient XBP1 splicing seen when cells were challenged with chemical stress (Figure 90A), which diminished on the scale of hours, consistent with previous reports (Lin et al. 2007).

[00849] The kinetics and magnitude of IRE1 activation in response to PCC gene repression varied with respect to the subunit targeted, which may reflect the fact that depletion of these factors also have variable effects on cell growth. *SEC61A1* is an essential gene with strong growth phenotypes in both CRISPRi and CRISPR cutting cell viability screens (Gilbert et al. 2014, Wang et al. 2015), however across both analyses, *SEC61B* depletion demonstrated at most a relatively mild effect, which included analysis of a CRISPRi sgRNA (*SEC61B-1*), for which

Applicants independently validated mRNA depletion (Figure 90B, 93B,C). Concordantly, Applicants observed a limited increase in *CHOP* expression upon depletion of *SEC61A1* (Figure 87D), which is indicative of ATF4 induction. The single-cell analyses nonetheless revealed that the majority of cells only induced the IRE1 α branch (Figure 87A, 93A), suggesting that the increase in *CHOP* expression might reflect general toxicity. Strikingly, depletion of *SEC61B*, which causes only mild growth defects, showed little to no *CHOP* induction over ten days of growth, while showing sustained markers of IRE1 α induction (Figure 87D). Applicants note that *SEC61B* appears to share a co-regulated promoter region with *ALG2*, a gene that functions in *N*-linked glycosylation, and that Applicants cannot formally separate the effects of repressing these genes (Figure 93B). Nonetheless, the consistent phenotypes of *SEC61A1*, *SEC61B*, and *SEC61G* repression suggest that genetic disruption of the PCC elicited selective activation of the IRE1 α branch. An alternative explanation for these observations is that depletion of PCC subunits induces general stress but impairs activation of the other branches of the UPR. ATF4 activation was intact, however, as cells transduced with sgRNAs targeting *SEC61A1*, *SEC61B*, or *SEC61G* showed robust *CHOP* induction in response to thapsigargin treatment (Figure 87E). Additionally, without stress these cells showed strong XBP1 splicing and low levels of *CHOP* as compared to cells transduced with an *HSPA5* targeting sgRNA, which demonstrated strong basal *CHOP* induction (Figure 87C,E). Collectively, these results confirm that depletion of PCC subunits selectively induces the IRE1 α branch.

[00850] Taken together with our single-cell epistasis data, these experiments suggest a selective role for the IRE1 α branch of the UPR in monitoring PCC function. Many of the strongest and most selective IRE1 α transcriptional targets in the epistasis experiment were PCC-associated genes (Figure 83F), and by qPCR analysis, Applicants confirmed upregulation of the two other PCC coding genes in response to depletion of *SEC61A1* or *SEC61B* (Figure 93B). Conversely, PCC coding and translocon-associated genes were hits in our unbiased genome-wide screen for IRE1 α activation and also showed preferential IRE1 α pathway activation in our perturb-seq analysis. Given that the translocon has been shown to selectively interact with IRE1 α (and not ATF6 or PERK) (Plumb et al. 2015, Kanda et al. 2016), these results suggest a model in which IRE1 α actively monitors the function and number of translocons and acts to increase them as needed (Figure 87F).

Example 31: DISCUSSION

[00851] Here Applicants present perturb-seq, an experimental and analytical platform that enables multiplexed profiling of many perturbations with single-cell resolution, and used it to systematically dissect the mammalian unfolded protein response (UPR). Perturb-seq is enabled by a strategy Applicants developed for robustly encoding the identity of a perturbation in an expressed transcript. Because this information is specifically encoded in the genome of cells bearing a particular perturbation, different perturbations can be interrogated together in pooled format via massively parallel droplet-based single-cell RNA-seq. This approach fills a gap in our ability to interrogate biological systems, allowing measurement of rich, unbiased phenotypes for many perturbations in one experiment.

[00852] Perturb-seq addresses a pressing problem in modern genetics: with genome-wide screens now routine (Shalem, Sanjana & Zhang 2015), strategies must be developed to rapidly and systematically understand hits. Perturb-seq enables this in two ways. First and most simply, it is a powerful platform for highly parallelized RNA-seq in homogeneous experimental conditions. Combined with our approach for combinatorial genetic perturbations, perturb-seq has the potential to interrogate tens of thousands of distinct combinations of perturbations in one pool. On the other end, the single-cell nature of the data provides deep and precise insights into how cells behave. Indeed the analytical approaches Applicants have developed provide general tools for exploiting single-cell data to identify co-regulated genes in a hypothesis-free manner, reveal discrete subpopulations that make bulk phenotypes misleading, and separate specific effects from confounding phenomena like the cell cycle.

[00853] Applicants here focused on genetic perturbations using CRISPRi, but the platform can be used to encode the identity of a wide range of perturbations or changes, such as CRISPR cutting-mediated loss of function (Wang et al. 2015, Shalem et al. 2014), gene activation (Tanenbaum et al. 2014, Konermann et al. 2015, Chavez et al. 2015), targeted mutation (Komor et al. 2016), or the identity of a particular variant in a library of systematically mutagenized proteins. Applicants have shown that CRISPRi gives strong, highly homogeneous, and simultaneous depletion of up to three targets and enables the study of essential genes. It also has some potential advantages when scaling to high-order combinations relative to CRISPR cutting, as genetic variability (e.g., partially functional alleles due to in-frame indels (Shi et al. 2015)) and non-specific toxicity due to DNA cutting (Wang et al. 2015, Aguirre et al. 2016, Munoz et al. 2016, Horlbeck et al. 2016) both increase with the number of cut sites. On the other hand,

CRISPR cutting provides complementary advantages including allowing domain-level analysis (Shi et al. 2015) and the generation of true null alleles (Shalem et al. 2014).

[00854] Perturb-seq is already a robust platform for in depth functional analysis of hit genes from screens. Scaling it further to genome-level coverage requires overcoming some obstacles, but none appear intractable. Cost per cell is a major consideration, but this will decline as technologies mature. With current techniques (Zheng et al. 2016), ~50,000 cells can be profiled in ten minutes. Additionally, perturb-seq will enable higher loading, which usually has the undesirable effect of increasing the frequency at which multiple cells are encapsulated per droplet, because our guide barcodes allow computational removal of cell doublets. Sequencing costs are high, but can be mitigated at the library preparation stage through specific amplification of select targets, as demonstrated with our barcode capture, or by depletion of uninteresting high abundance genes (Gu et al. 2016). A more subtle point is that provirus recombination during lentivirus transduction can potentially scramble barcode identities during standard methods of pooled library delivery (Sack et al. 2016). Applicants took careful steps to avoid these problems (see Methods) and expect that altered library preparation protocols will circumvent this issue.

[00855] By far the biggest barrier Applicants anticipate is on the analytical side. Perturb-seq generates massive amounts of data, but these data are intrinsically noisy due both to the low capture efficiencies of single-cell RNA-seq protocols and more fundamentally, the intrinsic stochasticity of transcription at the single-cell level. Rigorously extracting meaningful conclusions will require new tools and algorithms (Satija et al. 2015, Heimberg et al. 2016, Trapnell 2015, Grün, van Oudenaarden 2015), as well as new statistical approaches for assessing the significance of claims. Applicants have made some progress on this front, and showed how single-cell data can be used for example to objectively decouple the branches of the UPR, uncover subtle subpopulations within cells of the same type, and infer larger programs of gene expression using correlated expression. Along with previous successes (Macosko et al. 2015, Klein et al. 2015, Jaitin et al. 2014), and the novel analytical approaches presented in an accompanying manuscript (Dixit et al., co-submitted manuscript), Applicants anticipate that large-scale analyses of single cell behavior will enable systematic discovery of the complex and overlapping regulatory programs at work within cells. As perturb-seq scales, it will present computational challenges on the same scale as other “big data” fields.

[00856] Our experiments also provide important insights into our understanding of the logic used by the mammalian UPR to sense and respond appropriately to the diversity of challenges faced by the ER. A central question in the field of ER homeostasis is why metazoan cells have evolved to have three independent and mechanistically distinct sensors of protein misfolding in the ER. As expected from previous work (Acosta-Alvear et al. 2007, Han et al. 2013, Lee, Iwakoshi & Glimcher 2003, Shoulders et al. 2013), epistasis analysis using combinatorial depletion of the three main regulators of the UPR revealed both distinct and overlapping programs of gene expression controlled by PERK, ATF6, and IRE1 α . The clean divisions Applicants observed at the single-cell level highlight the importance of performing these experiments in a pooled format, as sensitive phenotypes like stress responses might otherwise be overwhelmed by experimental variation. PERK/ATF4 had the largest regulon in our data, while ATF6 and IRE1 α had more overlapping targets, with IRE1 having more unique targets such as the translocon. Additionally, Applicants observed that the three UPR branches were able to operate independently both at the bulk level, as different genetic perturbations showed differential activation of the three branches, and at the single-cell level.

[00857] Our genome-wide screen identified a broad diversity of genetic perturbations that activate signaling from the IRE1 α branch of the UPR, including some categories of genes that were not expected from analogous yeast screens (Jonikas et al. 2009). Subjecting these hits to perturb-seq showed that the screen in fact hit all three branches of the UPR, and that genes with similar functional roles induced the UPR in similar ways. Intriguingly, Applicants also found that the signatures of the three branches of the UPR were implicit in the correlations seen in the perturb-seq experiment, suggesting that in principle the hypothesis-driven epistasis experiment could be dispensable. Thus, perturb-seq could enable identification and characterization of novel pathways for which such *a priori* information is not available. The single-cell analysis also revealed a remarkable bifurcation of UPR branch activation within cells subject to identical perturbations. In particular, depletion of *HSPA5* resulted in some cells primarily inducing ATF6, others both IRE1 α and ATF6, and a distinct subset inducing all three UPR branches simultaneously. This example illustrates the utility of single-cell data: bulk RNA-seq would in this case describe an average state that no cell actually occupies, just as it obscures the effect of other factors like the cell cycle. As all *HSPA5*-perturbed cells were treated identically, the cause of such marked differences in behavior remains in question. The simplest answer, variability in

knockdown efficiency, does not appear to play a role, though the time course of *HSPA5* depletion might. But this possibility suggests a real, interesting biological hypothesis, namely that there is a temporal ordering to the activation of the three pathways in *HSPA5*-depleted cells (Lin et al. 2007).

[00858] Applicants also found that different perturbations show clear differences in their propensity to activate the different branches of the UPR. Particularly intriguing in this regard were the three components of the protein-conducting channel encoded by *SEC61A1*, *SEC61G*, and *SEC61B*, repression of which selectively activated the IRE1 α branch of the UPR. These results are notable in light of recent studies suggesting that IRE1 α acts in proximity to the translocon (Plumb et al. 2015, Kanda et al. 2016). Given that Applicants and others have observed regulation of translocon expression to be uniquely under IRE1 α control (Shoulders et al. 2013), this suggests a feedback model in which IRE1 α actively monitors the state of translocation. Interestingly, isolated IRE1 α induction would enable repair to or upregulation of the translocation machinery without a broader UPR induction, preventing potentially inappropriate responses including cell death (Hetz 2012). Additionally, further study of selective UPR induction using genetic perturbations, as Applicants have shown here, could also inform efforts to pharmacologically manipulate the UPR to enhance specific repair and prosurvival pathways (Shoulders et al. 2013). More broadly, this provides a concrete example of how the three distinct sensors of ER homeostasis allow the cell to monitor distinct types of ER stress and provide nuanced responses tailored to the specific challenge.

[00859] Applicants believe that their study of the mammalian UPR serves as a blueprint for the study of complex and overlapping transcriptional networks, in which a primary genome-wide screen serves as the input to more detailed analysis via perturb-seq. Applicants independently uncovered both well-known features of the UPR and new and unexpected findings. As the technology scales, it should be possible to conduct saturating analyses of every major pathway. Our success here, and the parallel success in understanding dendritic cell activation described herein, speak well to potential of the perturb-seq approach to become a standard strategy for understanding regulatory interactions in the cell.

Materials and Methods

[00860] Pool construction and transduction. A lentivirus backbone was constructed containing antiparallel cassettes of a PolIII promoter for sgRNA construction and a PolII promoter for expression of a selection marker and a polyadenylated RNA barcode. This construction avoided early termination of the viral transcript. A diverse library of RNA barcodes was inserted (approximately 1,000). Finally, sgRNAs were inserted. Association between RNA barcodes and sgRNAs was performed using Sanger sequencing. Lentivirus was made using 293T cells transfected with pBA439, and pMD2.G (Addgene 12259) at a 10:10:1 ratio, using Lipofectamine LTX and plus reagents according to the manufacturer's instructions.

[00861] For the mouse experiments, BMDCs were isolated from 6- to 8-week-old constitutive Cas9-expressing female mice and used as described previously (Platt et al., 2014). Cells were infected with the pooled lentiviral library at an MOI of approximately 0.2 at day 2. At day 9, BMDCs were stimulated with 100 ng/ml LPS. Cells were harvested for library preparation at 0hr and 3hr post stimulation.

[00862] For the human experiments, k562 transduced with Cas9 were purified by high GFP expression. These cells were transduced with the pooled lentiviral library at an MOI of approximately 0.2 for most experiments and an MOI of 1 for the high MOI experiment. Cells were harvested for library preparation 7 days post transduction for most experiments and 14 days post transduction for the second time point promoters experiment.

[00863] Library preparation. After cells were collected they were processed into single cell RNA-seq libraries following the manufacturer's protocol of the commercial 10X Chromium system.

[00864] Bone marrow derived dendritic cells (BMDCs). All animal protocols were reviewed and approved by the MIT / Whitehead Institute / Broad Institute Committee on Animal Care (CAC protocol 0609-058-12). To obtain sufficient number of cells, Applicants implemented a modified version of the DCs isolation protocol as previously described (Amit et al. 2009; Chevrier et al. 2011; Garber et al. 2012; Lutz et al. 1999; Rabani et al. 2011). Briefly, for all CRISPR knockout experiments six- to eight-week old constitutive Cas9-expressing female mice were used as described previously (Platt et al. 2014). RPMI medium (Invitrogen) supplemented with 10% heat inactivated FBS (Invitrogen), β -mercaptoethanol (50 μ M, Invitrogen), L-glutamine (2mM, VWR), penicillin/streptomycin (100U/ml, VWR), MEM non-essential amino acids (1X, VWR), HEPES (10mM, VWR), sodium pyruvate (1mM, VWR), and

GM-CSF (20 ng/ml; Peprotech) was used throughout the study. At day 0, cells were collected from femora and tibiae and plated in 100mm non tissue culture treated plastic dishes using 10ml medium per plate at concentration of 2×10^5 /ml. At day 2, cells were fed with another 10ml of medium per dish. At day 5, 12ml of the medium were carefully removed (to avoid removal of cells) and 10ml of fresh medium were added back to the original dish. Cells were fed with another 5ml medium at day 7. At day 8, all non-adherent and loosely bound cells were collected and harvested by centrifugation. Cells were then re-suspended with medium, plated at a concentration of 10×10^6 cells in 10ml medium per 100mm dish. At day 9, cells were stimulated for with LPS (100ng/ml, rough, ultrapure *E. coli* K12 strain, Invitrogen) and harvested. Cells were always plated at concentration of 2×10^5 /ml at day 0. Cells were harvested post stimulation after 0hr or 3hr and cells from cultures that contained 10% BFP positive cells were sorted for BFP+ and GFP+ (contain CAS9).

[00865] **Cell Culture.** Human K562 cells were purchased; murine BMDC were isolated from 6- to 8-week-old constitutive Cas9-expressing female mice.

[00866] K562 cells were grown to a confluence of 30-60% and spun down at 300 x g for 5 min. The supernatant was removed, and cells were resuspended in 1 mL of 1x PBS + 0.2% BSA (Sigma cat #A8806) and re-spun at 300 x g for 3 min. The supernatant was again removed, and the cells re- suspended in 1 mL of 1x PBS, passed through a 40-micron cell strainer (Falcon, VWR cat #21008- 949), and counted. For 10X, cells were diluted to the final concentration in 1x PBS + 200 µg/mL BSA (NEB, cat # B9000S).

[00867] **K562 cell cultures.** Applicants used transgenic K562 cells constitutively expressing Cas9. K562 cells were transduced using several titers of virus and cells were spin infected in 2000 rpm for 30 min. For the low MOI experiment Applicants used cultures that contained 10% BFP+ and for the high MOI 50% BFP+. Cells were grown in RPMI Medium 1640 + GlutaMAX (ThermoFisher) + 10% heat inactivated FBS (Invitrogen). Cells were grown to a confluence of 30-60% and spun down at 300 x g for 5 min. The supernatant was removed, and cells were suspended in 5 mL of 1x PBS + 0.2% BSA (Sigma cat #A8806) for sorting: BFP+ and GFP+ (CAS9 expressing) cells were sorted. Cells were harvested for library preparation 7 days post transduction for most experiments and 13 days post transduction for the second time point of the TF pool experiment. After sorting BFP+ GFP+ cells passed through a 40-micron cell strainer (Falcon, VWR cat #21008- 949), wash twice and counted.

[00868] Library preparation. After cells were collected, they were processed into single cell RNA-seq libraries following the manufacturer's protocol of the commercial 10X Chromium system.

[00869] Single cell library preparation. Prior to analysis, cells were diluted to the final concentration in 1x PBS + 200 µg/mL BSA (NEB, cat # B9000S). Sorted cells (BMDCs or K562 cells) were loaded on the 10X Chromium system (Zheng et al., 2016) (8,000 cells/channel) and single cell RNA-seq libraries were generated following the manufacturer's instructions.

[00870] Following WTA, a fraction of the WTA was used to amplify sgRNA barcode using a dial-out PCR strategy with the primer sequences below (the full primer sequence is a concatenation of the columns). The template material was approximately 5ng of WTA libraries. 25 cycles of PCR were performed using one of the dial-out primers below with the P7 Illumina reverse primer.

[00871] Preparation of RNA barcode sub-libraries. In order obtain high confidence information about the identities of which sgRNAs were present in each cell, a dial-out PCR strategy was performed using the primer sequences below (the full primer sequence is a concatenation of the columns). The template material was approximately 5ng of whole transcriptome amplification libraries obtained after a 10X library, 25 cycles of PCR was performed using one of the dial-out primers below with the a P7 Illumina primer reverse primer.

[00872] Primer sequences:

P5	Barcode	R1	Primer (BFP location of pBA439)
AATGATACGGCGACCACCGAGATC TACAC	NNNNNN NN	TCGTCGGCAGCGTCAGATGTGTATAAGA GACAG	TAGCAAAC TGGGGCAC AAGC
AATGATACGGCGACCACCGAGATC TACAC	TCGCCTTA	TCGTCGGCAGCGTCAGATGTGTATAAGA GACAG	TAGCAAAC TGGGGCAC AAGC
AATGATACGGCGACCACCGAGATC TACAC	CTAGTAC G	TCGTCGGCAGCGTCAGATGTGTATAAGA GACAG	TAGCAAAC TGGGGCAC AAGC
AATGATACGGCGACCACCGAGATC TACAC	TTCTGCCT	TCGTCGGCAGCGTCAGATGTGTATAAGA GACAG	TAGCAAAC TGGGGCAC AAGC
AATGATACGGCGACCACCGAGATC TACAC	GCTCAGG A	TCGTCGGCAGCGTCAGATGTGTATAAGA GACAG	TAGCAAAC TGGGGCAC AAGC
AATGATACGGCGACCACCGAGATC TACAC	CATGCCTA	TCGTCGGCAGCGTCAGATGTGTATAAGA GACAG	TAGCAAAC TGGGGCAC AAGC

AATGATACGGCGACCACCGAGATC TACAC	G TAGAGA G	TCGTCGGCAGCGTCAGATGTGTATAAGA GACAG	TAGCAAACCTGGGGCAC AAGC
-----------------------------------	------------------	---------------------------------------	---------------------------

[00873] P7 Illumina Reverse Primer:

[00874] CAAGCAGAAGACGGCATAACGAGAT

[00875] The PCR protocol is as follows:

	1 rxn
Q5 2X master mix	25
P7 primer @10uM	1.25
pBA439_rev @ 10uM	1.25
Template (5ng total)	x
water (qs up to final rxn of 50ul)	y
total	50

Temp	time	
98 C	10s	
98 C	2s	repeat for 25 cycles
65 C	5s	
72 C	10s	
72 C	1 min	

[00876] The resulting PCR product was purified with Agencourt AMPure XP beads at 0.7x and sequenced.

[00877] **Read alignment and Generation of Expression Matrix.** A digital expression matrix was obtained for each experiment using 10X's Cellranger pipeline with default parameters. The software uses STAR for alignment. All subsequence analysis information is available and maintained in current form in the following Git repository and documentation available through ReadtheDocs.

[00878] **Alignment of initial abundances.** Applicants compiled a dictionary of RNA barcodes associated with sgRNA from Sanger sequencing each part of the construct separately. Applicants devised an alternative cloning protocol for larger libraries to enable NGS pairing of

sgRNA/RNA barcodes, but did not use it for any of the results in this paper. The initial abundance of each plasmid in the pool was quantified using NGS of the RNA barcode. The RNA barcode/ sgRNA dictionary was used to convert the readout into a relative abundance estimate of sgRNA in the initial pool.

[00879] Alignment of Cell barcode / RNA barcodes libraries. In order to associate cell barcodes with guides, the same RNA barcode / sgRNA dictionary was used. Then NGS sequencing may have contained a complex library of molecules such as a full RNA-seq library reads. Paired end reads containing a cell barcode UMI on one side and an RNA barcode on another side were isolated and collapsed into unique molecules by first demultiplexing a sequencing run using *bcl2fastq2* with the followign options *--create-fastq-for-index-reads --barcode-mismatches 1 --no-lane-splitting --mask-short-adapter- reads 5 --minimum-trimmed-read-length 5* using a sample-sheet containing one line with polyG tracts in the index read columns (see below). The resulting Undetermined fastq files were split using kentools into two folders called split1 and split2 containing chunked R1 and R2 reads respectively.

Sample_ID	Sample_Name	I7_Index_ID	index	I5_Index_ID	index2
1	placeholder	silly1	GGGGGGGGGGGGGGG	silly2	GGGGGGGGG

[00880] Then the reads were concatenated and RNA barcode reads isolated using the following UNIX commands.

```
paste ./split1/$f1 ./split2/$f2 | awk '{if (NR%4==1) printf("$2"\t"); else if (NR%4==2)printf("$0"\n")}' | grep GGCACAAGCTTAATTAAGAATT
```

[00881] with GGCACAAGCTTAATTAAGAATT corresponding to a constant sequence within the RNA barcode, and \$f1 and \$f2 corresponding to input arguments containing the read1_splitfile and read2_splitfile respectively and creating an output split file.

[00882] All split files containing a particular index barcode are finally concatenated and collapsed to unique molecules using the following command.

```
cat *${inputbc}.txt | sort | uniq -c | sort -k1,1g | awk '{print $1"\t"$2"\t"$3"-"$4}'
```

[00883] the result file format looks as follows (grey highlighting for actual barcode portion of RNA barcode read):

2791:N:0:GGTGATACCTCATT+TCGCATAA
CGCAAAC**TGGGGCACAAGCTTAATTAAGAATT**CGATCAACGCAGAGACGGCCTAG-
GCGTATAAGT

2821:N:0:TCGAGCCTTATGGC+TCGCATAA
CGCAAAC**TGGGGCACAAGCTTAATTAAGAATT**GCTTGACTCGTTAGCGAGCCTAG-
CACGATACCG

2861:N:0:ACTCGAGAGTTCGA+TCGCATAA
CGCAAAC**TGGGGCACAAGCTTAATTAAGAATT**CGATCAACGCAGAGACGGCCTAG-
GTCTAATGAA

2941:N:0:GCACGTCTACTAGC+TCGCATAA
CGCAAAC**TGGGGCACAAGCTTAATTAAGAATT**C**TAACTCAGCGACTGGAGCCTAG-**
CATGTGCCCG

2971:N:0:ATAGATTGTCCGAA+TCGCATAA
CGCAAAC**TGGGGCACAAGCTTAATTAAGAATT**GCTTGACTCGTTAGCGAGCCTAG-
AAAAGGATGG

2991:N:0:GAGCAGGAGCTATG+TCGCATAA
CGCAAAC**TGGGGCACAAGCTTAATTAAGAATT**AAACCCTCACTGCCGACGCCTAG-
ACCTGTTACG

3331:N:0:GGAGGCCTGTTACG+TCGCATAA
CGCAAAC**TGGGGCACAAGCTTAATTAAGAATT**AGGGCTTGCAGTGCACGGCCTAG-
AACGTCAAGA

3391:N:0:CGACTCACGTTCAG+TCGCATAA
CGCAAAC**TGGGGCACAAGCTTAATTAAGAATT**CGATCAACGCAGAGACGGCCTAG-
CAGACTGGGC

[00884] Counts 1:N:0:Cellbarcode+Index/SampleBarcode RNAbarcoderead-UMI

[00885] The resulting file along with the preassociated RNA barcode, sgRNA dictionary is parsed using a custom python script to create a new table of probability estimates of which sgRNA are present in each cell. The probability estimate is thresholded to create a dictionary of which cell barcodes contain which sgRNAs.

[00886] **Fit of distribution of guides per cell.** The distribution of guides per cell was fit using a maximum likelihood approach.

[00887] **Linear model.** To fit the linear model Applicants compiled our covariate matrix **X** and our expression matrix (or one of the continuous covariates; as done for some assignments of cell states; below) as our matrix **Y**. Applicants then used the Python implementation of elastic net with the following parameters, to fit our model.

```
sklearn.linear_model.ElasticNet(l1_ratio=0.5,alpha=0.0005,max_iter=10000)
```

[00888] **Alternating descent fit of perturbation probability.** To account for the contribution of unperturbed cells in the population containing a particular sgRNA, Applicants constructed an approach in which the presence of sgRNA in a given cell was converted into a probability measure of that sgRNA having a phenotypic effect on the cell, as follows:

First, Applicants fit using the multivariate regression on
 $Y=X$

$$LL(y) = \log \left(\frac{\frac{1}{2\pi\sigma} e^{-\frac{(y-\mu_1)^2}{2\sigma}}}{\frac{1}{2\pi\sigma} e^{-\frac{(y-\mu_0)^2}{2\sigma}}} \right) = \frac{(y - \mu_0)^2 - (y - \mu_1)^2}{2\sigma}$$

$$P(X_j = 1|Y) = \frac{P(Y|X_j = 1)P(X_j = 1)}{P(Y)}$$

$$\frac{P(Y|X_j = 1)P(X_j = 1)}{P(Y|X_j = 0)P(X_j = 0) + P(Y|X_j = 1)P(X_j = 1)}$$

$$\frac{1}{1 + e^{-LL(Y)}}$$

[00889] Next, Applicants evaluate the fit with the guide covariate set to 0, X_0

$$P(X_i = 1) = \text{logistic}(\sum_i [Y_{ij} - X_o \beta_i]^2 - [Y_{ij} - \hat{Y}_{ij}]^2); \text{ where } \text{logistic}(x) = \frac{1}{1 + e^{-x}}$$

[00890] Finally, Applicants use the new covariate matrix X_{sg} to recompute β_{sg} .

[00891] **Significance testing for coefficients of linear model.** Applicants devised a permutation strategy to obtain an empirical null distribution of the coefficients associated with our sgRNA effects. Specifically, Applicants randomized the guide assignments to cells (such that co-occurrence between guides was preserved) and the linear model was recomputed with all other covariates being held constant. Applicants repeated this ten times. Applicants noticed that three significant as-yet latent factors impacted the empirical null distribution of coefficients: (1) the mean expression level of a gene; (2) the variance in expression of a gene; and (3) the number of cells a particular sgRNA was present in.

[00892] To control for these factors when assessing significance, each empirical null coefficient's value was assigned a point in 4-space: [Gene mean, gene variance, number of cells, value]. Applicants then estimated the multivariate density using a binning approach (*np.histogramdd* in Python). True nonzero coefficients were evaluated for significance relative to a matched set of bins (to create an empirical conditional cumulative probability distribution) conditioned on the first three factors.

[00893] A less stringent null was generated by obtaining a permuted distribution of coefficients based on a permuted cell to guide assignment, but only calculating the significance on a per guide basis without consideration of the mean expression level or variance of the gene being considered.

[00894] In both cases, Applicants used a Benjamini-Hochberg procedure to control for multiple hypothesis testing.

[00895] **Residuals analysis.** To determine the marginal effect of each covariate in explaining the observed gene expression variation, Applicants estimated the model R^2 by cross-validated (trained on 80% of the data and tested against 20%) for the addition of each of the covariates.

[00896] To determine the extent to which our covariates explained the major axes of variation in our data, Applicants decomposed the residuals using the same randomized PCA approach described in the **Definition of Cell States** section. Two major metrics were evaluated: (1) the eigenvalue distribution, and (2) the extent to which the top loadings were enriched for biological terms.

[00897] **Cross validated R^2 .** To estimate the generalizability of the model, Applicants determined a cross-validated R^2 by training our model on 80% of our data and determining the fit on the remaining 20%.

[00898] **Definition of cell states.** Cell states were defined using cells that did not have sgRNAs introduced. Starting from an expression matrix Y , variable genes are selected based on fitting a nonparametric loess regression (using a moving window of 25 percent of the data) to the relationship between the average expression of a gene and its respective coefficient of variation (after normalizing each cell for complexity). Genes with high residuals (i.e more variable than genes at comparable expression levels) were selected (approximately 1000 genes).

[00899] Next, the cell normalized expression matrix (the sum of the number of transcripts is renormalized to 10,000) and then log2 transformed with a pseudocount of 1. The final normalization is Z-transforming the genes of the matrix. Randomized PCA is performed on the Z transformed expression matrix using Facebook's implementation through the python package *fbpca* retaining the top 50 components.

[00900] A combination of the elbow method looking at the eigenvalue gap of each component, GO enrichment of each component using jackstraw (cite Storey paper), and a PC robustness analysis (in which increasing amounts of random noise is added to the data and the stability of each principal component with respect to the original components is evaluated) was used to determine the number of principal components (which in general was approximately 10)

[00901] A two dimensional visualization of the relationship between cells was performed using TSNE. Clustering was performed using Infomap with k refined so that slightly more clusters are created than one would expect. Clusters are subsequently merged in an iterative fashion that have fewer than 100 differentially expressed genes between them.

[00902] Differential expression is evaluated using a Welch's t-test on the Z-transformed values between each cluster and the rest of the cells.

[00903] A Benjamini-Hochberg procedure was used to control for multiple hypothesis testing.

[00904] The clusters are evaluated for GO enrichment using FDR corrected p-values (see Interpretation of Results).

[00905] **Relation of perturbed cells to unperturbed states.** In order to define the relationship between the cell states in the unperturbed cells and the perturbed cells, applicants projected the perturbed cells onto the same significant principal component vectors derived from the unperturbed cells. The projection onto these components was used as a covariate by itself, especially with the K562 cells, where the major axes of variation such as cell cycle describe somewhat continuous processes. Alternatively, for the PBMDs, discrete cell types are readily discernable. As such, a random forest classifier was trained using class labels obtained by the merged Infomap clusters with features consisting of PC scores. from sklearn.ensemble import RandomForestClassifier

```
clf = RandomForestClassifier(n_estimators=100,n_jobs=-1,oob_score=True,class_weight='balanced')
```

[00906] The out of bag probability estimates were used to generate ROC curves to determine the sensitivity and specificity of classification per cluster.

[00907] Finally, the random forest was applied to the projected PC scores of the perturbed cells to obtain class membership estimates.

[00908] **Testing for significant association of guides with cell states.** Applicants used several tests to assess guide-state associations and our reported results are coherent across the tests.

[00909] In a continuous test Applicants used the Wilcoxon-ranksum test to examine if the cells with a specific guide have significantly higher or lower cell-state scores compared to the cells with non-targeting or intergenic guides. The test is performed separately for each batch (to avoid batch-effects), and then the p-values are merged using Fisher's combined probability test.

[00910] In a discrete test, each cell is assigned with a single state (the one it scored highest in). Cells with low cell-state scores (<0.5) are classified as unclustered. For each guide-state pair <g,s> Applicants examine if the number of cells with guide g in state s is significantly larger or smaller than expected (using a binomial test). The expected is based on the probability of the cell-state in the cells with intergenic/non-targeting guides.

[00911] A guide-cell association is defined as significant if its p-value <0.05 after correcting for multiple hypotheses testing using the Benjamini–Hochberg procedure.

[00912] **Tests of sgRNA effect on Covariate.** In order to evaluate the effect of a sgRNA on a covariate such as genes detected, transcripts detected, or cell state, an Welch's t-test is performed between the covariate and the set of control sgRNAs.

[00913] **Expectation Maximization.** In order to account for the contribution of unperturbed cells in the population containing a particular sgRNA, an expectation maximization approach was constructed in which the presence of sgRNA in a given cell was converted into a probability measure. The approach is as follows:

[00914] First, fit β using the multivariate regression on

$$[00915] \quad Y = X\beta$$

[00916] Then evaluate the fit with guide covariate set to 0 X_o

$$[00917] \quad P(X_j = 1) = \text{logistic}(\sum_i [Y_{ij} - X_{oi}\beta]^2 - [Y_{ij} - \bar{Y}_{ij}]^2); \quad \text{where}$$

$$\text{logistic}(x) = \frac{1}{1 + e^{-x}}$$

Finally, the new covariate matrix X_M is used to recompute β_M

[00918] **Significance Testing for Coefficients of Linear Model.** A permutation strategy was devised to get an empirical null distribution of the coefficients associated with our sgRNA effects. Specifically, the guide assignments to cells were randomized (such that co-occurrence between guides was preserved) and the linear model was recomputed with all other covariates being held constant ten times. Applicants noticed that three significant as-yet latent factors impacted the empirical null distribution of coefficients: 1) the mean expression level of a gene 2) the variance in expression of a gene 3) the number of cells a particular sgRNA was present in.

[00919] In order to control for these factors in assessing significance, each empirical null coefficient's value was assigned a point in 4-space: [Gene mean, Gene variance, number of cells, value] A multivariate density estimation was subsequently performed using a binning approach (*np.histogramdd* in Python). True coefficients were evaluated for significance relative to a matched set of bins (to create an empirical conditional cumulative probability distribution) conditioned on the first three factors.

[00920] A Benjamini-Hochberg procedure was used to control for multiple hypothesis testing.

[00921] **Residuals Analysis.** In order to determine the marginal effect of each covariate in explaining the observed gene expression variation a cross-validated estimation (trained on 80%

of the data and tested against 20%) of the model R^2 was obtained for the addition of each of the covariates.

[00922] In order to determine the extent to which our covariates explained the major axes of variation in our data, applicants decomposed the residuals using the same randomized PCA approach described in the **Definition of Cell States** section. Two major metrics were evaluated: 1) the eigenvalue distribution 2) the extent to which the top loadings were enriched for biological terms using the aforementioned approach.

[00923] **Interpretation of Results.** GO analysis was performed using the *goatools* python package with an FDR threshold of 0.1.

[00924] **Power Analysis / Design of Experiments.** Power analysis was performed to determine as a function of observed effect size and baseline expression of a gene, how many cells are required to observe a signal.

[00925] **Fitness effects of sgRNAs.** Applicants used several approaches to assess the fitness effects of sgRNAs and all yielded qualitatively comparable results.

[00926] In the most straightforward approach, Applicants obtained estimates of the initial abundances of each sgRNA in the pool. The initial abundance of each plasmid in the pool was quantified using NGS of the RNA barcode. The RNA barcode / sgRNA dictionary was used to convert the readout into a relative abundance estimate of sgRNA in the initial pool. Then, Applicants calculated the fold change of the observed abundance of cells containing a particular sgRNA compared to its respective abundance in the original pool.

[00927] In an alternative approach, Applicants developed two complementary Bayesian probabilistic models that compute the expected probability of each guide in the resulting cell population. The first model assumes that none of the guides affect cell fitness and is suitable in situations where the perturbations are unlikely to effect cell survival or proliferations, as in BMDC TFs. The expected guide frequency is modeled as a function of the MOI, the guide frequency in the initial library, and the guide detection rate, such that multiple infections are considered to be independent Poisson processes. Applicants applied this model and examined if the observed guide frequencies deviate from the expected ones. Indeed, in K562 cells the observed guide distribution deviates significantly from the expected distribution, indicating that certain perturbations affect cell fitness. Applicants then extended the probabilistic model to account for these effects by introducing fitness parameters, which are set by solving a maximum

likelihood optimization problem. In this setting Applicants model the fitness effects of multiple perturbations as additive, and identify perturbations that show different co-occurrence frequencies as indicators of genetic interactions.

[00928] In both Bayesian models the infection, selection, and detection processes are modeled such that: (1) the guide detection rate, denoted as d , is the same for all guides, (2) there are no falsely detected guides, (3) the infection with different guides is independent and, (4) due to BFP selection, denoted as B , all cells have at least one guide by the end of the screen. Let M, f_g , and d_g denote the MOI, the frequency of guide g in the initial pool of the screen, and the event of detecting guide g , respectively. Applicants model the infection as a Poisson process, such that i_g denotes the number of times the cell was infected with guide g , $i_g \sim \text{Poisson}(f_g \cdot M)$. Let $I_{c,g}$, N_0 , and N_g be the indicator that cell c was detected with guide g in the screen, the number of cells in which no guide was detected, and the number of cells in which guide g was detected, respectively.

[00929] To find the most likely parameters and compute the expected guide probabilities without fitness effects Applicants solve the following maximum likelihood optimization problem.

[00930]

$$L(M, p_d) = \prod_c P(d_1, \dots, d_k | B)$$

[00931]

[00932]

$$P(d_1, \dots, d_k | B) = P(B | d_1, \dots, d_k) P(d_1, \dots, d_k) = \begin{cases} P(d_1, \dots, d_k) & \text{if } \sum_g d_g > 0 \\ P(B | \forall_g d_g = 0) P(\forall_g d_g = 0) & \text{otherwise} \end{cases}$$

[00933]

[00934]

$$L(M, p_d) = \prod_c \left[P(B | \forall_g d_g = 0)^{I_{c,0}} \prod_g P(d_g = 1)^{I_{c,g}} \cdot P(d_g = 0)^{1 - I_{c,g}} \right] = P(B | \forall_g d_g = 0)^{N_0} \left[\prod_g P(d_g = 1)^{N_g} \cdot P(d_g = 0)^{N - N_g} \right] = (1 - e^{-M(1-p_d)})^{N_0} \left[\prod_g P(d_g = 1)^{N_g} \cdot P(d_g = 0)^{N - N_g} \right]$$

[00936]

[00937] The log likelihood is

[00938]

[00939]
$$\log[L(M, p_d)] = N_0 \log(1 - e^{-M(1-p_d)}) + \sum_g \{N_g \log[P(d_g = 1)] + (N - N_g) \log[P(d_g = 0)]\} =$$

[00940]
$$N_0 \log(1 - e^{-M(1-p_d)}) + \sum_g \{N_g \log(1 - e^{-f_g p_d M}) - (N - N_g) f_g p_d M\}$$

[00941]

[00942] Below Applicants provide the equations used to compute the maximum likelihood function above.

[00943]

[00944] (1.1) $P(d_g = 0) = e^{-f_g M p_d}$

[00945] $P(d_g = 0 | I_g) = (1 - p_d)^{I_g} = q^{I_g}$, $P(d_g = 0) = \sum_n q^n P(I_g = n) = E[q^{I_g}]$

[00946]

[00947] Based on the moment-generating function

[00948] $E[e^{tX}] = e^{\lambda(e^t - 1)}$, where $X \sim \text{Poisson}(\lambda)$ and $t \in \mathbb{R}$

[00949] $t = \ln q, \lambda = f_g \cdot M, E[q^{I_g}] = E[e^{I_g \ln q}] = e^{f_g M (q - 1)} = e^{-f_g M p_d}$

[00950]

[00951] (1.2) $P(B) = 1 - e^{-M}$

[00952] $P(B) = 1 - P(\neg B)$

[00953]
$$P(\neg B) = P(\forall_g I_g = 0) = \prod_g P(I_g = 0) = \prod_g e^{-M f_g} = e^{-M \sum_g f_g} = e^{-M}$$

[00954]

[00955] (1.3) $P(B | \forall_g d_g = 0) = 1 - e^{-M(1-p_d)}$

[00956] $P(B | \forall_g d_g = 0) = 1 - P(\neg B | \forall_g d_g = 0)$

[00957] $P(\neg B | \forall_g d_g = 0) = P(\forall_g I_g = 0 | \forall_g d_g = 0) = \prod_g P(I_g = 0 | d_g = 0) = \prod_g \frac{P(d_g = 0 | I_g = 0) P(I_g = 0)}{P(d_g = 0)}$

$= \prod_g \frac{1 - e^{-M f_g}}{1 - e^{-f_g M p_d}} = \prod_g e^{-M f_g (1 - p_d)} = e^{-M(1 - p_d) \sum_g f_g} = e^{-M(1 - p_d)}$

[00958]

[00959] To account for fitness effects Applicants introduce a fitness parameter for each guide, denoted as δ_g , and define the fitness normalization factor Z , such that

[00960]
$$Z = E_{P(I_g)}[\delta_g] = \delta_0 \prod_g [(e^{-f_g M} (1 - \delta_g) + \delta_g)] = \delta_0 \prod_g \alpha_g$$

$$\begin{aligned}
 [00961] \quad Z &= E_{P(i_g)}[\delta_g] = \sum_{(n_2, \dots, n_k)} P(i_1 = n_1, \dots, i_k = n_k) \delta_0 \delta_1^{i_1 > 0} \dots \delta_k^{i_k > 0} \\
 &= \delta_0 \sum_{(n_2, \dots, n_k)} \left[P(i_1 = 0, i_2 = n_2, \dots, i_k = n_k) \delta_2^{i_2 > 0} \dots \delta_k^{i_k > 0} \right. \\
 [00962] \quad &\quad \left. + P(i_1 > 0, i_2 = n_2, \dots, i_k = n_k) \delta_1 \delta_2^{i_2 > 0} \dots \delta_k^{i_k > 0} \right] = \\
 [00963] \quad &
 \end{aligned}$$

$$[00964] \quad \delta_0 \sum_{(n_2, \dots, n_k)} (P(i_1 = 0) + P(i_1 > 0) \delta_1) \prod_{2 \leq g \leq k} P(i_g = n_g) \delta_1^{i_g > 0} =$$

$$[00965] \quad \delta_0 \prod_g (P(i_g = 0) + P(i_g > 0) \delta_g) =$$

$$[00966] \quad \delta_0 \prod_g (e^{-f_g M} (1 - \delta_g) + \delta_g) = \delta_0 \prod_g \alpha_g$$

$$[00967] \quad \text{Where } \alpha_g = e^{-f_g M} (1 - \delta_g) + \delta_g$$

[00968]

[00969] Applicants adjust the probability that a cell that was infected i_g times with guide g to account for the fitness effects.

$$[00970] \quad P_f(i_g = n_g) = \frac{\delta_g^{n_g > 0} P(i_g = n_g)}{\alpha_g}$$

$$\begin{aligned}
 [00971] \quad P_f(i_1 = n_1) &= \frac{1}{Z} \sum_{(n_2, \dots, n_k)} P(i_1 = n_1, i_2 = n_2, \dots, i_k = n_k) \delta_0 \delta_1^{n_1 > 0} \dots \delta_k^{i_k > 0} \\
 &= \frac{1}{\prod_g (e^{-f_g M} (1 - \delta_g) + \delta_g)} \delta_1^{n_1 > 0} P(i_1 = n_1) \prod_{g=1} (e^{-f_g M} (1 - \delta_g) + \delta_g)
 \end{aligned}$$

$$[00972] \quad = \frac{\delta_1^{n_1 > 0} P(i_1 = n_1)}{e^{-f_1 M} (1 - \delta_1) + \delta_1}$$

[00973]

[00974] The probability that guide g was not detected in a specific cell is now

$$[00975] \quad P(d_g = 0) = \frac{e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}}{e^{-f_g M} (1 - \delta_g) + \delta_g}$$

[00976]

$$[00977] \quad P(d_g = 0 | i_g) = (1 - p_d)^{i_g} = q^{i_g}$$

[00978]

[00979]
$$P_f(d_g = 0) = \sum_n q^n P_f(i_g = n) = \sum_n q^n \frac{\delta_g^{n>0} P(i_g = n)}{\alpha_g} =$$

[00980]
$$\frac{1}{\alpha_g} \sum_n q^n \delta_g^{n>0} P(i_g = n) = \frac{1}{\alpha_g} \left[P(i_g = 0) - \delta_g P(i_g = 0) + \delta_g \sum_n q^n P(i_g = n) \right] =$$

[00981]
$$\frac{1}{\alpha_g} [e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}] = \frac{e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}}{e^{-f_g M} (1 - \delta_g) + \delta_g}$$

[00982]

[00983] The probability of detecting guide g in a cell

[00984]
$$P_f(d_g = 1) = \frac{\delta_g (1 - e^{-f_g M p_d})}{\alpha_g}$$

[00985] The probability that a cell survived BFP selection given that none of the guides was detected

[00986]
$$P(B | \forall_g d_g = 0) = 1 - \frac{1}{\prod_g (1 - \delta_g + \delta_g e^{-f_g M (p_d - 1)})}$$

[00987]

[00988]
$$P(\neg B | \forall_g d_g = 0) = P(\forall_g i_g = 0 | \forall_g d_g = 0) = \frac{P(\forall_g d_g = 0 | \forall_g i_g = 0) P(\forall_g i_g = 0)}{P(\forall_g d_g = 0)}$$

[00989]
$$\prod_g \frac{\frac{P(i_g = 0)}{\alpha_g}}{e^{-f_g M} [e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}]} = \prod_g \frac{e^{-f_g M}}{e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}} =$$

[00990]
$$\prod_g \frac{e^{-f_g M}}{e^{-f_g M} (1 - \delta_g + \delta_g e^{-f_g M (p_d - 1)})} = \frac{1}{\prod_g (1 - \delta_g + \delta_g e^{-f_g M (p_d - 1)})}$$

[00991]

[00992] The probability of the BFP selection

[00993]
$$P(B) = 1 - \frac{e^{-M}}{\prod_g e^{-f_g M} (1 - \delta_g) + \delta_g}$$

[00994]

[00995]
$$P(\neg B) = P(\forall_g i_g = 0) = \prod_g \frac{P(i_g = n_1)}{e^{-f_g M} (1 - \delta_g) + \delta_g} = \frac{e^{-M}}{\prod_g e^{-f_g M} (1 - \delta_g) + \delta_g}$$

[00996]

[00997]

[00998] The new maximum likelihood function that accounts for the fitness parameters is

$$\log[L(M, p_d, \delta)]$$

$$= N_0 \log P(\mathcal{E} | \forall_g \alpha_g = 0) + \sum_g^y \{N_g \log[P(\alpha_g = 1)] + (N - N_g) \log[P(\alpha_g = 0)]\}$$

$$\begin{aligned} y &= \sum_g \left\{ N_g \log \frac{\delta_g (1 - e^{-f_g M p_d})}{\alpha_g} + (N - N_g) \log \frac{e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}}{\alpha_g} \right\} \\ &= \sum_g \{N_g \log[\delta_g (1 - e^{-f_g M p_d})] - N_g \log \alpha_g \\ &\quad + (N - N_g) \log(e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}) - (N - N_g) \log \alpha_g\} \\ &= \sum_g \{N_g \log(1 - e^{-f_g M p_d}) + N_g \log \delta_g \\ &\quad + (N - N_g) \log(e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}) - N \log \alpha_g\} \end{aligned}$$

$$\log[L(M, p_d, \delta)]$$

$$\begin{aligned} &= N_0 \log \left(1 - \frac{1}{\prod_g (1 - \delta_g + \delta_g e^{f_g M (1 - p_d)})} \right) \\ &\quad + \sum_g \{N_g \log(1 - e^{-f_g M p_d}) + N_g \log \delta_g \\ &\quad + (N - N_g) \log(e^{-f_g M} (1 - \delta_g) + \delta_g e^{-f_g M p_d}) - N \log \alpha_g\} \end{aligned}$$

[00999] **Analysis of perturbation effects on individual genes and gene modules.** The most variable genes from each Perturb-Seq experiment were filtered by using the *jackstraw* approach (Chung and Storey, 2015) to identify the most significant genes (q-value < 0.01) in the top 20 PCs of the coefficient matrix. The genes were then clustered using k-means clustering by their coefficients. Optimal *k* was chosen by visual inspection of clustering results. Gene ontology (GO) enrichment analysis was performed on each cluster using *goatools* (Tang et al., 2015) with FDR threshold of 0.05.

[001000] **Comparison to ChIP-seq binding profiles.** Applicants analyzed assignments of TF binding in gene promoters in BMDCs following LPS stimulation across four time points (0 minutes, 30 minutes, 60 minutes, and 120 minutes) (Garber et al., 2012). To test for significant binding, two tests were used. First, the regulatory coefficients of bound genes were compared to

those of unbound genes using a non-parametric Mann-Whitney test to identify significantly different means. Because of the possibility that this significance was driven by skewed covariates from unbound genes, Applicants also tested whether the coefficients of bound genes significantly deviated from 0, using the non-parametric one-sample Wilcoxon test. Finally, because TFs could both activate and repress genes, Applicants examined the number of bound genes significantly up- or down-regulated. To do this, Applicants used the distribution of covariates of unbound genes to define thresholds at the 5th percentile of lowest negative and highest positive coefficients. Any bound genes with coefficients that surpassed the thresholds were considered significant. Between the set of genes with significant positive and negative coefficients, Applicants used the larger set to infer whether the transcription factor was activating or inhibiting gene expression. Only BMDC expressed genes were considered in all ChIP-Seq analysis.

[001001] Power analysis and experimental design considerations. Power analysis was performed to determine how many cells are required to observe a signal as a function of observed effect size and baseline expression of a gene.

[001002] As an estimate of required read depth, Applicants downsampled at the UMI level. For example, for a vector of gene expression for a cell with the following values: [2,0,1,6] Applicants convert it into the following vector [1,1,3,4,4,4,4,4] on which downsampling is performed with equal probability without replacement. It is reconstructed into the original probability space by binning the observed integer counts. Applicants also downsampled cells without replacement from our observed set.

[001003] Applicants performed our regression analysis on the downsampled expression matrix for various amounts of downsampling and recomputed resulting regulatory matrix. For each level of downsampling, 10 instances are averaged. Applicants compared between the original regulatory matrix and the matrix that results after downsampling using either a Pearson or Spearman correlation.

[001004] Plasmid design and construction The perturb-seq expression vector (pBA439) was derived from a previously described CRISPRi expression vector (herein referred to as the “original CRISPRi expression vector”) (Addgene plasmid #60955) (Gilbert et al. 2014). To construct pBA439, the mU6-sgRNA-EF1a-PURO-BFP region from the original CRISPRi vector and a BGH polyadenylation sequence amplified by PCR from pcDNA3.1(+) (Invitrogen, V790-

20) were inserted in reverse orientation between the XbaI and EcoRI sites of the original CRISPRi expression vector. A random 18 nucleotide barcode was then inserted between the BFP and BGH polyA sequences (using disrupted EcoRI and AvrII sites) by Gibson assembly to construct the perturb-seq expression library (pBA571). The perturb-seq library was prepared with an estimated barcode diversity of >100,000 essentially as previously described (Kampmann, Bassik & Weissman 2014). Guide RNA protospacer sequences were individually cloned into both the original CRISPRi expression vector and the pBA571 library (between the BstXI and BlnI sites) by ligation. Each vector was then verified by Sanger sequencing of the protospacer and, if applicable, its corresponding barcode. Final guide expression vectors containing barcodes that introduced the conserved polyadenylation signal AATAAA were discarded. To construct pMH0001, a minimal ubiquitous chromatin opening element (UCOE) (Müller-Kuller et al. 2015) was inserted upstream of the SFFV promoter in the lentiviral dCas9-BFP-KRAB expression vector (Gilbert et al. 2014).

[001005] The UPRE reporter was built into a backbone for lentiviral expression that has been previously described (Addgene plasmid #44012) (Meerbrey et al. 2011). This parental vector was digested with AgeI and religated to remove unwanted functional cassettes, and the UPRE promoter region or EF1a promoter were inserted between the BamHI and XhoI site of the resulting product. The UPRE promoter region contains 5 UPR elements (UPREs, 5'-TGACGTGG-3') upstream of the *c-fos* minimal promoter (-53 to +45 of the human *c-fos* promoter) (Wang et al. 2000). Lastly, mCherry and sfGFP were cloned adjacent to UPRE and EF1a promoters, respectively (into an HpaI site). The resulting vectors are pBA407 (UPRE-mCh-Ubc-Neo) and pBA409 (EF1a-sfGFP-Ubc-Neo).

[001006] For testing of constant region variants in K562 cells, constant region variants fused to a GFP-targeting protospacer (EGFP-NT2, sequence GACCAGGATGGGCACCACCC) or a negative control protospacer were PCR-amplified and inserted into BstXI/XhoI-digested pBA439 (perturb-seq expression vector) by Gibson assembly. For testing of U6 promoters, U6 promoters from cow (bU6-2, GenBank DQ150531 and bU6-3, GenBank DQ150532 (Lambeth et al. 2006)), sheep (sU6-1, GenBank HM641427 and sU6-2, GenBank HM641426 (Hu et al. 2011)), buffalo (buU6, GenBank JN417659 (Zhang et al. 2014)), and pig (pU6, GenBank EU520423 (Chuang et al. 2009)) spanning ~400-500 bp upstream of the TSS, modified to contain a BstXI site at the TSS, and fused to EGFP-NT2 and the original constant region were

obtained as synthetic DNA segments (Integrated DNA technologies) and inserted into HpaI/XhoI-digested pBA439 (perturb-seq expression vector) by Gibson assembly.

[001007] Three-guide vectors were assembled by a two-step cloning procedure (Figure 88F). First, guide RNAs to be included were cloned into the corresponding single guide RNA expression vectors. Briefly, complementary oligonucleotides (Integrated DNA Technologies) containing the protospacer sequence and ligation overhangs were annealed and ligated into BstXI/BlnI-digested guide RNA expression vectors containing specific primer binding sites flanking the guide RNA expression cassette. The three-guide RNA expression cassettes were then PCR-amplified and assembled into HpaI/XhoI-digested pBA571 (perturb-seq expression library) by a single four-piece Gibson assembly step. Vectors were validated and barcodes were determined as described above. For three-guide vectors targeting the UPR branches, the bU6, mU6, and hU6 cassettes were designed to either express an sgRNA targeting *ATF6*, *EIF2AK3* (PERK), or *ERN1* (IRE1 α), respectively, or a non-targeting negative control sgRNA. The following protospacer sequences were used: *ATF6*-targeting, gGGGATCTGAGAATGTACCA; *EIF2AK3*-targeting, gCGGGCTGAGACGTGGCCAG; *ERN1*-targeting, gAGAACTGACTAGGCAGCGG; non-targeting sgRNA in bU6 cassette, gACGACTAGTTAGGCGTGTA; non-targeting sgRNA in mU6 cassette, gGCCAAACGTGCCCTGACGG; non-targeting sgRNA in hU6 cassette, gCCTTGGCTAAACCGCTCCC.

[001008] **Cell culture, DNA transfections, viral production, and construction of reporter cell lines** K562 cells were grown in RPMI-1640 with 25mM HEPES, 2.0 g/L NaHCO₃, 0.3 g/L L-Glutamine supplemented with 10% FBS, 2 mM glutamine, 100 units/mL penicillin and 100 μ g/mL streptomycin. HEK293T cells were grown in Dulbecco's modified eagle medium (DMEM) in 10% FBS, 100 units/mL penicillin and 100 μ g/mL streptomycin. Cells were treated with tunicamycin or thapsigargin (Sigma, T9033) solubilized in DMSO. Lentivirus was produced by transfecting HEK293T with standard packaging vectors using *TransIT*®-LTI Transfection Reagent (Mirus, MIR 2306). Viral supernatant was harvested at least 2 days after transfection and filtered through a PVDF syringe filter and / or frozen prior to infection.

[001009] To construct the UPRE reporter cell line, K562 cells stably expressing dCas9-KRAB (Gilbert et al. 2014), originally constructed from K562 cells obtained from ATCC 536 (RRID:CVCL_0004), were stably transduced with pBA407 and selected in media supplemented

with 500 µg/mL Geneticin (Gibco, 10131-035). The clonal line cBA010 was then selected by limiting dilution. cBA011 is a derivative of cBA010 containing pBA409. cBA011 was made by stable transduction and selection of GFP positive cells using fluorescence activated cell sorting on a BD FACSAria2. The GFP reporter cell line was constructed by infecting K562 cells stably expressing dCas9-KRAB with a Murine Stem Cell Virus (MSCV) retrovirus expressing GFP from the SV40 promoter. MSCV retrovirus was produced by transfecting amphotropic Phoenix packaging cell lines with standard packaging vectors. K562 cells stably expressing GFP were sorted to purity by flow cytometry using a BD FACS Aria2. To construct the GFP+ K562 UCOE-dCas9-KRAB cell line, the GFP reporter cell line was transduced with pMH0001 at a multiplicity of infection of ~3. Transduced cells were sorted for BFP expression (top 33%) by flow cytometry on a BD FACS Aria2. BFP fluorescence was monitored for several generations and found to be stable.

[001010] Design and cloning of constant region variants for testing in *E. coli*. Constant region bases to mutate were identified by inspection of the crystal structure of Cas9 bound to guide RNA and target DNA (Figure 88D, PDB ID codes 4OO8, 4UN3, 4ZT0 (Nishimasu et al. 2014, Anders et al. 2014, Jiang et al. 2015)). Bases that did not form direct contacts with Cas9 or with other nucleotides of the constant region were deemed amenable for mutation. If applicable, sequence conservation patterns of the base in crRNAs/tracrRNAs of *Streptococcus* species and previous reports of constant regions carrying changes at the base (Briner et al. 2014, Dang et al. 2015) were used to determine the type of mutation. In this fashion, 15 constant region variants with mutations in different parts of the constant region were designed (Figure 82D). The most diverse constant region variants cr2 and cr3 were designed by combining multiple individual mutations (Figure 82D).

[001011] To rapidly assess the activity of the variant constant regions, the variants were fused to a mRFP-targeting guide RNA (mRFP-NT1, sequence AACTTTCAGTTTAGCGGTCT) (Qi et al. 2013) and tested in an *E. coli* CRISPRi reporter strain for knockdown of mRFP (see below). To eliminate variability from copy number variation, guide RNAs were cloned into a plasmid for site-specific integration into the *E. coli* genome at *attL* and expressed from single copy from an IPTG-inducible P_{LlacO-1} promoter (Lutz, Bujard 1997). To construct the integrating guide RNA expression plasmid, a guide RNA expression cassette was first PCR-amplified from pgRNA-bacteria Addgene plasmid #44251) (Qi et al. 2013), modified to be flanked by the strong

synthetic terminators L3S3P22 and L3S2P21 (Chen et al. 2013) and inserted into pCAH63 (Haldimann, Wanner 2001) at the *ClaI/NheI* sites. The constitutive promoter from pgRNA-bacteria was replaced with the IPTG-inducible P_{LacO-1} promoter, generating pCs-550r. Then, pCs-550r was further modified to include the constant region used in mammalian CRISPRi (Gilbert et al. 2014), PCR-amplified with an mRFP-targeting protospacer and inserted into pCs-550r at the *SpeI* and *KpnI* sites to generate pMJ020. Finally, constant region variants 1-15 as well as cr2 and cr3 were cloned into pMJ020 by inverse PCR with mutations encoded in primer overhangs, by site-directed mutagenesis following standard procedures, or by insertion of a synthetic DNA segment encoding the constant region (Integrated DNA Technologies) into *SpeI/KpnI*-digested pMJ020 by Gibson assembly.

[001012] Construction of *E. coli* CRISPRi reporter strain and testing of constant region variants The *E. coli* CRISPRi reporter strain was constructed by sequential insertion of a construct for IPTG-inducible expression of dCas9, a construct for constitutive expression of mRFP, and a construct for IPTG-inducible guide RNA expression (described above) into the *E. coli* genome. First, a *lacIq-t0-P_{LacO-1}-dCas9* cassette (*lacIq* for strong expression of the Lac repressor; t0, a transcription terminator; P_{LacO-1} -*dCas9*; for IPTG-inducible expression of *S. pyogenes* D10A/H840A *Cas9* (dCas9)) was inserted into the chromosome of *E. coli* BW25113 at +19 *attL* via of lambda Red recombinase-mediated recombineering (Thomason et al. 2014). Then, a *nfsA::mRFP-kan* cassette for expression of mRFP from the J23119 promoter, a strong synthetic constitutive promoter from the Anderson promoter collection (<http://parts.igem.org/Promoters/Catalog/Anderson>), was inserted into an *E. coli* MG1655-derived strain by lambda Red recombinase-mediated recombineering as described previously (Qi et al. 2013), and moved from the MG1655-derived strain into the dCas9-expressing BW25113 strain by P1 transduction and selection on kanamycin following a published protocol (Thomason, Costantino & Court 2007). Plasmids for expression of mRFP-NT1 with the different constant region variants were integrated into the dCas9- and mRFP-expressing strain at *attL* using the helper plasmid pINT-ts (Haldimann, Wanner 2001), selecting for chloramphenicol resistance.

[001013] Single colonies of strains with the integrated guide RNA expression plasmids were inoculated into LB and grown overnight in deep 96-well blocks at 37 °C with shaking at 900 rpm. Stationary-phase cultures were back-diluted 1:30 and grown into mid-exponential phase, at which point they were back-diluted 1:10000 into LB with 1 mM IPTG for induction of guide

RNA and dCas9 expression. Induced cultures were grown at 37 °C with shaking until OD_{600 nm} reached ~0.4-0.7 (approximately 5 hrs), at which point they were diluted 1:30 in PBS in a 96-well plate. RFP fluorescence was recorded on a LSR-II flow cytometer (BD Biosciences) equipped with a 96-well high-throughput sampler. Each experiment was carried out using three individual colonies for each constant region variant. RFP levels were normalized to those of a strain expressing a non-targeting guide RNA.

[001014] Testing of single- and three-guide vectors in K562 cells by GFP knockdown Vectors for expression of EGFP-NT2 in different contexts were delivered into GFP+ K562 cells with dCas9-KRAB or with UCOE-dCas9-KRAB by lentiviral transduction at MOI of 0.1-0.5. For all experiments using GFP+ K562 cells with UCOE-dCas9-KRAB, transduced cells were allowed to recover for 2 d, then selected to purity using 2 µg/mL puromycin for 3 d, and allowed to recover for another 2 d before GFP levels were recorded by flow cytometry on a LSR-II flow cytometer (BD Biosciences). For experiments involving only GFP+ K562 cells with dCas9-KRAB, cells were grown out for 7-9 d after transduction and GFP levels were recorded by flow cytometry, using BFP expression to gate for transduced cells. Flow cytometry data were analyzed using FlowCytometryTools v0.4.5 (<http://eyurtsev.github.io/FlowCytometryTools/>). For plotting, flow cytometry events were normalized to population size and the histograms were smoothed by kernel density estimation. For estimating knockdowns, GFP levels of wild-type (GFP-) K562 cells were subtracted.

[001015] Perturb-seq screening Viruses were individually packaged and harvested in preparation for perturb-seq screening. Individual packaging of the lentivirus and pooling at the step of virus or cells was done to avoid intermolecular recombination of proviral genomes and to ensure maintenance of paired barcode-sgRNA coupling (Sack et al. 2016). For the pilot experiment represented in Figure 81C, 1D, 1F, cBA010 cells were individually spininfected with virus (at 33°C for 2 hours at 1000xg) in media supplemented with 8 µg/mL polybrene; 5 hours post spinfection, virus was removed by centrifugation and cells were resuspended in fresh media. Three days later, a transduction efficiency of 20-30%, as determined by percentage of BFP positive (BFP+) cells, was measured by flow cytometry and cells were pooled with equal numbers of guide RNA-containing (BFP+) cells. Control cells were included in the pool at 3-fold coverage. Pooled cells were then grown in the presence of puromycin (3 µg/mL) for 5 additional days. Seven days post transduction cells were sorted on a BD FACSAria2 to near purity and

eight days post transduction the sorted cells were separated into droplet emulsion using the Chromium™ Single Cell 3' Solution according to manufacturer's instructions (10X Genomics).

[001016] For the perturb-seq epistasis experiment, seven three-guide vectors targeting every possible combination of ATF6, IRE1, and PERK as well as two independent three-guide vectors with three negative control guide RNAs and different barcodes were individually packaged into lentiviruses. Freshly produced (i.e. not frozen) lentiviruses were spinfected into cBA007 cells (at 33°C for 2 h at 1000 × g) in media supplemented with 8 µg/mL polybrene. The virus was removed by centrifugation and cells were resuspended in fresh media. Three days after infection, transduction efficiencies of 5-10% were measured by flow cytometry. Cells were combined into a pool with equal numbers of transduced (BFP+) cells for each vector (resulting in 2-fold excess of negative control vectors) and the combined cells were then sorted on a BD FACS Aria2 to near purity. To limit heterogenous effects of cell microenvironments caused by cell settling, the sorted cells were grown with continuous agitation on an orbital shaker. Five days after infection, the pooled and sorted cells were split into three populations, which were treated as follows: 1) DMSO control treatment for 6 hr; 2) treatment with 4 µg/mL tunicamycin for 6 hr; and 3) treatment with 100 nM thapsigargin for 4 hr. At the end of the treatment, the cells were separated into droplet emulsion using the Chromium™ Single Cell 3' Solution according to manufacturer's instructions (10X Genomics). Cells loaded onto the device were 90.4%, 87.9%, and 85.3% viable for the different treatment conditions, respectively.

[001017] For the large-scale perturb-seq screen of UPR-inducing guide RNAs, viruses were individually titered by test infections into cBA011 cells and then pooled evenly. To account for varied effects on cell viability across the guide RNA sub-library and minimize cell number difference, pooling titers were determined by the percentage of BFP+ cells remaining 6 days post transduction. Two negative control guides were included, NegCtrl-2 and NegCtrl-3. NegCtrl-2 and select guides (those encoded by pDS002, pDS017, pDS026, pDS032, pDS033, pDS044, pDS052, pDS088, pDS091, pDS160, pDS186) were included at higher representation within the lentivirus pool, 8-fold and 2-fold, respectively. The lentivirus library pool was then used to infect cBA010 cells (performed by spinfection at 33°C for 3 hours at 1000xg) so that a single pooled cell population with all perturbations would be carried through subsequent steps. Post centrifugation, cells were immediately removed from virus and transferred to a spinner flask for growth in fresh media. Three days later, a transduction efficiency of 15% was measured by flow

cytometry and BFP⁺ cells were sorted to near purity on a BD FACSAria2. To limit heterogenous effects of cell microenvironments caused by cell settling, the sorted cells were grown with continuous agitation on an orbital shaker. Approximately seven days post transduction cells were separated into droplet emulsion using the Chromium™ Single Cell 3' Solution across two separate runs totaling 10 lanes on the device according to manufacturer's instructions (10X Genomics). Cells loaded onto the device were 92% BFP⁺ and 93-94% viable, as determined by flow cytometry.

[001018] For all perturb-seq experiments single-cell RNA-seq libraries were prepared according to the Single Cell 3' Reagent Kits User Guide (10X Genomics). However, this protocol produces libraries that are not compatible with the HiSeq 4000 due to the presence of some sort of toxic byproducts that it is uniquely sensitive to. To remove this issue, Applicants implemented a short cleanup protocol taking place after library preparation. 120 – 200 ng of library material was split into parallel PCR reactions containing 0.3 μM each of the Illumina P5 and P7 primers, and amplified using Kapa HiFi ReadyMix according to the following protocol: (1) 95 °C 80 sec (2) 98 °C 20 sec/65 °C 30 sec/72 °C 20 sec for 6 cycles (3) 72 °C 1 min. PCR products were then SPRI-purified at 1X ratio, repooled during elution, and then fragments of length 350 – 525 bp were selected using the BluePippin (Sage Science).

[001019] Genome-scale CRISPRi screening Reporter screens were conducted using protocols similar to those previously described (Gilbert et al. 2014, Sidrauski et al. 2015). The hCRISPRi-v1 (Gilbert et al. 2014) or the compact (5 sgRNA/gene) hCRISPRi-v2 (Horlbeck et al. 2016) sgRNA libraries were transduced into cBA011 cells at an MOI < 1 (percent BFP⁺ cells was ~45% and 26%, respectively). For the hCRISPRi-v1 screen, cells were grown in spinner flasks for 2 days without selection, followed by 3 days of selection with 1 μg/mL puromycin. Screen replicates were split post infection and carried separately throughout the remainder of the experiment. One replicate arm of the hCRISPRi-v1 screen was carried with media supplemented with 88-150nM ISRIB throughout, although differences observed between the replicates at the level of both sgRNAs and genes (Table S3) were negligible. For the hCRISPRi-v2 screen, cells were grown in spinner flasks for 2 days without selection, followed by 5 days of selection with 1-3 μg/mL puromycin. Screen replicates were split into separate spinner flasks on day 3. For both screens, cells were separated into those with the highest (~28-33%) and lowest (~30-35%) mCherry/GFP ratio eight days post transduction by fluorescence-activated cell sorting. Cell

pellets were frozen after collection. Approximately 23-30 million cells were collected per bin during screening of the hCRISPRi-v1 library (a representation of ~450) and 19-22 million cells per bin for hCRISPRi-v2 (a representation of ~600). Genomic DNA was isolated from frozen cells essentially as described previously (hCRISPRi-v2, 18239126) and the sgRNA-encoded regions were enriched, amplified, and prepared for sequencing (Gilbert et al. 2014).

[001020] Sequenced protospacer sequences were aligned and data were processed as described (Gilbert et al. 2014, Horlbeck et al. 2016) with custom Python scripts (available at <https://github.com/mhorlbeck/ScreenProcessing>). Reporter phenotypes (referred to as Reporter signal) for library sgRNAs were calculated as the log₂ enrichment of sgRNA sequences identified within the high mCherry/GFP cells over the low mCherry/GFP cells. Phenotypes for each transcription start site were then calculated as the average reporter phenotype of the 3 sgRNAs with the strongest phenotype by absolute value (most active sgRNAs). Mann-Whitney test p-values were calculated by comparing all sgRNAs targeting a given TSS to the full set of negative control sgRNAs. For data presented in Figures 84D-F and 90B,C, genes with multiple targeted TSSs were collapsed such that only the TSS with the lowest p-value was used. Screen hits were defined as those genes (or separately those TSSs) with a discriminate score, defined as the absolute value of a calculated reporter phenotype over the standard deviation of all evaluated phenotypes multiplied by the log₁₀ of the Mann-Whitney p-value for given candidate, greater than 7. Growth screen data in Figure 90C has been reported elsewhere, and the unreported screen was conducted in parallel essential as described (Horlbeck et al. 2016). Gene ontology analysis was conducting using select databases (GOTERM_BP_FAT, GOTERM_CC_FAT, GOTERM_MF_FAT, KEGG_PATHWAY) and hits (calculated from all TSSs) with a phenotype of greater than 1 using DAVID Bioinformatic Resources 6.8 Beta (<https://david-d.ncifcrf.gov/>) with (Huang, Sherman & Lempicki 2009a, Huang, Sherman & Lempicki 2009b). Biological classifications reported in Figure 84F-G were manually assembled from the literature and using resources from the HUGO Gene Nomenclature Committee (www.genenames.org), AmiGO, the GO Consortium's annotation and ontology toolkit (Carbon et al. 2009) (<http://amigo.geneontology.org>), DAVID Bioinformatic Resources (<https://david.ncifcrf.gov/>) (Huang, Sherman & Lempicki 2009a, Huang, Sherman & Lempicki 2009b).

[001021] **Individual re-test of sgRNA reporter phenotypes** Viruses were individually packaged and harvested as described above. UPRE reporter-containing K562 cells (cBA011)

cells were infected with thawed virus. Additionally, parental K562 dCas9-KRAB cells (Gilbert et al. 2014) were transduced with negative controls. Flow cytometer readings of the mCherry UPRE signal and GFP EF1a signal were taken periodically and 8 days post transduction. Median fluorescence signals were analyzed by subtracting an average background signal from control-transduced K562 dCas9-KRAB cells and normalizing the mCherry, GFP, or mCherry/GFP measurement from guide-containing cells (as determined by BFP fluorescence) in each well to untransduced cells. Data from wells with fewer than 500 transduced or untransduced cells or with lower than expected BFP signal (3 standard deviations below the mean of BFP median from all other experimental wells) were systematically discarded from further analysis. For experiments where a flow cytometer reading was taken on the second day post transduction, data was also filtered for a minimum day 2 viability percentage.

[001022] RT-qPCR and semi-quantitative PCR for XBP1 mRNA splicing Cells were harvested and total RNA was isolated using TRIzol® Reagent (ThermoFisher Scientific, 15596-018) and Phase Lock Gel tubes (VWR, 10052-170) or NucleoSpin® RNA (Macherey-Nagel, 740955.50) essentially according to manufacturers' instructions. RNA prepared by TRIzol® extraction was treated with TURBO™ DNase (ThermoFisher Scientific). RNA was converted to cDNA using SuperScript® II or SuperScript® III Reverse Transcriptase (ThermoFisher Scientific) under standard conditions with oligo(dT) primers or random hexamers with or without RNaseOUT™ Recombinant Ribonuclease Inhibitor (ThermoFisher Scientific). Quantitative PCR reactions were prepared with 1X master mix containing 1X Colorless GoTaq® Reaction Buffer (Promega, M792A), MgCl₂ (0.7 mM), dNTPs (0.2 mM each), primers (0.75 μM each), and 1000X SYBR Green with GoTaq® DNA polymerase (Promega, M830B) in 22 μL reactions. Reactions were run on a LightCycler® 480 Instrument (Roche). Semi-quantitative XBP1-specific PCR reactions were prepared with 2 μL of cDNA diluted 1:10 using a master mix containing 0.9X Colorless GoTaq® Reaction Buffer (Promega, M792A), dNTPs (0.23 mM each), primers (0.45 μM each) with GoTaq® DNA polymerase (Promega, M830B) in 22.1 μL reactions. These reactions were run on a standard thermocycler program with 30 second at 60.5°C for annealing and 28 cycles. PCR products were visualized on 8% TBE gels.

[001023] Quantification and Statistical Analysis. Applicants will first provide an overview of the methods used, and then describe their specific application to each figure.

[001024] Pipeline overview All analysis was performed in Python, using a combination of Numpy, Pandas, scikit-learn, and a custom-made perturb-seq library. The general outline is presented in Figure 89A, and Applicants will outline the steps below.

[001025] Sequencing: Reads from 10X single-cell RNAseq experiments were aligned and collapsed to unique molecular identifier (UMI) counts using 10X's cellranger software (version 1.1). The result is a large digital expression matrix with cell barcodes as rows and gene identities as columns.

[001026] Perturbation identity mapping Specifically amplified guide barcode libraries were described as above and either sequenced as spike-ins or independently. The specific amplification strategy Applicants used preserved the 3' end of the transcript (and thus the cell barcode and UMI of a given captured molecule) and introduced an Illumina read 1 primer upstream of the guide barcode sequence. These reads were aligned using bowtie (flags: -v2 -q -m1) to a library of expected barcode sequences. Applicants then collapsed all reads with common cell barcode, UMI, and read identity (as some reads were not mapped by bowtie due to low quality scores) to produce a table consisting of possible guide identities for each cell, and the number of molecules attributing a given guide identity to that cell. Applicants defined the coverage of a given proposed identity as the number of reads divided by the number of UMIs, and defined a proposed identity as having good coverage if it: (1) had a coverage level at or above the mean coverage level minus two times the standard deviation in coverage (2) had at least 50 raw reads and (3) had at least 3 UMIs. Any cell that had only a single identity that met these criteria was assigned that perturbation identity. Any cell that had two or more identities meeting these criteria was assigned as a multiple (either a multiple infection, or a multiple encapsulation during emulsion generation). Any cell that had no identities meeting these criteria was assigned as unidentifiable.

[001027] Expression normalization To normalize for differences in sequencing capture and coverage across emulsion droplets, Applicants rescaled all cells to have the median number of total UMIs (i.e. each row of the raw digital expression matrix is normalized to the same sum). Expression of each gene was then z-normalized with respect to the mean and standard deviation of that gene in the control (unperturbed) population:

$$x_{\text{normalized}} = \frac{x - \mu_{\text{control}}}{\sigma_{\text{control}}}$$

[001028] This normalization means that control cells always have mean normalized expression of 0 for all genes and standard deviation 1, so that the units of expression are “standard deviations above/below the control distribution.”

[001029] In the epistasis experiment, the control population was the DMSO-treated cells. In the perturb-seq experiment, they were the cells containing the NegCtrl-2 guide. In the perturb-seq experiment, the mixed population was run in ten separate pools that were treated independently during library preparation (corresponding to lanes on the 10X Chromium instrument and on the Illumina sequencer). To avoid any lane-dependent batch effects, cells were normalized to the control cell distribution within the same lane.

[001030] **Low cell count/inviable cell removal** While developing the low rank ICA method described below, Applicants observed that all experiments always contained two subpopulations that were peculiar in that they contained roughly equal membership from all perturbations. Further investigation showed that these were a group of cells with systematically lower total UMI counts (visible as a small second mode in the distribution of total UMIs per cell) and a group of cells that contained markers of activation of apoptotic programs. Applicants attributed the first population to partly failed RNAseq library preparation occurring in a small number of emulsion droplets, and the second to inviable cells (which Applicants knew were present in the cells placed used in the 10X experiment). Neither population composed more than a few percent of the total number of cells. Though low rank ICA always isolated these in an unbiased way, Applicants generally excluded them from analysis. The low UMI count cells were simply removed using a threshold. To remove the apoptotic cells, Applicants trained a random forest regressor (described in more detail below in the section on UPR branch activation scoring below) to recognize them using the cells in our epistasis experiment as training data.

[001031] **Identification of differentially expressed genes** The end result of the previous steps is a normalized gene expression matrix where each cell has been assigned a perturbation identity. In general, Applicants were interested in analyzing differences between populations, and used two distinct strategies for isolating interesting genes.

[001032] *Kolmogorov-Smirnov test/metric*: The Kolmogorov-Smirnov test is a nonparametric test for equality of probability distributions based on a metric defined on their cumulative

distribution functions. Specifically, if $F_{\text{perturbed}}$ and F_{control} are the CDFs for a given gene in the perturbed and control distribution, the test statistic is

$$D = \sup_x |F_{\text{perturbed}}(x) - F_{\text{control}}(x)|$$

[001033] This can be assigned a p -value in a standard way. However, the large scale of single-cell data means that many genes were often significantly perturbed without being interestingly perturbed, simply because of small differences detected by great sampling depth. Thus in some cases Applicants placed a direct threshold on the test statistic D itself, which ensured that changes were both significant (in the statistical sense) and also of reasonable magnitude, as it is valid metric on the space of CDFs.

[001034] *Random forest classifier* An advantage of perturb-seq is that the cell populations are known, which means that supervised learning methods can be brought to bear. Our strategy here was motivated by the idea that a gene is likely important for a given perturbation if its expression level can be used to accurately predict that perturbation's identity. This idea is particularly useful when many perturbations are being compared, as what you want then are the genes that best distinguish all of the perturbations from each other.

[001035] To leverage this idea, Applicants used random forest classifiers. Given a set of perturbations, Applicants would train a random forest classifier to predict perturbation identity using a subset of genes. Specifically, Applicants used the implementation of extremely randomized trees implemented in scikit-learn, generally with 1000 trees in the forest. Applicants performed a two-stage fitting process for a given number of desired features N_{genes} . First, Applicants set aside 20% of the cells. The remaining 80% were used to train a random forest classifier (usually with 1000 estimators) to predict the perturbation identity using the normalized expression profile for each cell (with some threshold on gene expression level) as the set of features. The random forest assigns importances to features during training based on their predictive value, and Applicants would then take the top N_{genes} sorted by importance as the set of most informative genes. To evaluate how informative these genes were, Applicants would then retrain the classifier using only these genes, and predict the perturbation present in the 20% of cells Applicants had initially set aside. For sets of perturbations with large differences, Applicants routinely saw accuracies of 80-90%. The genes chosen by the random forest essentially always showed marked differences by the Kolmogorov-Smirnov approach outlined

above, and the forests had the advantages that they scaled to an arbitrary number of perturbations, and that the selected genes were known to vary informatively across perturbations instead of simply having a difference in distribution.

[001036] **Low rank ICA** Single-cell data are intrinsically very noisy, either due to real biological variation or problems in capture efficiency. As described in the main text, these effects can affect the sensitivity of methods like principal components analysis, which is intrinsically variance-maximizing and hence very sensitive to outliers. To isolate larger trends within the data, Applicants developed a simple two-step approach called low rank ICA. The first step consists of isolating a low rank approximation of the dynamics within the experiment. To do this, Applicants used Robust PCA (Candès et al. 2011), which seeks a decomposition of the form $X = L + S$

[001037] where X is the normalized expression matrix, L is a low rank matrix, and S is a sparse matrix (most entries are zero). Specifically, Robust PCA solves the optimization problem

$$\min_{L,S} \|L\|_* + \lambda \|S\|_1 \quad \text{subject to} \quad X = L + S$$

[001038] where $\|\cdot\|_*$ is the nuclear norm (sum of singular values) and $\|\cdot\|_1$ is the sum of the absolute values of the entries of the matrix. These constraints naturally induce L to be low rank, and S to be sparse. In implementations, Applicants used the augmented Lagrangian multiplier method (Lin, Chen & Ma 2010), which was fast and efficient.

[001039] Applicants should note that our interpretation of this optimization problem is slightly different from that seen in some other instances, where S is regarded as capturing noise corrupting the “true” dynamics seen in L . In single-cell data the “noise” may actually be biological in origin, but our primary intent is to isolate the low rank approximation L , which is effectively a smoothed version of the population’s dynamics that leaves major trends intact. The advantage of the decomposition of course is that the S matrix is still available afterward, and it may in fact carry useful information about highly stochastic processes within the population.

[001040] Our next goal was to isolate the major trends within the low rank dynamics of the population. To do this Applicants applied independent components analysis (ICA). ICA posits a model in which the expression of a given gene (y_j) can be decomposed as a linear sum of various effects (s_1 to s_n) that are statistically independent of each other:

$$y_j = a_{j1}s_1 + a_{j2}s_2 + \dots + a_{jn}s_n$$

[001041] Solving this problem is beyond the scope of this section, but our interest lies primarily in the vector version of this formula,

$$\mathbf{y} = A\mathbf{s}$$

[001042] in which a cell's expression profile \mathbf{y} (over all genes) is viewed as a linear sum of independent effects, and the equivalent matrix version

$$Y = AS$$

[001043] in which Applicants decompose all of the dynamics of the cells within our population (the columns of Y) into sums of independent components (ICs). The matrix A above is called the mixing matrix, and in our context describes which genes contribute to which effects. As noted above, the key difference in this case from decompositions like principal components analysis is that the \mathbf{s} components are derived in a way to make them as statistically independent as possible. Once the matrix A is estimated, Applicants can then "unmix" the dynamics of each cell in the population by applying the inverse operation (denoted here by W) to its expression profile:

$$\mathbf{s} = W\mathbf{y}$$

[001044] This yields a low-dimensional description of what each cell is doing in terms of the independent factors given by \mathbf{s} .

[001045] In our case Applicants apply ICA to the low rank matrix L , i.e. $Y = L^T$ above. Thus Applicants try to separate the population's low rank dynamics into independent factors. As the ICA minimization problem posed in the strongest form cannot practically be solved, different algorithms will give somewhat different answers based on the tradeoffs they make. After trying several methods, Applicants settled on the ProDenICA algorithm (Friedman, Hastie & Tibshirani 2001), which Applicants found to frequently give the highest quality components.

[001046] In general Applicants applied low rank ICA in two ways. First, it can be used to partition cells into subpopulations. Strong trends often lead to independent components that are bimodal, so simply thresholding the value of a component is a means of clustering. Applicants note however that an advantage of this method of subpopulation identification is that it can also identify continuous trends, rather than enforcing discrete categories that may not exist like in other methods of clustering. Secondly, the mixing matrix A is very informative, as it determines the extent to which each gene contributes to a given component. This can be useful both in understanding what the component is measuring (if the most heavily weighted genes have a clear common function) and in identifying groups of genes that are co-expressed in an unbiased way.

[001047] Interpretation of independent components does have some caveats. First, they have no natural sign (so an “enriched” effect may appear as a low value of an independent component) or scale: thus there is no natural order where the first IC is somehow more informative than the next, consistent with the fact that they are meant to represent independent effects. Applicants do note that one pragmatic solution is to order the components by the norm of the corresponding column in the mixing matrix, which tends to place the most interesting components first.

[001048] **t-sne visualization** To obtain two-dimensional projections of the population’s dynamics, Applicants first further reduce the dimensionality of the low rank matrix L using classical PCA (with the number of components determined from a scree plot), and then further reduce these components via t -distributed stochastic neighbor embedding. Applicants occasionally directly visualize the ICs in this way as well, but because they lack intrinsic scale like principal components, dominant effects can be crowded out by minor ones.

[001049] **Hierarchical clustering of genes** Several of the analyses in the paper use single-cell co-expression information to cluster genes. For a given list of genes, Applicants perform this clustering by first calculating the gene-gene correlation matrix ρ over all cells in the population. This is then converted to a dissimilarity matrix π via the transformation $\pi = \sqrt{2(1 - \rho)}$. The dissimilarity matrix is then clustered using Ward’s method. For visualization purposes, Applicants then apply the optimal leaf ordering algorithm in MATLAB. This reorders the leaves in the dendrogram by flipping tree branches to maximize the similarity between adjacent leaves, but without dividing any branches (i.e. the clustering is unchanged, but the dendrogram ordering is in some sense optimal). Applicants then reorder the columns and rows of the correlation matrix via the resulting ordering, so that groups of genes with correlated expression appear as blocks along the diagonal.

[001050] **Cell cycle position** Applicants used an approach previously described, in which the expression of sets of experimentally-derived genes specific for each cell cycle phase is used for each cell to score cell cycle phase (Macosko et al. 2015).

[001051] **Average expression profiles** Applicants often create synthetic bulk profiles for different populations. These are created by averaging the normalized expression profile of each cell within that population together.

[001052] **Analytical steps for each figure.** Applicants now describe the analysis behind each figure in the paper, with references as necessary to the above sections.

[001053] **Single-cell analysis in Figure 83.** Applicants formed a population consisting of cells treated with 100 nM thapsigargin in each of our 8 genetic backgrounds, along with DMSO-treated control cells. As outlined in the “Low cell count/inviable cell removal” section, Applicants removed cells with substantially lower than average UMI counts or that scored strongly for inviability markers from analysis, as these groups partitioned away from the rest of the population in preliminary analyses. For each perturbation, Applicants then looked for genes that were differentially expressed relative to the control, as described in the “Identification of differentially expressed genes.” Applicants made a list of all genes that had a mean expression of at least 0.5 UMI per cell in the population and for which the Kolmogorov-Smirnov test statistic $D > 0.15$ in at least one perturbation. This led to a group of 1,711 differentially expressed genes. Applicants formed a reduced gene expression matrix containing only these genes, and performed low rank ICA to reduce the population’s dynamics therein to 16 ICs (Figure 89B). Applicants examined the raw trends in the population by reducing the low rank matrix to 16 components via PCA (16 components) and then to two dimensions via t-sne, revealing a general breakdown by perturbation and by cell cycle within each perturbation (Figure 89B). Applicants then looked for ICs whose average value varied either across the perturbation, or across the cell cycle position. For each category, four components showed clear trends at the average level and in the t-sne plots (Figure 89B). For example, several of the components clearly showed the expected epistasis patterns for PERK, ATF6, and IRE1 (Figure 89B). The plots made in Figure 83B of the main text were then made by furthering reducing only the ICs that varied across perturbation (IC1 – IC4 in Figure 89B) or across the cell cycle (IC5 – IC8 in Figure 89B) to two dimensions using t-sne. (I.e., Applicants constructed matrices with cells as rows and the given ICs as columns and reduced those matrices to two dimensions with t-sne.)

To make the plots in Figure 83C, Applicants then subsampled our population to only look at cells treated with thapsigargin with or without depletion of PERK, and the DMSO-treated control. Applicants applied the same methodology as above, though with 12 ICs instead of 16. The “G1 cell” IC described in the main text was bimodal within each subpopulation (see inset in right panel of Figure 83C), but with varying distances between the two modes (note that the IC takes a substantially lower value in the +Tg population than in any of the others, Figure 83C). Applicants split each population at the location of the mode. The cell cycle position histograms were made as described above. To make Figure 83E, Applicants took the 25 genes that most

positively influenced the IC and the 25 genes that most negatively influenced the IC (by sorting the mixing matrix column for that IC by coefficient value) and then clustered them based on coexpression as described in the “Hierarchical clustering of genes” section. The meaning of each cluster was discerned by the pattern of up- and down-regulation observed within.

[001054] Branch epistasis analysis in Figure 83F. Applicants created two populations: (1) consisting of cells treated with 100 nM thapsigargin in each of our 8 genetic backgrounds, along with DMSO-treated control cells, or (2) consisting of cells treated with 4 µg/mL tunicamycin in each of our 8 genetic backgrounds, along with DMSO-treated control cells. To identify informative differentially regulated genes, Applicants used the random forest classifier method described in the “Identification of differentially expressed genes,” limiting the random forest to pick 100 genes for each of the two populations. Applicants then combined these two lists in one and discarded any duplicate genes. Applicants created average profiles of expression of these genes for each of the nine conditions present in the two populations, as visualized in Figure 83F. The average epistatic phenotype of a gene can then be viewed as a 9-vector in either the thapsigargin- or tunicamycin-treated populations. Applicants discarded any genes where the correlation between these two conditions was less than 0.9, as Applicants were only interested in factors that showed the same regulation in response to both conditions. The end result was the 104 genes presented in Figure 83F. These were then clustered based on their co-expression pattern as described in the “Hierarchical clustering of genes” section, with the exception that Spearman correlation was used instead of Pearson correlation (to emphasize the large shifts in expression across the population). Rough meanings were ascribed to clusters based on the average pattern of gene expression across perturbations, but Applicants emphasize that many targets show some degree of cross-regulation. To assess this in an unbiased way, Applicants constructed a matrix consisting of the average expression of the 104 assayed genes across the 17 unique conditions present in the experiment, and reduced it to four independent components using FastICA. Three of the components clearly corresponded to ATF6, IRE1, and PERK perturbations, as they showed banded patterns in the reduced matrix matching the pattern of epistasis for those regulators seen in Figure 83F (e.g. the PERK component was high in all PERK conditions, and low everywhere else). The fourth component was low in the DMSO and all tunicamycin-treated conditions, and high in the thapsigargin-treated condition, so Applicants discarded it as representing the difference between the two perturbations. The panel at the

bottom of Figure 83F plots the mixing matrix coefficients for each gene in the indicated component, and thus determines how much that gene affects that component's value.

[001055] Genome-wide CRISPRi screen in Figure 84. Analysis of the screen is described above along-side the experimental details above.

[001056] Clustering of guides and perturbations in Figure 85. Applicants first split our large perturb-seq population into subpopulations based on guide identity and created average expression profiles (see "Average expression profiles" section) for each perturbations of all genes with mean representation >1 UMI per cell. Applicants calculated the perturbation-perturbation correlation matrix between all average expression profiles, and then clustering it using the same methodology described in the "Hierarchical clustering of genes." The ordering is seen in Figure 91A. Because guides targeting the same gene behaved similarly in this analysis, in subsequent analyses Applicants instead split the population into subpopulations based on guide target (thus merging subpopulations that had different guides that targeted the same gene). Applicants clustered these profiles using the same criteria, and optimally ordered the resulting dendrogram and correlation matrix (as described in "Hierarchical clustering of genes") to produce Figure 85A.

[001057] Assessing guide homogeneity and knockdown in Figure 85. Most guide targets were too low abundance to interrogate directly at single-cell resolution. Applicants first directly visualized the shift in guide target expression induced by the guide, comparing the distribution of expression in control cells to cells perturbed for a given target (Figure 91B). Applicants calculated mean knockdown per guide (Figure 85C), and assigned 95% confidence intervals to our estimates via bootstrapping.

[001058] Applicants also attempted to assess to what extent knockdown varied throughout the population based on phenotype. To do this, Applicants needed an unbiased means of assessing deviation in behavior from the control cells. Applicants leveraged a method called OneClassSVM (Schölkopf et al. 2001) which is a means of novelty detection. Given a set of training exemplars, a OneClassSVM learns an estimate of how those points are distributed (potentially in a high-dimensional space). When given new observations, the OneClassSVM then estimates how likely it is that those observations came from the same distribution as the training set, or if they are outliers (potentially novel). In our case Applicants trained the OneClassSVM using control cells, and thus scored the extent to which perturbed cells scored as outliers, or if

they fell within the expected range of behavior for unperturbed cells. Specifically, for each guide target, Applicants performed the following algorithm:

1. Form a population of all cells perturbed for that target, and an equal number of randomly sampled control cells.
2. Find all genes that are expressed at an average level of 0.5 UMI per cell or higher and that are differentially expressed between control and perturbed cells by the Kolmogorov-Smirnov test (as described in “Identification of differentially expressed genes”) at $P < 0.01$.
3. Form a reduced gene expression matrix consisting only of the differentially expressed genes. Create a low-dimensional picture of the dynamics within this population by reducing this matrix to 8 dimensions via PCA.
4. To form an estimate of “normal” behavior, train a OneClassSVM model to estimate the support of the control cells in this 8-dimensional space. The model was trained assuming a contamination rate with outliers of 5%.
5. Score each cell in the perturbed population using the OneClassSVM model to estimate the extent it deviates from control behavior.

[001059] These scores generally assigned most or all of the perturbed cells outlier status, except in guides where very few genes were perturbed to begin with (bottom panel of Figure 85D). Ordering the cells by score, Applicants split each perturbed cell population into top third and bottom thirds (i.e. the most and least perturbed cells) and assessed the difference in average knockdown in each of these populations (Figure 91C), with a difference of ~8% on average.

[001060] Applicants also reported the number of differentially expressed genes measured above in the bottom panel of Figure 85D.

[001061] **Scoring branch activation in Figure 85D.** As outlined in the main text, Applicants adopted a data-driven strategy to score activation of each of the UPR branches using the epistasis experiment as training data. To do this, Applicants assigned the label “ATF6 active”, “IRE1 active”, or “PERK active” to each cell in the epistasis experiment based on whether a given branch was present (i.e. not depleted) and induced (tunicamycin or thapsigargin had been added). For example, cells treated with thapsigargin and depleted for IRE1 would have ATF6 and PERK active, but not IRE1. Applicants converted these labels to scores of 0, for inactive, and 1, for active, and then trained three random forest regressors to predict activation of each branch. The

training strategy was the same as outlined in the “Identification of differentially expressed genes” section: each cell was regarded as a training data point, with every gene of mean > 1 UMI initially regarded as a possible feature for predicting branch activation. Each regressor was constrained to use the top 25 genes for predicting branch activation, as Applicants found no performance improvement when more genes were included. The genes isolated as most important by the three regressors for scoring activation of the three branches all appear in the epistasis analysis in Figure 83F.

[001062] To assess performance, Applicants compared this approach to scoring based on two other strategies:

1. *Gene list approach*: A list of hand-picked branch-specific genes were chosen from Figure 83F, and a score was defined as the sum of the normalized expression of those genes.
2. *ICA approach*: To allow for more complicated logic than simple sums, Applicants used the ICA decomposition seen in Figure 83F and computed the value of each IC for every cell in the population.

[001063] With each scoring system, Applicants subtracted the median of the DMSO-treated control cells’ and thresholded all cells with negative scores to zero. Applicants then assessed the overlap of score distributions between cells expected to have a given branch active or inactive. Both the gene list approach and the ICA approach performed worse than the random forests in this analysis (Figure 91F).

[001064] The branch scores seen in Figure 85D are the result of applying the random forest regressor scoring system to each cell in the perturb-seq experiment, and then averaging the results within cells knocked down for the same gene. The average scores assigned by the ICA method agree well (cf. Figure 85D and Figure 91E).

[001065] **Single-cell analysis in Figure 86.** Applicants formed a population of cells containing either our two guides targeting *HSPA5*, or the NegCtrl2 guide. Applicants found all genes that had mean abundance >0.5 UMI per cell and that were differentially expressed between the two populations by Kolmogorov-Smirnov test ($P < 0.01$), resulting in ~2,100 genes. Applicants formed a reduced gene expression matrix consisting only of these genes and applied low rank ICA to reduce the population’s dynamics therein to 12 ICs. The t-sne plots were made by reducing the low rank matrix to 16 components using PCA and then applying t-sne (see “t-sne

visualization” section). Branch activation scores in Figure 86C were assigned as described above in the “Scoring branch activation in Figure 85D” section.

[001066] Two ICs varied substantially in average value between the control and perturbed cells (Figure 86B). The first, IC1, had a two-phase distribution in which all control cells and the majority of HSPA5-perturbed cells fell in the large lower peak, and a subpopulation of HSPA5-perturbed cells fell into a long tail of higher values (Figure 86B). Applicants defined the HSPA5 IC1 HIGH cells to be the ones that fell within this tail (Figure 86B). Figure 86D shows the normalized expression of genes found in our epistasis analysis (Figure 83F) as columns, and the HSPA5-perturbed cells as rows, ordered by increasing IC1. Figure 86E was created simply by averaging the expression of HSPA5 within the subpopulations defined in Figure 86B. Figure 86F was created using the cell cycle positions called in the “Cell cycle position” section.

[001067] **Gene clustering analysis in Figure 86H.** Applicants first needed an unbiased approach to find programs of gene expression induced in the perturb-seq experiment. To do this Applicants separated the population into control cells (containing our two control guides) and perturbed cells (containing any guide). Applicants constructed average expression profiles (see “Average expression profiles” section) of each, and then restricted our analysis to genes of mean expression > 0.5 UMI per cell on average in the perturbed population, and whose normalized expression was > 0.5 . (Control cells by definition have mean normalized expression 0 for all genes, see “Expression normalization” section.) Applicants then used a random forest classifier approach to select 200 of these induced genes that varied informatively across all of the perturbations in the perturb-seq experiment (see “Identification of differentially expressed genes” section). The genes were then clustered based on their co-expression throughout the population, with the dendrogram leaves optimally reordered (see “Hierarchical clustering of genes” section). Our assumption was that many of these “induced genes” were involved in the unfolded protein response. Applicants evaluated UPR dependence by examining the expression pattern of the induced genes within thapsigargin- and tunicamycin-treated cells (Figure 92B). Applicants also assigned identities to some other clusters based on clear functional connections (as seen in Figure 86H).

[001068] **Comparison of clustering of UPR genes in Figure 86I and Figure 92A.** As many UPR genes fell out of the previous analysis, Applicants wanted to evaluate the ability to go the opposite direction, and cluster known interactions. Applicants thus reexamined the list of UPR-

regulated genes found in Figure 83F. Applicants separated the perturb-seq population into control cells (containing our two control guides) and perturbed cells (containing any guide). Applicants constructed average expression profiles (see “Average expression profiles” section) of each, and then restricted our analysis to the UPR-regulated genes that showed the same pattern of induction or repression in the perturbed cell population as they did in the cells treated with thapsigargin in the epistasis experiment that had all branches of the UPR intact. Applicants then performed hierarchical clustering of these genes (see “Hierarchical clustering of genes” section) using co-expression information from either (1) all cells in the epistasis experiment, (2) all cells in the perturb-seq, and (3) only control cells in the perturb-seq experiment. Applicants assessed the similarity among clusterings using the cophenetic correlation coefficient, i.e. the correlation coefficient between dendrogram distances taken over all possible pairs of genes. Closeness in cophenetic correlation thus implies that the dendrograms tend to place the same genes close to each other. The figure is meant only as a visual aid, as the cophenetic correlation relies on information beyond the linear order. The genes were roughly grouped based on their epistasis pattern in the epistasis experiment (as in Figure 83F), and then color was preserved as they were shuffled by the other two clusterings.

[001069] Enrichment of cholesterol genes in Figure 86K. Our unbiased analysis in Figure 86H contained a cluster of genes involved in cholesterol biosynthesis: *ACAT2*, *FDPS*, *FADS1*, *INSIG1*, *TMEM97*. Applicants made a “cholesterol score” by summing the normalized expression of this group of genes in each cell, and then created a subpopulation containing (1) cells with cholesterol scores at or above the 95% of the control cell population and (2) control cells. This gave ~9,000 cells. Within this subpopulation, Applicants then correlated the cholesterol score with the normalized expression of all genes with mean > 0.25 UMI per cell. Applicants then selected all genes that had a correlation of 0.15 or higher with the cholesterol score for further analysis. Applicants clustered the genes by co-expression within the population (see “Hierarchical clustering of genes” section), and then selected a group of 23 genes that clustered together with the original five and that appeared as a distinct block on the diagonal of the gene-gene correlation matrix. To demonstrate the improvement in correlation obtained by this “fishing” approach, Applicants compared correlation matrices composed of these 23 genes and 23 random genes of similar average abundance between our enriched population, and control cells (seen in Figure 86K). Finally, Applicants used Enrichr (Kuleshov et al. 2016) to obtain

Reactome annotations and Encode SREBP binding state. Note that some of the genes that don't have annotations nevertheless are almost certainly cholesterol-related, such as the lncRNA RP11-660L16 which is directly next to *DHC7R*.

[001070] **Enrichment of heat shock genes in Figure 86L.** Applicants followed an identical approach to the above, except starting with the genes *HSPA1A* and *HSPA1B*. In this case Applicants present all of the genes that had correlations of 0.15 or greater. Applicants again used Enrichr to find the top 3 most enriched transcription factor binding sites among the set of genes, as presented in Figure 6L.

[001071] **Single-cell analysis in Figure 87.** Applicants formed populations of cells containing guides targeting either *SEC61A1* or *SEC61B*, along with cells containing the NegCtrl2 guide. Applicants found all genes that had mean abundance >0.5 UMI per cell and that were differentially expressed between the two populations by Kolmogorov-Smirnov test setting a threshold of $D > 0.15$ for *SEC61A1*, and $D > 0.1$ for *SEC61B*, which is a weaker perturbation (see "Identification of differentially expressed genes" section). The different thresholds were chosen largely for esthetic reasons: lowering the threshold with *SEC61A1*, which is a strong perturbation, resulted in the inclusion of a number of cell cycle genes that caused the control population to fragment into subpopulations by cell cycle phase, which Applicants felt was distracting. In each case Applicants formed a reduced gene expression matrix consisting only of differentially expressed genes, then applied robust PCA (see "Low rank ICA" section) to these matrices, and then visualized the dynamics using t-sne plots generated using the first 16 principal components (see "t-sne visualization" section). Branch activation scores in Figure 87A,B, S6A were assigned as described above in the "Scoring branch activation in Figure 85D" section.

References

Acosta-Alvear, D., Zhou, Y., Blais, A., Tsikitis, M., Lents, N.H., Arias, C., Lennon, C.J., Kluger, Y. & Dynlacht, B.D. 2007, "XBP1 controls diverse cell type- and condition-specific transcriptional regulatory networks", *Molecular Cell*, vol. 27, no. 1, pp. 53-66.

Adiconis, X., Borges-Rivera, D., Satija, R., DeLuca, D. S., Busby, M. A., Berlin, A. M., Sivachenko, A., Thompson, D. A., Wysoker, A., Fennell, T., Gnirke, A., Pochet, N., Regev, A. & Levin, J. Z. Comparative analysis of RNA sequencing methods for degraded or low-input samples. *Nat Methods*. 10, 623-629, doi:10.1038/nmeth.2483 (2013). PMID:3821180.

Aguirre, A.J., Meyers, R.M., Weir, B.A., Vazquez, F., Zhang, C., Ben-David, U., Cook, A., Ha, G., Harrington, W.F., Doshi, M.B., et al 2016, "Genomic Copy Number Dictates a Gene-Independent Cell Response to CRISPR/Cas9 Targeting", *Cancer Discovery*, vol. 6, no. 8, pp. 914-929.

Altshuler, D., Daly, M.J., and Lander, E.S. (2008). Genetic Mapping in Human Disease. *Science* (80-.). 322, 881–888.

Amit, I., Garber, M., Chevrier, N., Leite, A.P., Donner, Y., Eisenhaure, T., Guttman, M., Grenier, J.K., Li, W., Zuk, O., et al. (2009). Unbiased Reconstruction of a Mammalian Transcriptional Network Mediating Pathogen Responses. *Science* (80-.). 326, 257–263.

Anders, C., Niewoehner, O., Duerst, A. & Jinek, M. 2014, "Structural basis of PAM-dependent target DNA recognition by the Cas9 endonuclease", *Nature*, vol. 513, no. 7519, pp. 569-573.

Assarsson, E., Lundberg, M., Holmquist, G., Bjorkestén, J., Thorsen, S.B., Ekman, D., Eriksson, A., Rennel Dickens, E., Ohlsson, S., Edfeldt, G., et al. (2014). Homogenous 96-plex PEA immunoassay exhibiting high sensitivity, specificity, and excellent scalability. *PloS one* 9, e95192.

Bamshad, M.J., Ng, S.B., Bigham, A.W., Tabor, H.K., Emond, M.J., Nickerson, D.A., and Shendure, J. (2011). Exome sequencing as a tool for Mendelian disease gene discovery. *Nat. Rev. Genet.* 12, 745–755.

Bandyopadhyay, S., Mehta, M., Kuo, D., Sung, M.-K., Chuang, R., Jaehnig, E.J., Bodenmiller, B., Licon, K., Copeland, W., Shales, M., et al. (2010). Rewiring of Genetic Networks in Response to DNA Damage. *Science* (80-.). 330, 1385–1389.

Bao, X.R., Ong, S., Goldberger, O., Peng, J., Sharma, R., Thompson, D.A., Vafai, S.B., Cox, A.G., Marutani, E., Ichinose, F., et al 2016, "Mitochondrial dysfunction remodels one-carbon metabolism in human cells", *eLife*, vol. 5.

Bassik, M.C., Kampmann, M., Lebbink, R.J., Wang, S., Hein, M.Y., Poser, I., Weibezahn, J., Horlbeck, M. a, Chen, S., Mann, M., et al. (2013a). A systematic mammalian genetic interaction map reveals pathways underlying ricin susceptibility. *Cell* 152, 909–922.

Bassik, M.C., Kampmann, M., Lebbink, R.J., Wang, S., Hein, M.Y., Poser, I., Weibezahn, J., Horlbeck, M.A., Chen, S., Mann, M., et al. (2013b). A Systematic Mammalian Genetic Interaction Map Reveals Pathways Underlying Ricin Susceptibility. *Cell* 152, 909–922.

Beerenwinkel, N., Pachter, L., and Sturmfels, B. (2007). Epistasis and Shapes of Fitness Landscapes. *Stat. Sin.* 17, 1317–1342.

Bendall, S.C., Simonds, E.F., Qiu, P., Amir el, A.D., Krutzik, P.O., Finck, R., Bruggner, R.V., Melamed, R., Trejo, A., Ornatsky, O.I., *et al.* (2011). Single-cell mass cytometry of differential immune and drug responses across a human hematopoietic continuum. *Science* 332, 687-696.

Berger, A.H., Brooks, A.N., Wu, X., Shrestha, Y., Chouinard, C., Piccioni, F., Bagul, M., Kamburov, A., Imielinski, M., Hogstrom, L., *et al.* (2016). High-throughput Phenotyping of Lung Cancer Somatic Mutations. *Cancer Cell* 0, 248–249.

Blecher-Gonen, R., Barnett-Itzhaki, Z., Jaitin, D., Amann-Zalcenstein, D., Lara-Astiaso, D. & Amit, I. High-throughput chromatin immunoprecipitation for genome-wide mapping of *in vivo* protein-DNA interactions and epigenomic states. *Nat Protoc.* 8, 539-554, doi:10.1038/nprot.2013.023 (2013).

Bochkis, I. M., Przybylski, D., Chen, J. & Regev, A. Changes in nucleosome occupancy associated with metabolic alterations in aged mammalian liver. *Cell reports.* 9, 996-1006, doi:10.1016/j.celrep.2014.09.048 (2014). PMID:4250828.

Boone, C., Bussey, H., and Andrews, B.J. (2007). Exploring genetic interactions and networks with yeast. *8*, 437–449.

Bornstein, C., Winter, D., Barnett-Itzhaki, Z., David, E., Kadri, S., Garber, M. & Amit, I. A negative feedback loop of transcription factors specifies alternative dendritic cell chromatin States. *Mol Cell.* 56, 749-762, doi:10.1016/j.molcel.2014.10.014 (2014). PMID:4412443.

Botstein, D., and Risch, N. (2003). Discovering genotypes underlying human phenotypes: past successes for mendelian disease, future approaches for complex disease. *Nat. Genet.* 33, 228–237.

Briner, A.E., Donohoue, P.D., Gomaa, A.A., Selle, K., Slorach, E.M., Nye, C.H., Haurwitz, R.E., Beisel, C.L., May, A.P. & Barrangou, R. 2014, "Guide RNA functional modules direct Cas9 activity and orthogonality", *Molecular Cell*, vol. 56, no. 2, pp. 333-339.

Buettner, F., Natarajan, K.N., Casale, F.P., Proserpio, V., Scialdone, A., Theis, F.J., Teichmann, S.A., Marioni, J.C., and Stegle, O. (2015). Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat. Biotechnol.* 33, 155–160.

Cabili, M. N., Dunagin, M. C., McClanahan, P. D., Biaesch, A., Padovan-Merhar, O., Regev, A., Rinn, J. L. & Raj, A. Localization and abundance analysis of human lncRNAs at single-cell and single-molecule resolution. *Genome Biol.* 16, 20, doi:10.1186/s13059-015-0586-4 (2015). PMID:4369099.

Cabili, M. N., Trapnell, C., Goff, L., Koziol, M., Tazon-Vega, B., Regev, A. & Rinn, J. L. Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes Dev.* 25, 1915-1927, doi:10.1101/gad.17446611 (2011). PMID:3185964.

Calfon, M., Zeng, H., Urano, F., Till, J.H., Hubbard, S.R., Harding, H.P., Clark, S.G. & Ron, D. 2002, "IRE1 couples endoplasmic reticulum load to secretory capacity by processing the XBP-1 mRNA", *Nature*, vol. 415, no. 6867, pp. 92-96.

Candès, E.J., and Recht, B. (2009). Exact Matrix Completion via Convex Optimization. *Found. Comput. Math.* 9, 717–772.

Candès, E.J., Li, X., Ma, Y. & Wright, J. 2011, "Robust principal component analysis?", *Journal of the ACM (JACM)*, vol. 58, no. 3, pp. 11.

Capaldi, A.P., Kaplan, T., Liu, Y., Habib, N., Regev, A., Friedman, N., and O'Shea, E.K. (2008). Structure and function of a transcriptional network activated by the MAPK Hog1. *Nat. Genet.* 40, 1300–1306.

Carter, G.W., Prinz, S., Neou, C., Shelby, J.P., Marzolf, B., Thorsson, V., and Galitski, T. (2007). Prediction of phenotype and gene expression for combinations of mutations. *Mol. Syst. Biol.* 3, 96.

Carbon, S., Ireland, A., Mungall, C.J., Shu, S., Marshall, B. & Lewis, S. 2009, "AmiGO: online access to ontology and annotation data", *Bioinformatics (Oxford, England)*, vol. 25, no. 2, pp. 288-289.

Cartwright, T., Perkins, N.D., and L Wilson, C. (2016). NFKB1: a suppressor of inflammation, ageing and cancer. *FEBS J.* 283, 1812–1822.

Chan, M. M., Smith, Z. D., Egli, D., Regev, A. & Meissner, A. Mouse ooplasm confers context-specific reprogramming capacity. *Nature genetics.* 44, 978-980, doi:10.1038/ng.2382 (2012). PMID:3432711.

Chavez, A., Scheiman, J., Vora, S., Pruitt, B.W., Tuttle, M., P R Iyer, E., Lin, S., Kiani, S., Guzman, C.D., Wiegand, D.J., et al 2015, "Highly efficient Cas9-mediated transcriptional

programming", *Nature Methods*, vol. 12, no. 4, pp. 326-328.

Chen, Y., Liu, P., Nielsen, A.A.K., Brophy, J.A.N., Clancy, K., Peterson, T. & Voigt, C.A. 2013, "Characterization of 582 natural and synthetic terminators and quantification of their design constraints", *Nature Methods*, vol. 10, no. 7, pp. 659-664.

Chen, S., Sanjana, N.E., Zheng, K., Shalem, O., Lee, K., Shi, X., Scott, D.A., Song, J., Pan, J.Q., Weissleder, R., et al. (2015). Genome-wide CRISPR Screen in a Mouse Model of Tumor Growth and Metastasis. *Cell* 160, 1246–1260.

Cheng, C. S., Rai, K., Garber, M., Hollinger, A., Robbins, D., Anderson, S., Macbeth, A., Tzou, A., Carneiro, M. O., Raychowdhury, R., Russ, C., Hacohen, N., Gershenwald, J. E., Lennon, N., Nusbaum, C., Chin, L., Regev, A. & Amit, I. Semiconductor-based DNA sequencing of histone modification states. *Nat Commun.* 4, 2672, doi:10.1038/ncomms3672 (2013). PMID:3917140.

Chevrier, N., Mertins, P., Artyomov, M. N., Shalek, A. K., Iannaccone, M., Ciaccio, M. F., Gat-Viks, I., Tonti, E., DeGrace, M. M., Clauser, K. R., Garber, M., Eisenhaure, T. M., Yosef, N., Robinson, J., Sutton, A., Andersen, M. S., Root, D. E., von Andrian, U., Jones, R. B., Park, H., Carr, S. A., Regev, A., Amit, I. & Hacohen, N. Systematic discovery of TLR signaling components delineates viral-sensing circuits. *Cell.* 147, 853-867, doi:10.1016/j.cell.2011.10.022 (2011). PMID:3809888.

Chuang, C., Lee, K., Fan, C. & Su, Y. 2009, "Porcine type III RNA polymerase III promoters for short hairpin RNA expression", *Animal Biotechnology*, vol. 20, no. 1, pp. 34-39.

Chung, K., Wallace, J., Kim, S.Y., Kalyanasundaram, S., Andalman, A.S., Davidson, T.J., Mirzabekov, J.J., Zalocusky, K.A., Mattis, J., Denisin, A.K., et al. (2013). Structural and molecular interrogation of intact biological systems. *Nature* 497, 332-337.

Chung, N.C., and Storey, J.D. (2015). Statistical significance of variables driving systematic variation in high-dimensional data. *Bioinformatics* 31, 545–554.

Cohen, A., and Sheva, B. - (1998). HIDDEN MARKOV MODELS IN BIOMEDICAL SIGNAL PROCESSING. 20, 1145–1150.

Cong, L., Ran, F.A., Cox, D., Lin, S., Barretto, R., Habib, N., Hsu, P.D., Wu, X., Jiang, W., Marraffini, L. a, et al. (2013). Multiplex genome engineering using CRISPR/Cas systems. *Science* 339, 819–823.

Costanzo, M., Baryshnikova, A., Bellay, J., Kim, Y., Spear, E.D., Sevier, C.S., Ding, H.,

Koh, J.L.Y., Toufighi, K., Mostafavi, S., et al. (2010). The Genetic Landscape of a Cell. *Science* (80-.). 327, 425–431.

Dang, Y., Jia, G., Choi, J., Ma, H., Anaya, E., Ye, C., Shankar, P. & Wu, H. 2015, "Optimizing sgRNA structure to improve CRISPR-Cas9 knockout efficiency", *Genome Biology*, vol. 16, pp. 280.

Duan, Q., Flynn, C., Niepel, M., Hafner, M., Muhlich, J.L., Fernandez, N.F., Rouillard, A.D., Tan, C.M., Chen, E.Y., Golub, T.R., et al. (2014). LINCS Canvas Browser: interactive web app to query, browse and interrogate LINCS L1000 gene expression signatures. *Nucleic Acids Res.* 42, W449–W460.

Elsharkawy, A.M., Oakley, F., Lin, F., Packham, G., Mann, D.A., and Mann, J. (2010). The NF-kappaB p50:p50:HDAC-1 repressor complex orchestrates transcriptional inhibition of multiple pro-inflammatory genes. *J. Hepatol.* 53, 519–527.

Engreitz, J. M., Pandya-Jones, A., McDonel, P., Shishkin, A., Sirokman, K., Surka, C., Kadri, S., Xing, J., Goren, A., Lander, E. S., Plath, K. & Guttman, M. The Xist lncRNA exploits three-dimensional genome architecture to spread across the X chromosome. *Science*. 341, 1237973, doi:10.1126/science.1237973 (2013). PMID:3778663.

Fan, H.C., Fu, G.K., and Fodor, S.P. a. (2015). Combinatorial labeling of single cells for gene expression cytometry. *Science*. 347, 1258367–1258367.

Fogli, A. & Boespflug-Tanguy, O. 2006, "The large spectrum of eIF2B-related diseases", *Biochemical Society Transactions*, vol. 34, no. Pt 1, pp. 22-29.

Friedman, J., Hastie, T. & Tibshirani, R. 2001, *The elements of statistical learning*, Springer series in statistics Springer, Berlin.

Galonska, C., Smith, Z. D. & Meissner, A. In Vivo and in vitro dynamics of undifferentiated embryonic cell transcription factor 1. *Stem Cell Reports*. 2, 245-252, doi:10.1016/j.stemcr.2014.01.007 (2014). PMID:3964277.

Garber, M., Yosef, N., Goren, A., Raychowdhury, R., Thielke, A., Guttman, M., Robinson, J., Minie, B., Chevrier, N., Itzhaki, Z., et al. (2012). A High-Throughput Chromatin Immunoprecipitation Approach Reveals Principles of Dynamic Gene Regulation in Mammals. *Mol. Cell* 47, 810–822.

Gat-Viks, I., Chevrier, N., Wilentzik, R., Eisenhaure, T., Raychowdhury, R., Steuerman, Y., Shalek, A. K., Hacohen, N., Amit, I. & Regev, A. Deciphering molecular circuits from

genetic variation underlying transcriptional responsiveness to stimuli. *Nature biotechnology*. 31, 342-349, doi:10.1038/nbt.2519 (2013). PMID:3622156.

Gilbert, L.A., Horlbeck, M.A., Adamson, B., Villalta, J.E., Chen, Y., Whitehead, E.H., Guimaraes, C., Panning, B., Ploegh, H.L., Bassik, M.C., et al. (2014). Genome-Scale CRISPR-Mediated Control of Gene Repression and Activation. *Cell* 159, 647–661.

Gilbert, L. a, Larson, M.H., Morsut, L., Liu, Z., Brar, G. a, Torres, S.E., Stern-Ginossar, N., Brandman, O., Whitehead, E.H., Doudna, J. a, et al. (2013). CRISPR-mediated modular RNA-guided regulation of transcription in eukaryotes. *Cell* 154, 442–451.

Gomez, D., Shankman, L.S., Nguyen, A.T., and Owens, G.K. (2013). Detection of histone modifications at specific gene loci in single cells in histological sections. *Nature methods* 10, 171-177.

Goodwin, E.C. & Rottman, F.M. 1992, "The 3'-flanking sequence of the bovine growth hormone gene contains novel elements required for efficient and accurate polyadenylation", *The Journal of Biological Chemistry*, vol. 267, no. 23, pp. 16330-16334.

Grün, D. & van Oudenaarden, A. 2015, "Design and Analysis of Single-Cell Sequencing Experiments", *Cell*, vol. 163, no. 4, pp. 799-810.

Gu, W., Crawford, E.D., O'Donovan, B.D., Wilson, M.R., Chow, E.D., Retallack, H. & DeRisi, J.L. 2016, "Depletion of Abundant Sequences by Hybridization (DASH): using Cas9 to remove unwanted high-abundance species in sequencing libraries and molecular counting applications", *Genome Biology*, vol. 17, pp. 41.

Guttman, M., Donaghey, J., Carey, B. W., Garber, M., Grenier, J. K., Munson, G., Young, G., Lucas, A. B., Ach, R., Bruhn, L., Yang, X., Amit, I., Meissner, A., Regev, A., Rinn, J. L., Root, D. E. & Lander, E. S. lincRNAs act in the circuitry controlling pluripotency and differentiation. *Nature*. 477, 295-300, doi:10.1038/nature10398 (2011). PMID:3175327.

Haas, B. J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P. D., Bowden, J., Couger, M. B., Eccles, D., Li, B., Lieber, M., Macmanes, M. D., Ott, M., Orvis, J., Pochet, N., Strozzi, F., Weeks, N., Westerman, R., William, T., Dewey, C. N., Henschel, R., Leduc, R. D., Friedman, N. & Regev, A. De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nat Protoc*. 8, 1494-1512, doi:10.1038/nprot.2013.084 (2013). PMID:3875132.

Haber, J.E., Braberg, H., Wu, Q., Alexander, R., Haase, J., Ryan, C., Lipkin-Moore, Z.,

Franks-Skiba, K.E., Johnson, T., Shales, M., et al. (2013). Systematic triple-mutant analysis uncovers functional connectivity between pathways involved in chromosome regulation. *Cell Rep.* 3, 2168–2178.

Hacisuleyman, E., Goff, L. A., Trapnell, C., Williams, A., Henao-Mejia, J., Sun, L., McClanahan, P., Hendrickson, D. G., Sauvageau, M., Kelley, D. R., Morse, M., Engreitz, J., Lander, E. S., Guttman, M., Lodish, H. F., Flavell, R., Raj, A. & Rinn, J. L. Topological organization of multichromosomal regions by the long intergenic noncoding RNA Firre. *Nat Struct Mol Biol.* 21, 198-206, doi:10.1038/nsmb.2764 (2014). PMID:3950333.

Haldimann, A. & Wanner, B.L. 2001, "Conditional-replication, integration, excision, and retrieval plasmid-host systems for gene structure-function studies of bacteria", *Journal of Bacteriology*, vol. 183, no. 21, pp. 6384-6393.

Hamanaka, R.B., Bennett, B.S., Cullinan, S.B. & Diehl, J.A. 2005, "PERK and GCN2 Contribute to eIF2 α Phosphorylation and Cell Cycle Arrest after Activation of the Unfolded Protein Response Pathway", *Molecular Biology of the Cell*, vol. 16, no. 12, pp. 5493-5501.

Han, J., Back, S.H., Hur, J., Lin, Y., Gildersleeve, R., Shan, J., Yuan, C.L., Krokowski, D., Wang, S., Hatzoglou, M., et al 2013, "ER-stress-induced transcriptional regulation increases protein synthesis leading to cell death", *Nature Cell Biology*, vol. 15, no. 5, pp. 481-490.

Hartl, D.L. (2014). What can Applicants learn from fitness landscapes? *Curr. Opin. Microbiol.* 21, 51–57.

Heckl, D., Kowalczyk, M. S., Yudovich, D., Belizaire, R., Puram, R. V., McConkey, M. E., Thielke, A., Aster, J. C., Regev, A. & Ebert, B. L. Generation of mouse models of myeloid malignancy with combinatorial genetic lesions using CRISPR-Cas9 genome editing. *Nature biotechnology.* 32, 941-946, doi:10.1038/nbt.2951 (2014). PMID:4160386.

Heimberg, G., Bhatnagar, R., El-Samad, H., and Thomson, M. (2016). Low Dimensionality in Gene Expression Data Enables the Accurate Extraction of Transcriptional Programs from Shallow Sequencing. *Cell Syst.* 2, 239–250.

Helft, J., Böttcher, J., Chakravarty, P., Zelenay, S., Huotari, J., Schraml, B.U., Goubau, D., and Reis e Sousa, C. (2015). GM-CSF Mouse Bone Marrow Cultures Comprise a Heterogeneous Population of CD11c+MHCII+ Macrophages and Dendritic Cells. *Immunity* 42, 1197–1211.

Hetz, C. 2012, "The unfolded protein response: controlling cell fate decisions under ER

stress and beyond", *Nature Reviews. Molecular Cell Biology*, vol. 13, no. 2, pp. 89-102.

Horlbeck, M.A., Gilbert, L.A., Villalta, J.E., Adamson, B., Pak, R.A., Chen, Y., Fields, A.P., Park, C.Y., Corn, J.E. & Kampmann, M. 2016, "Compact and highly active next-generation libraries for CRISPR-mediated gene repression and activation", *eLife*, vol. 5, pp. e19760.

Hu, S., Ni, W., Hazi, W., Zhang, H., Zhang, N., Meng, R. & Chen, C. 2011, "Cloning and functional analysis of sheep U6 promoters", *Animal Biotechnology*, vol. 22, no. 3, pp. 170-174.

Huang, D.W., Sherman, B.T. & Lempicki, R.A. 2009a, "Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists", *Nucleic Acids Research*, vol. 37, no. 1, pp. 1-13.

Huang, D.W., Sherman, B.T. & Lempicki, R.A. 2009b, "Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources", *Nature Protocols*, vol. 4, no. 1, pp. 44-57.

Hughes, T.R., Marton, M.J., Jones, A.R., Roberts, C.J., Stoughton, R., Armour, C.D., Bennett, H. a, Coffey, E., Dai, H., He, Y.D., et al. (2000). Functional Discovery via a Compendium of Expression Profiles. *Cell* 102, 109–126.

Jaitin, D. A., Kenigsberg, E., Keren-Shaul, H., Elefant, N., Paul, F., Zaretsky, I., Mildner, A., Cohen, N., Jung, S., Tanay, A. & Amit, I. Massively parallel single-cell RNA-seq for marker-free decomposition of tissues into cell types. *Science*. 343, 776-779, doi:10.1126/science.1247651 (2014). PMID:4412462.

Janssen, K.P., Knez, K., Spasic, D., and Lammertyn, J. (2013). Nucleic acids for ultra-sensitive protein detection. *Sensors* 13, 1353-1384.

Jiang, F., Zhou, K., Ma, L., Gressel, S. & Doudna, J.A. 2015, "STRUCTURAL BIOLOGY. A Cas9-guide RNA complex preorganized for target DNA recognition", *Science (New York, N.Y.)*, vol. 348, no. 6242, pp. 1477-1481.

Jin, F., Hazbun, T., Michaud, G.A., Salcius, M., Predki, P.F., Fields, S., and Huang, J. (2006). A pooling-deconvolution strategy for biological network elucidation. *Nat Methods* 3, 183–189.

Joensson, H.N., and Andersson Svahn, H. (2012). Droplet Microfluidics-A Tool for Single-Cell Analysis. *Angew. Chemie Int. Ed.* 51, 12176–12192.

Jonikas, M.C., Collins, S.R., Denic, V., Oh, E., Quan, E.M., Schmid, V., Weibezahn, J.,

Schwappach, B., Walter, P., Weissman, J.S., et al 2009, "Comprehensive characterization of genes required for protein folding in the endoplasmic reticulum", *Science (New York, N.Y.)*, vol. 323, no. 5922, pp. 1693-1697.

Jovanovic, M., Rooney, M. S., Mertins, P., Przybylski, D., Chevrier, N., Satija, R., Rodriguez, E. H., Fields, A. P., Schwartz, S., Raychowdhury, R., Mumbach, M. R., Eisenhaure, T., Rabani, M., Gennert, D., Lu, D., Delorey, T., Weissman, J. S., Carr, S. A., Hacohen, N. & Regev, A. Dynamic profiling of the protein life cycle in response to pathogens. *Science*. 347, 1259038, doi:10.1126/science.1259038 (2015). PMID:PMC Journal - In Process.

Jovanovic, M., Rooney, M.S., Mertins, P., Przybylski, D., Chevrier, N., Satija, R., Rodriguez, E.H., Fields, A.P., Schwartz, S., Raychowdhury, R., et al. (2015). Dynamic profiling of the protein life cycle in response to pathogens. *Science* (80-.). 347, 1259038–1259038.

Kabadi, A.M., Ousterout, D.G., Hilton, I.B., and Gersbach, C.A. (2014). Multiplex CRISPR/Cas9-based genome engineering from a single lentiviral vector. *Nucleic Acids Res.* 42, 1–11.

Kampmann, M., Bassik, M.C., and Weissman, J.S. (2014). Functional genomics platform for pooled screening and generation of mammalian genetic interaction maps. *Nat. Protoc.* 9, 1825–1847.

Kanda, S., Yanagitani, K., Yokota, Y., Esaki, Y. & Kohno, K. 2016, "Autonomous translational pausing is required for XBP1u mRNA recruitment to the ER via the SRP pathway", *Proceedings of the National Academy of Sciences of the United States of America*, vol. 113, no. 40, pp. E5895.

Kantlehner, M., Kirchner, R., Hartmann, P., Ellwart, J.W., Alunni-Fabbroni, M., and Schumacher, A. (2011). A high-throughput DNA methylation analysis of a single cell. *Nucleic acids research* 39, e44.

Kearns, N. A., Genga, R. M., Enuameh, M. S., Garber, M., Wolfe, S. A. & Maehr, R. Cas9 effector-mediated regulation of transcription and differentiation in human pluripotent stem cells. *Development*. 141, 219-223, doi:10.1242/dev.103341 (2014). PMID:3865759.

Kelley, D. & Rinn, J. Transposable elements reveal a stem cell-specific class of long noncoding RNAs. *Genome Biol.* 13, R107, doi:10.1186/gb-2012-13-11-r107 (2012). PMID:3580499.

Kemmeren, P., Sameith, K., Van De Pasch, L.A.L., Benschop, J.J., Lenstra, T.L.,

Margaritis, T., O'Duibhir, E., Apweiler, E., Van Wageningen, S., Ko, C.W., et al. (2014). Large-scale genetic perturbations reveal regulatory networks and an abundance of gene-specific repressors. *Cell* 157, 740–752.

Kharchenko, P. V., Silberstein, L., and Scadden, D.T. (2014). Bayesian approach to single-cell differential expression analysis. *Nat. Methods* 11, 740–742.

Klein, A.M., Mazutis, L., Akartuna, I., Tallapragada, N., Veres, A., Li, V., Peshkin, L., Weitz, D.A., and Kirschner, M.W. (2015). Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell* 161, 1187–1201.

Komor, A.C., Kim, Y.B., Packer, M.S., Zuris, J.A. & Liu, D.R. 2016, "Programmable editing of a target base in genomic DNA without double-stranded DNA cleavage", *Nature*, vol. 533, no. 7603, pp. 420-424.

Konermann, S., Brigham, M.D., Trevino, A.E., Joung, J., Abudayyeh, O.O., Barcena, C., Hsu, P.D., Habib, N., Gootenberg, J.S., Nishimasu, H., et al. (2014). Genome-scale transcriptional activation by an engineered CRISPR-Cas9 complex. *Nature* 517, 583–588.

Kowalczyk, M.S., Tirosh, I., Heckl, D., Rao, T.N., Dixit, A., Haas, B.J., Schneider, R.K., Wagers, A.J., Ebert, B.L., and Regev, A. (2015). Single-cell RNA-seq reveals changes in cell cycle and differentiation programs upon aging of hematopoietic stem cells. *Genome Res.* 25, 1860–1872.

Kuleshov, M.V., Jones, M.R., Rouillard, A.D., Fernandez, N.F., Duan, Q., Wang, Z., Koplev, S., Jenkins, S.L., Jagodnik, K.M., Lachmann, A., et al 2016, "Enrichr: a comprehensive gene set enrichment analysis web server 2016 update", *Nucleic Acids Research*, vol. 44, no. W1, pp. 90.

Kumar, R. M., Cahan, P., Shalek, A. K., Satija, R., DaleyKeyser, A. J., Li, H., Zhang, J., Pardee, K., Gennert, D., Trombetta, J. J., Ferrante, T. C., Regev, A., Daley, G. Q. & Collins, J. J. Deconstructing transcriptional heterogeneity in pluripotent stem cells. *Nature*. 516, 56-61, doi:10.1038/nature13920 (2014). PMID:4256722.

Labzin, L.I., Schmidt, S. V, Masters, S.L., Beyer, M., Krebs, W., Klee, K., Stahl, R., Lütjohann, D., Schultze, J.L., Latz, E., et al. (2015). ATF3 Is a Key Regulator of Macrophage IFN Responses. *J. Immunol.* 195, 4446–4455.

Lamb, J., Crawford, E.D., Peck, D., Modell, J.W., Blat, I.C., Wrobel, M.J., Lerner, J., Brunet, J., Subramanian, A., Ross, K.N., et al. (2006). The Connectivity Map : Using. *Science*

(80-). 313, 1929–1935.

Lambeth, L.S., Wise, T.G., Moore, R.J., Muralitharan, M.S. & Doran, T.J. 2006, "Comparison of bovine RNA polymerase III promoters for short hairpin RNA expression", *Animal Genetics*, vol. 37, no. 4, pp. 369-372.

Lara-Astiaso, D., Weiner, A., Lorenzo-Vivas, E., Zaretsky, I., Jaitin, D. A., David, E., Keren-Shaul, H., Mildner, A., Winter, D., Jung, S., Friedman, N. & Amit, I. Immunogenetics. Chromatin state dynamics during blood formation. *Science*. 345, 943-949, doi:10.1126/science.1256271 (2014). PMID:4412442.

Laufer, C., Fischer, B., Billmann, M., Huber, W., and Boutros, M. (2013). Mapping genetic interactions in human cancer cells with RNAi and multiparametric phenotyping. *Nat. Methods* 10, 427–431.

Lawrence, M.S., Stojanov, P., Mermel, C.H., Robinson, J.T., Garraway, L.A., Golub, T.R., Meyerson, M., Gabriel, S.B., Lander, E.S., and Getz, G. (2014). Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature* 505, 495–501.

Lee, A., Iwakoshi, N.N. & Glimcher, L.H. 2003, "XBP-1 regulates a subset of endoplasmic reticulum resident chaperone genes in the unfolded protein response", *Molecular and Cellular Biology*, vol. 23, no. 21, pp. 7448-7459.

Lee, M. N., Ye, C., Villani, A. C., Raj, T., Li, W., Eisenhaure, T. M., Imboywa, S. H., Chipendo, P. I., Ran, F. A., Slowikowski, K., Ward, L. D., Raddassi, K., McCabe, C., Lee, M. H., Frohlich, I. Y., Hafler, D. A., Kellis, M., Raychaudhuri, S., Zhang, F., Stranger, B. E., Benoist, C. O., De Jager, P. L., Regev, A. & Hacohen, N. Common genetic variants modulate pathogen-sensing responses in human dendritic cells. *Science*. 343, 1246980, doi:10.1126/science.1246980 (2014). PMID:4124741.

Liang, S., Zhang, W., McGrath, B.C., Zhang, P. & Cavener, D.R. 2006, "PERK (eIF2alpha kinase) is required to activate the stress-activated MAPKs and induce the expression of immediate-early genes upon disruption of ER calcium homeostasis", *The Biochemical Journal*, vol. 393, no. Pt 1, pp. 201-209.

Liberali, P., Snijder, B., and Pelkmans, L. (2014). Single-cell and multivariate approaches in genetic perturbation screens. *Nat. Rev. Genet.* 16, 18–32.

Lin, J.H., Li, H., Yasumura, D., Cohen, H.R., Zhang, C., Panning, B., Shokat, K.M., Lavail, M.M. & Walter, P. 2007, "IRE1 signaling affects cell fate during the unfolded protein

response", *Science (New York, N.Y.)*, vol. 318, no. 5852, pp. 944-949.

Lin, Z., Chen, M. & Ma, Y. 2010, "The augmented lagrange multiplier method for exact recovery of corrupted low-rank matrices", *arXiv preprint arXiv:1009.5055*, .

Lorthongpanich, C., Cheow, L.F., Balu, S., Quake, S.R., Knowles, B.B., Burkholder, W.F., Solter, D., and Messerschmidt, D.M. (2013). Single-cell DNA-methylation analysis reveals epigenetic chimerism in preimplantation embryos. *Science* 341, 1110-1112.

Lutz, R. & Bujard, H. 1997, "Independent and tight regulation of transcriptional units in *Escherichia coli* via the LacR/O, the TetR/O and AraC/I1-I2 regulatory elements", *Nucleic Acids Research*, vol. 25, no. 6, pp. 1203-1210.

Macosko, E.Z., Basu, A., Satija, R., Nemes, J., Shekhar, K., Goldman, M., Tirosh, I., Bialas, A.R., Kamitaki, N., Martersteck, E.M., et al. (2015). Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* 161, 1202–1214.

Makarova, K.S., Haft, D.H., Barrangou, R., Brouns, S.J.J., Charpentier, E., Horvath, P., Moineau, S., Mojica, F.J.M., Wolf, Y.I., Yakunin, A.F., et al. (2011). Evolution and Classification of the CRISPR–Cas Systems. *Nat. Rev. Microbiol.* 9, 467–477.

Manolio, T.A., Collins, F.S., Cox, N.J., Goldstein, D.B., Hindorff, L.A., Hunter, D.J., McCarthy, M.I., Ramos, E.M., Cardon, L.R., Chakravarti, A., et al. (2009). Finding the missing heritability of complex diseases. *Nature* 461, 747–753.

Martincorena, I., and Campbell, P.J. (2015). Somatic mutation in cancer and normal cells. *Science (80-.)*. 349, 1483–1489.

Meerbrey, K.L., Hu, G., Kessler, J.D., Roarty, K., Li, M.Z., Fang, J.E., Herschkowitz, J.I., Burrows, A.E., Ciccio, A., Sun, T., et al 2011, "The pINDUCER lentiviral toolkit for inducible RNA interference in vitro and in vivo", *Proceedings of the National Academy of Sciences of the United States of America*, vol. 108, no. 9, pp. 3665-3670.

Meier, J.A., and Larner, A.C. (2014). Toward a new STATE: the role of STATs in mitochondrial function. *Semin. Immunol.* 26, 20–28.

Melnikov, A., Murugan, A., Zhang, X., Tesileanu, T., Wang, L., Rogov, P., Feizi, S., Gnirke, A., Callan, C.G., Kinney, J.B., et al. (2012). Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay. *Nat. Biotechnol.* 30, 271–277.

Müller-Kuller, U., Ackermann, M., Kolodziej, S., Brendel, C., Fritsch, J., Lachmann, N.,

Kunkel, H., Lausen, J., Schambach, A., Moritz, T., et al 2015, "A minimal ubiquitous chromatin opening element (UCOE) effectively prevents silencing of juxtaposed heterologous promoters by epigenetic remodeling in multipotent and pluripotent stem cells", *Nucleic Acids Research*, , pp. gkv019.

Munoz, D.M., Cassiani, P.J., Li, L., Billy, E., Korn, J.M., Jones, M.D., Golji, J., Ruddy, D.A., Yu, K., McAllister, G., et al 2016, "CRISPR Screens Provide a Comprehensive Assessment of Cancer Vulnerabilities but Generate False-Positive Hits for Highly Amplified Genomic Regions", *Cancer Discovery*, vol. 6, no. 8, pp. 900-913.

Na, Y. R., Kim, S. Y., Gaublomme, J. T., Shalek, A. K., Jorgolli, M., Park, H. & Yang, E. G. Probing enzymatic activity inside living cells using a nanowire-cell "sandwich" assay. *Nano Lett.* 13, 153-158, doi:10.1021/nl3037068 (2013). PMID:3541459.

Nagano, T., Lubling, Y., Stevens, T.J., Schoenfelder, S., Yaffe, E., Dean, W., Laue, E.D., Tanay, A., and Fraser, P. (2013). Single-cell Hi-C reveals cell-to-cell variability in chromosome structure. *Nature* 502, 59-64.

Neumann, B., Walter, T., Hériché, J.-K., Bulkescher, J., Erfle, H., Conrad, C., Rogers, P., Poser, I., Held, M., Liebel, U., et al. (2010). Phenotypic profiling of the human genome by time-lapse microscopy reveals cell division genes. *Nature* 464, 721–727.

Nishimasu, H., Ran, F.A., Hsu, P.D., Konermann, S., Shehata, S.I., Dohmae, N., Ishitani, R., Zhang, F. & Nureki, O. 2014, "Crystal structure of Cas9 in complex with guide RNA and target DNA", *Cell*, vol. 156, no. 5, pp. 935-949.

Nissim, L., Perli, S.D., Fridkin, A., Perez-Pinera, P. & Lu, T.K. 2014, "Multiplexed and programmable regulation of gene networks with an integrated RNA and CRISPR/Cas toolkit in human cells", *Molecular Cell*, vol. 54, no. 4, pp. 698-710.

Okabe, Y., and Medzhitov, R. (2014). Tissue-Specific Signals Control Reversible Program of Localization and Functional Polarization of Macrophages. *Cell* 157, 832–844.

Pardon, E., Laeremans, T., Triest, S., Rasmussen, S.G., Wohlkonig, A., Ruf, A., Muyldermans, S., Hol, W.G., Kobilka, B.K., and Steyaert, J. (2014). A general protocol for the generation of Nanobodies for structural biology. *Nature protocols* 9, 674-693.

Parnas, O., Jovanovic, M., Eisenhaure, T.M., Herbst, R.H., Dixit, A., Ye, C.J., Przybylski, D., Platt, R.J., Tirosh, I., Sanjana, N.E., et al. (2015). A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks. *Cell* 162, 675–686.

Perfetto, S.P., Chattopadhyay, P.K., and Roederer, M. (2004). Seventeen-colour flow cytometry: unravelling the immune system. *Nature reviews Immunology* 4, 648-655.

Phillips, P.C. (2008). Epistasis — the essential role of gene interactions in the structure and evolution of genetic systems.

Platt, R.J., Chen, S., Zhou, Y., Yim, M.J., Swiech, L., Kempton, H.R., Dahlman, J.E., Parnas, O., Eisenhaure, T.M., Jovanovic, M., et al. (2014). CRISPR-Cas9 Knockin Mice for Genome Editing and Cancer Modeling. *Cell* 159, 440–455.

Plumb, R., Zhang, Z., Appathurai, S. & Mariappan, M. 2015, "A functional link between the co-translational protein translocation pathway and the UPR", *eLife*, vol. 4.

Pollen, A.A., Nowakowski, T.J., Shuga, J., Wang, X., Leyrat, A.A., Lui, J.H., Li, N., Szpankowski, L., Fowler, B., Chen, P., et al 2014, "Low-coverage single-cell mRNA sequencing reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex", *Nature Biotechnology*, vol. 32, no. 10, pp. 1053-1058.

Qi, L.S., Larson, M.H., Gilbert, L.A., Doudna, J.A., Weissman, J.S., Arkin, A.P., and Lim, W.A. (2013). Repurposing CRISPR as an RNA-Guided Platform for Sequence-Specific Control of Gene Expression. *Cell* 152, 1173–1183.

Rabani, M., Raychowdhury, R., Jovanovic, M., Rooney, M., Stumpo, D. J., Pauli, A., Hacohen, N., Schier, A. F., Blackshear, P. J., Friedman, N., Amit, I. & Regev, A. High-resolution sequencing and modeling identifies distinct dynamic RNA regulatory strategies. *Cell*. 159, 1698-1710, doi:10.1016/j.cell.2014.11.015 (2014). PMID:4272607.

Rajagopal, N., Srinivasan, S., Kooshesh, K., Guo, Y., Edwards, M.D., Banerjee, B., Syed, T., Emons, B.J.M., Gifford, D.K., and Sherwood, R.I. (2016). High-throughput mapping of regulatory DNA. *Nat. Biotechnol.* 34, 167–174.

Ramsauer, K., Farlik, M., Zupkovitz, G., Seiser, C., Kröger, A., Hauser, H., and Decker, T. (2007). Distinct modes of action applied by transcription factors STAT1 and IRF1 to initiate transcription of the IFN-gamma-inducible *gfp2* gene. *Proc. Natl. Acad. Sci. U. S. A.* 104, 2849–2854.

Ram, O., Goren, A., Amit, I., Shores, N., Yosef, N., Ernst, J., Kellis, M., Gymrek, M., Issner, R., Coyne, M., Durham, T., Zhang, X., Donaghey, J., Epstein, C. B., Regev, A. & Bernstein, B. E. Combinatorial patterning of chromatin regulators uncovered by genome-wide location analysis in human cells. *Cell*. 147, 1628-1639, doi:10.1016/j.cell.2011.09.057 (2011).

PMCID:3312319.

Ran, F.A., Cong, L., Yan, W.X., Scott, D.A., Gootenberg, J.S., Kriz, A.J., Zetsche, B., Shalem, O., Wu, X., Makarova, K.S., et al. (2015). In vivo genome editing using *Staphylococcus aureus* Cas9. *Nature* 520, 186–191.

Ron, D. & Walter, P. 2007, "Signal integration in the endoplasmic reticulum unfolded protein response", *Nature Reviews. Molecular Cell Biology*, vol. 8, no. 7, pp. 519-529.

Rosvall, M., and Bergstrom, C.T. (2008). Maps of random walks on complex networks reveal community structure. *Proc. Natl. Acad. Sci.* 105, 1118–1123.

Sack, L.M., Davoli, T., Xu, Q., Li, M.Z. & Elledge, S.J. 2016, "Sources of Error in Mammalian Genetic Screens", *G3 (Bethesda, Md.)*, vol. 6, no. 9, pp. 2781-2790.

Sackton, T.B., and Hartl, D.L. (2016). Perspective Genotypic Context and Epistasis in Individuals and Populations. *Cell* 166, 279–287.

Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nature biotechnology*. 33, 495-502, doi:10.1038/nbt.3192 (2015).

Sauvageau, M., Goff, L. A., Lodato, S., Bonev, B., Groff, A. F., Gerhardinger, C., Sanchez-Gomez, D. B., Hacisuleyman, E., Li, E., Spence, M., Liapis, S. C., Mallard, W., Morse, M., Swerdel, M. R., D'Ecclessis, M. F., Moore, J. C., Lai, V., Gong, G., Yancopoulos, G. D., Friendewey, D., Kellis, M., Hart, R. P., Valenzuela, D. M., Arlotta, P. & Rinn, J. L. Multiple knockout mouse models reveal lincRNAs are required for life and brain development. *eLife*. 2, e01749, doi:10.7554/eLife.01749 (2013). PMCID:3874104.

Sawyers, C. (2004). Targeted cancer therapy. *Nature* 432, 294–297.

Schölkopf, B., Platt, J.C., Shawe-Taylor, J., Smola, A.J. & Williamson, R.C. 2001, "Estimating the support of a high-dimensional distribution", *Neural computation*, vol. 13, no. 7, pp. 1443-1471.

Schwartz, S., Agarwala, S. D., Mumbach, M. R., Jovanovic, M., Mertins, P., Shishkin, A., Tabach, Y., Mikkelsen, T. S., Satija, R., Ruvkun, G., Carr, S. A., Lander, E. S., Fink, G. R. & Regev, A. High-resolution mapping reveals a conserved, widespread, dynamic mRNA methylation program in yeast meiosis. *Cell*. 155, 1409-1421, doi:10.1016/j.cell.2013.10.047 (2013). PMCID:3956118.

Schwartz, S., Bernstein, D. A., Mumbach, M. R., Jovanovic, M., Herbst, R. H., Leon-

Ricardo, B. X., Engreitz, J. M., Guttman, M., Satija, R., Lander, E. S., Fink, G. & Regev, A. Transcriptome-wide mapping reveals widespread dynamic-regulated pseudouridylation of ncRNA and mRNA. *Cell*. 159, 148-162, doi:10.1016/j.cell.2014.08.028 (2014). PMID:4180118.

Schwartz, S., Mumbach, M. R., Jovanovic, M., Wang, T., Maciag, K., Bushkin, G. G., Mertins, P., Ter-Ovanesyan, D., Habib, N., Cacchiarelli, D., Sanjana, N. E., Freinkman, E., Pacold, M. E., Satija, R., Mikkelsen, T. S., Hacohen, N., Zhang, F., Carr, S. A., Lander, E. S. & Regev, A. Perturbation of m6A writers reveals two distinct classes of mRNA methylation at internal and 5' sites. *Cell reports*. 8, 284-296, doi:10.1016/j.celrep.2014.05.048 (2014). PMID:4142486.

Shahni, R., Cale, C.M., Anderson, G., Osellame, L.D., Hambleton, S., Jacques, T.S., Wedatilake, Y., Taanman, J.-W., Chan, E., Qasim, W., et al. (2015). Signal transducer and activator of transcription 2 deficiency is a novel disorder of mitochondrial fission. *Brain* 138, 2834–2846.

Shakya, A., Callister, C., Goren, A., Yosef, N., Garg, N., Khoddami, V., Nix, D., Regev, A. & Tantin, D. Pluripotency transcription factor Oct4 mediates stepwise nucleosome demethylation and depletion. *Mol Cell Biol*. 35, 1014-1025, doi:10.1128/MCB.01105-14 (2015). PMID:4333097.

Shalek, A. K., Gaublomme, J. T., Wang, L., Yosef, N., Chevrier, N., Andersen, M. S., Robinson, J. T., Pochet, N., Neuberg, D., Gertner, R. S., Amit, I., Brown, J. R., Hacohen, N., Regev, A., Wu, C. J. & Park, H. Nanowire-mediated delivery enables functional interrogation of primary immune cells: application to the analysis of chronic lymphocytic leukemia. *Nano Lett*. 12, 6498-6504, doi:10.1021/nl3042917 (2012). PMID:3573729.

Shalek, A.K., Satija, R., Adiconis, X., Gertner, R.S., Gaublomme, J.T., Raychowdhury, R., Schwartz, S., Yosef, N., Malboeuf, C., Lu, D., et al. (2013). Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells. *Nature* 498, 236–240.

Shalek, A.K., Satija, R., Shuga, J., Trombetta, J.J., Gennert, D., Lu, D., Chen, P., Gertner, R.S., Gaublomme, J.T., Yosef, N., et al. (2014). Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. *Nature* 510, 363–369.

Shalem, O., Sanjana, N.E., Hartenian, E., Shi, X., Scott, D.A., Mikkelsen, T.S., Heckl, D., Ebert, B.L., Root, D.E., Doench, J.G., et al. (2014). Genome-scale CRISPR-Cas9 knockout

screening in human cells. *Science* 343, 84–87.

Shalem, O., Sanjana, N.E. & Zhang, F. 2015, "High-throughput functional genomics using CRISPR-Cas9", *Nature Reviews Genetics*, vol. 16, no. 5, pp. 299-311.

Shao, H., Burrage, L.C., Sinasac, D.S., Hill, A.E., Ernest, S.R., O'Brien, W., Courtland, H.-W., Jepsen, K.J., Kirby, A., Kulbokas, E.J., et al. (2008). Genetic architecture of complex traits: large phenotypic effects and pervasive epistasis. *Proc. Natl. Acad. Sci. U. S. A.* 105, 19910–19914.

Shendure, J., and Akey, J.M. (2015). The origins, determinants, and consequences of human mutations. *Science* (80-.). 349, 1478–1483.

Shendure, J., and Fields, S. (2016). Massively Parallel Genetics. *Genetics* 203, 617–619.

Shi, J., Wang, E., Milazzo, J.P., Wang, Z., Kinney, J.B. & Vakoc, C.R. 2015, "Discovery of cancer drug targets by CRISPR-Cas9 screening of protein domains", *Nature biotechnology*, vol. 33, no. 6, pp. 661-667.

Shoulders, M.D., Ryno, L.M., Genereux, J.C., Moresco, J.J., Tu, P.G., Wu, C., Yates, J.R., Su, A.I., Kelly, J.W. & Wiseman, R.L. 2013, "Stress-independent activation of XBP1s and/or ATF6 reveals three functionally diverse ER proteostasis environments", *Cell Reports*, vol. 3, no. 4, pp. 1279-1292.

Sidrauski, C., Tsai, J.C., Kampmann, M., Hearn, B.R., Vedantham, P., Jaishankar, P., Sokabe, M., Mendez, A.S., Newton, B.W., Tang, E.L., et al 2015, "Pharmacological dimerization and activation of the exchange factor eIF2B antagonizes the integrated stress response", *eLife*, vol. 4, pp. e07314.

Sisler, J.D., Morgan, M., Raje, V., Grande, R.C., Derecka, M., Meier, J., Cantwell, M., Szczepanek, K., Korzun, W.J., Lesnefsky, E.J., et al. (2015). The Signal Transducer and Activator of Transcription 1 (STAT1) Inhibits Mitochondrial Biogenesis in Liver and Fatty Acid Oxidation in Adipocytes. *PLoS One* 10, e0144444.

Smith, R.P., Taher, L., Patwardhan, R.P., Kim, M.J., Inoue, F., Shendure, J., Ovcharenko, I., and Ahituv, N. (2013). Massively parallel decoding of mammalian regulatory sequences supports a flexible organizational model. *Nat. Genet.* 45, 1021–1028.

Smith, Z. D., Chan, M. M., Humm, K. C., Karnik, R., Mekhoubad, S., Regev, A., Eggan, K. & Meissner, A. DNA methylation dynamics of the human preimplantation embryo. *Nature*. 511, 611-615, doi:10.1038/nature13581 (2014). PMID:4178976.

Smith, Z. D., Chan, M. M., Mikkelsen, T. S., Gu, H., Gnirke, A., Regev, A. & Meissner, A. A unique regulatory phase of DNA methylation in the early mammalian embryo. *Nature*. 484, 339-344, doi:10.1038/nature10960 (2012). PMID:3331945.

Smyth, R.P., Davenport, M.P. & Mak, J. 2012, "The origin of genetic diversity in HIV-1", *Virus Research*, vol. 169, no. 2, pp. 415-429.

Snijder, B., Sacher, R., Rämö, P., Damm, E., Liberali, P. & Pelkmans, L. 2009, "Population context determines cell-to-cell variability in endocytosis and virus infection", *Nature*, vol. 461, no. 7263, pp. 520-523.

Sokolov, A., Carlin, D.E., Paull, E.O., Baertsch, R., and Stuart, J.M. (2016). Pathway-Based Genomics Prediction using Generalized Elastic Net. *PLOS Comput. Biol.* 12, e1004790.

Stegle, O., Teichmann, S.A., and Marioni, J.C. (2015). Computational and analytical challenges in single-cell transcriptomics. *Nat. Rev. Genet.* 16, 133–145.

Sun, L., Goff, L. A., Trapnell, C., Alexander, R., Lo, K. A., Haciosuleyman, E., Sauvageau, M., Tazon-Vega, B., Kelley, D. R., Hendrickson, D. G., Yuan, B., Kellis, M., Lodish, H. F. & Rinn, J. L. Long noncoding RNAs regulate adipogenesis. *Proceedings of the National Academy of Sciences of the United States of America*. 110, 3387-3392, doi:10.1073/pnas.1222643110 (2013). PMID:3587215.

Tanenbaum, M.E., Gilbert, L.A., Qi, L.S., Weissman, J.S. & Vale, R.D. 2014, "A protein-tagging system for signal amplification in gene expression and fluorescence imaging", *Cell*, vol. 159, no. 3, pp. 635-646.

Tang, H., Klopfenstein, D., Pedersen, B., Flick, P., Sato, K., Ramirez, F., Yunes, J., and Mungall, C. (2015). GOATOOLS: Tools for Gene Ontology.

Theile, C.S., Witte, M.D., Blom, A.E., Kundrat, L., Ploegh, H.L., and Guimaraes, C.P. (2013). Site-specific N-terminal labeling of proteins using sortase-mediated reactions. *Nature protocols* 8, 1800-1807.

Thomason, L.C., Costantino, N. & Court, D.L. 2007, "E. coli genome manipulation by P1 transduction", *Current Protocols in Molecular Biology*, vol. Chapter 1, pp. Unit 1.17.

Thomason, L.C., Sawitzke, J.A., Li, X., Costantino, N. & Court, D.L. 2014, "Recombineering: genetic engineering in bacteria using homologous recombination", *Current Protocols in Molecular Biology*, vol. 106, pp. 39.

Tong, A.H.Y. (2004). Global Mapping of the Yeast Genetic Interaction Network. *Science*

(80-). 303, 808–813.

Trapnell, C. 2015, "Defining cell types and states with single-cell genomics", *Genome Research*, vol. 25, no. 10, pp. 1491-1498.

Trapnell, C., Cacchiarelli, D., Grimsby, J., Pokharel, P., Li, S., Morse, M., Lennon, N. J., Livak, K. J., Mikkelsen, T. S. & Rinn, J. L. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature biotechnology*. 32, 381-386, doi:10.1038/nbt.2859 (2014). PMID:4122333.

Trapnell, C., Hendrickson, D. G., Sauvageau, M., Goff, L., Rinn, J. L. & Pachter, L. Differential analysis of gene regulation at transcript resolution with RNA-seq. *Nature biotechnology*. 31, 46-53, doi:10.1038/nbt.2450 (2013). PMID:3869392.

Trombetta, J. J., Gennert, D., Lu, D., Satija, R., Shalek, A. K. & Regev, A. Preparation of Single-Cell RNA-Seq Libraries for Next Generation Sequencing. *Curr Protoc Mol Biol*. 107, 4 22 21-24 22 17, doi:10.1002/0471142727.mb0422s107 (2014). PMID:4338574.

Tsumura, A., Hayakawa, T., Kumaki, Y., Takebayashi, S., Sakaue, M., Matsuoka, C., Shimotohno, K., Ishikawa, F., Li, E., Ueda, H.R., et al. (2006). Maintenance of self-renewal ability of mouse embryonic stem cells in the absence of DNA methyltransferases Dnmt1, Dnmt3a and Dnmt3b. *Genes to Cells* 11, 805–814.

Tussiwand, R., Lee, W.-L., Murphy, T.L., Mashayekhi, M., KC, W., Albring, J.C., Satpathy, A.T., Rotondo, J.A., Edelson, B.T., Kretzer, N.M., et al. (2012). Compensatory dendritic cell development mediated by BATF-IRF interactions. *Nature* 490, 502–507.

Tyynismaa, H., Carroll, C.J., Raimundo, N., Ahola-Erkkilä, S., Wenz, T., Ruhanen, H., Guse, K., Hemminki, A., Peltola-Mjøsund, K.E., Tulkki, V., et al 2010, "Mitochondrial myopathy induces a starvation-like response", *Human Molecular Genetics*, vol. 19, no. 20, pp. 3948-3958.

Van Der Maaten, L. 2014, "Accelerating t-SNE using tree-based algorithms.", *Journal of machine learning research*, vol. 15, no. 1, pp. 3221-3245.

Visscher, P.M., Brown, M.A., McCarthy, M.I., and Yang, J. (2012). Five Years of GWAS Discovery. *Am. J. Hum. Genet.* 90, 7–24.

Walter, P. & Ron, D. 2011, "The unfolded protein response: from stress pathway to homeostatic regulation", *Science (New York, N.Y.)*, vol. 334, no. 6059, pp. 1081-1086.

Wang, L., Shalek, A. K., Lawrence, M., Ding, R., Gaublomme, J. T., Pochet, N.,

Stojanov, P., Sougnez, C., Shukla, S. A., Stevenson, K. E., Zhang, W., Wong, J., Sievers, Q. L., MacDonald, B. T., Vartanov, A. R., Goldstein, N. R., Neuberg, D., He, X., Lander, E., Hacohen, N., Regev, A., Getz, G., Brown, J. R., Park, H. & Wu, C. J. Somatic mutation as a mechanism of Wnt/beta-catenin pathway activation in CLL. *Blood*. 124, 1089-1098, doi:10.1182/blood-2014-01-552067 (2014). PMID:4133483.

Wang, T., Wei, J.J., Sabatini, D.M., and Lander, E.S. (2014). Genetic screens in human cells using the CRISPR-Cas9 system. *Science* 343, 80–84.

Wang, T., Birsoy, K., Hughes, N.W., Krupczak, K.M., Post, Y., Wei, J.J., Lander, E.S., and Sabatini, D.M. (2015). Identification and characterization of essential genes in the human genome. *Science* (80-.). 350, 1096–1101.

Wang, Y., Shen, J., Arenzana, N., Tirasophon, W., Kaufman, R.J. & Prywes, R. 2000, "Activation of ATF6 and an ATF6 DNA binding site by the endoplasmic reticulum stress response", *The Journal of Biological Chemistry*, vol. 275, no. 35, pp. 27013-27020.

Wei, L., Fan, M., Xu, L., Heinrich, K., Berry, M.W., Homayouni, R., and Pfeffer, L.M. (2008). Bioinformatic analysis reveals cRel as a regulator of a subset of interferon-stimulated genes. *J. Interferon Cytokine Res.* 28, 541–551.

Washietl, S., Kellis, M. & Garber, M. Evolutionary dynamics and tissue specificity of human long noncoding RNAs in six mammals. *Genome Res.* 24, 616-628, doi:10.1101/gr.165035.113 (2014). PMID:3975061.

Weinberger, E.D. (1991). Fourier and Taylor series on fitness landscapes. *Biol. Cybern.* 65, 321–330.

Wong, A.S.L., Choi, G.C.G., Cui, C.H., Pregernig, G., Milani, P., Adam, M., Perli, S.D., Kazer, S.W., Gaillard, A., Hermann, M., et al. (2016). Multiplexed barcoded CRISPR-Cas9 screening enabled by CombiGEM. *Proc. Natl. Acad. Sci.* 113, 2544–2549.

Wu, C., Yosef, N., Thalhamer, T., Zhu, C., Xiao, S., Kishi, Y., Regev, A. & Kuchroo, V. K. Induction of pathogenic TH17 cells by inducible salt-sensing kinase SGK1. *Nature*. 496, 513-517, doi:10.1038/nature11984 (2013). PMID:3637879.

Yosef, N., and Regev, A. (2016). Writ large: Genomic dissection of the effect of cellular environment on immune response. *Science* (80-.). 354, 64–68.

Yosef, N., Shalek, A. K., Gaublomme, J. T., Jin, H., Lee, Y., Awasthi, A., Wu, C., Karwacz, K., Xiao, S., Jorgolli, M., Gennert, D., Satija, R., Shakya, A., Lu, D. Y., Trombetta, J.

J., Pillai, M. R., Ratcliffe, P. J., Coleman, M. L., Bix, M., Tantin, D., Park, H., Kuchroo, V. K. & Regev, A. Dynamic regulatory network controlling TH17 cell differentiation. *Nature*. 496, 461-468, doi:10.1038/nature11981 (2013). PMID:3637864.

Yu, C., Mannan, A.M., Yvone, G.M., Ross, K.N., Zhang, Y.-L., Marton, M.A., Taylor, B.R., Crenshaw, A., Gould, J.Z., Tamayo, P., et al. (2016). High-throughput identification of genotype-specific cancer vulnerabilities in mixtures of barcoded tumor cell lines. *Nat. Biotechnol.* 34, 419–423.

Zalatan, J.G., Lee, M.E., Almeida, R., Gilbert, L.A., Whitehead, E.H., La Russa, M., Tsai, J.C., Weissman, J.S., Dueber, J.E., Qi, L.S., et al. (2015). Engineering Complex Synthetic Transcriptional Programs with CRISPR RNA Scaffolds. *Cell* 160, 339–350.

Zetsche, B., Gootenberg, J.S., Abudayyeh, O.O., Slaymaker, I.M., Makarova, K.S., Essletzbichler, P., Volz, S.E., Joung, J., van der Oost, J., Regev, A., et al. (2015). Cpf1 Is a Single RNA-Guided Endonuclease of a Class 2 CRISPR-Cas System. *Cell* 163, 759–771.

Zetsche, B., Heidenreich, M., Mohanraju, P., Fedorova, I., Kneppers, J., DeGennaro, E.M., Winblad, N., Choudhury, S.R., Abudayyeh, O.O., Gootenberg, J.S., et al 2016, "Multiplex gene editing by CRISPR-Cpf1 through autonomous processing of a single crRNA array", *bioRxiv*, .

Zhang, X., Liu, Q., Luo, C., Deng, Y., Cui, K. & Shi, D. 2014, "Identification and characterization of buffalo 7SK and U6 pol III promoters and application for expression of short hairpin RNAs", *International Journal of Molecular Sciences*, vol. 15, no. 2, pp. 2596-2607.

Zheng, G.X.Y., Terry, J.M., Belgrader, P., Ryvkin, P., Bent, Z.W., Wilson, R., Ziraldo, S.B., Wheeler, T.D., McDermott, G.P., Zhu, J., et al 2016, "Massively parallel digital transcriptional profiling of single cells", *bioRxiv*, .

Zuk, O., Hechter, E., Sunyaev, S.R., and Lander, E.S. (2012). The mystery of missing heritability: Genetic interactions create phantom heritability. *Proc. Natl. Acad. Sci. U. S. A.* 109, 1193–1198.

[001072] The invention is further described by the following numbered paragraphs:

1. A method of reconstructing a cellular network or circuit, comprising introducing at least 1, 2, 3, 4 or more single-order or combinatorial perturbations to each cell in a population of cells;

measuring genomic, genetic and/or phenotypic differences of each cell and coupling combinatorial perturbations with measured differences to infer intercellular and/or intracellular networks or circuits.

2. The method of paragraph 1, wherein the single-order or combinatorial perturbations comprise 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98 or 99 or massively parallel perturbations.

3. The method of paragraph 1 or 2 wherein the perturbation(s) comprise one or more genetic perturbation.

4. The method of paragraph 1, 2 or 3 wherein the perturbation(s) comprise one or more epigenetic or epigenomic perturbation.

5. The method of any of the preceding paragraphs wherein at least one perturbation is introduced with RNAi- or a CRISPR-Cas system.

6. The method of any of the preceding paragraphs wherein at least one perturbation is introduced via a chemical agent, an intracellular spatial relationship between two or more cells, an increase or decrease of temperature, addition or subtraction of energy, electromagnetic energy, or ultrasound.

7. The method of any one of the preceding paragraphs wherein the cell comprises a cell in a model non-human organism, a model non-human mammal that expresses a Cas protein, a mouse that expresses a Cas protein, a cell *in vivo* or a cell *ex vivo* or a cell *in vitro*.

8. The method of any one of the preceding paragraphs wherein measuring or measured differences comprises measuring or measured differences of DNA, RNA, protein or post translational modification; or measuring or measured differences of protein or post translational modification correlated to RNA and/or DNA level(s).

9. The method of any preceding paragraph wherein the method includes sequencing, and prior to sequencing:

(a) perturbing and isolating a single cell with at least one labeling ligand specific for binding at one or more target RNA transcripts, or isolating a single cell with at least one labeling ligand specific for binding at one or more target RNA transcripts and perturbing the cell; and/or

(b) lysing the cell under conditions wherein the labeling ligand binds to the target RNA transcript(s).

10. The method of paragraph 9 wherein step (a) comprises perturbing and isolating a single cell with at least one labeling ligand specific for binding at one or more target RNA transcripts.

11. The method of paragraph 9 wherein step (a) comprises isolating a single cell with at least one labeling ligand specific for binding at one or more target RNA transcripts and perturbing the cell.

12. The method of any preceding paragraph wherein perturbing or perturbation(s) comprise(s) genetic perturbing.

13. The method of any preceding paragraph wherein perturbing or perturbation(s) comprise(s) single-order perturbations.

14. The method of any preceding paragraph wherein perturbing or perturbation(s) comprise(s) combinatorial perturbations.

15. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises gene knock-down, gene knock-out, gene activation, gene insertion, or regulatory element deletion.

16. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises genome-wide perturbation.

17. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises performing CRISPR-Cas-based perturbation.

18. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises performing pooled single or combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

19. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises performing pooled combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

20. The method of any one of paragraphs 18 or 19, wherein each sgRNA comprises a unique molecular identifier.

21. The method of any one of paragraphs 18 or 19, wherein each sgRNA is co-delivered with a reporter mRNA.

22. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises subjecting the cell to an increase or decrease in temperature.

23. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises subjecting the cell to a chemical agent.
24. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises subjecting the cell to a chemical agent and/or temperature increase or decrease across a gradient.
25. The method of any preceding paragraph wherein the cell is in a microfluidic system.
26. The method of any preceding paragraph wherein the cell is in a droplet.
27. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises transforming or transducing the cell or a population that includes and from which the cell is isolated with one or more genomic sequence-perturbation constructs that perturbs a genomic sequence in the cell.
28. The method of any preceding paragraph wherein perturbing or perturbation(s) comprises multiplex transformation with a plurality of genomic sequence-perturbation constructs.
29. The method of any one of the preceding paragraphs, wherein the labeling ligand comprises a unique perturbation identifier (UPI) sequence attached to a perturbation-sequence-capture sequence, and sequencing includes isolating via microbeads comprising a perturbation-sequence-capture-binding-sequence having specific binding affinity for the perturbation-sequence-capture sequence attached to the UPI sequence.
30. The method of paragraph 29, wherein the UPI sequence is attached to a universal ligation handle sequence, whereby a unique source identifier USI may be generated by split-pool ligation.
31. The method of any preceding paragraph wherein the labeling ligand comprises an oligonucleotide label comprising a regulatory sequence configured for amplification by T7 polymerase.
32. The method of any preceding paragraph wherein the labeling ligands comprises oligonucleotide sequences configured to hybridize to a transcript specific region.
33. The method of any preceding paragraph wherein the labeling ligands comprises an oligonucleotide label.
34. The method of paragraph 33, wherein the oligonucleotide label further comprises a photocleavable linker.
35. The method of paragraph 33 or 34 wherein the oligonucleotide label further comprises a restriction enzyme site between the labeling ligand and unique constituent identifier (UCI).

36. The method of any preceding paragraph, wherein isolating comprises forming discrete unique-identifier-transfer compositions, each comprising the cell and a transfer particle, wherein:

- (a) an oligonucleotide label further comprises a capture sequence, and unique constituent identifier (UCI) and capture sequence are together releasably attached to the labeling ligand;

- (b) the labelling ligand is bound to the target cellular constituent; and,

- (c) a transfer particle comprises:

- (i) a capture-binding-sequence having specific binding affinity for the capture sequence attached to the UCI, and,

- (ii) a unique source identifier (USI) sequence that is unique to each transfer particle.

37. The method of paragraph 36 wherein, the USI comprises 4-15 nucleotides.

38. The method of paragraph 36 or 37, further comprising releasing the UCI from the labeled ligand, under conditions within the unique-identifier-transfer composition so that the released capture sequence binds to the capture-binding-sequence on the transfer particle, thereby transferring the UCI to the transfer particle.

39. The method of any preceding paragraph wherein the ligation handle comprises a restriction site for producing an overhang complementary with a first index sequence overhang, and wherein the method further comprises digestion with a restriction enzyme.

40. The method of method of any preceding paragraph wherein wherein the ligation handle comprises a nucleotide sequence complementary with a ligation primer sequence and wherein the overhang complementary with a first index sequence overhang is produced by hybridization of the ligation primer to the ligation handle.

41. The method of any one of the preceding paragraphs further comprising quantitating relative amount of UCI sequence associated with a first cell to the amount of the same UCI sequence associated with a second cell, whereby the relative differences of a cellular constituent between cell(s) are determined.

42. The method of any one of the preceding paragraphs wherein the labeling ligand comprises an antibody or an antibody fragment.

43. The method of any one of the preceding paragraphs wherein the antibody fragment is a nanobody, Fab, Fab', (Fab')₂, Fv, ScFv, diabody, triabody, tetrabody, Bis-scFv, minibody, Fab₂, or Fab₃ fragment.

44. The method of any preceding paragraph wherein the labeling ligand comprises an aptamer.

45. The method of any preceding paragraph wherein the labeling ligand is a nucleotide sequence complementary to a target sequence.

46. The method of any method of any preceding paragraph wherein the cell or a population that includes wherein the cell(s) are a member of a cell population, and the method further comprises transforming or transducing the cell population with one or more genomic sequence-perturbation constructs that perturb a genomic sequence in the cells, wherein each distinct genomic sequence-perturbation construct comprises a unique-perturbation-identified (UPI) sequence unique to that construct.

47. The method of paragraph 46, wherein the genomic sequence-perturbation construct comprises a sequence encoding a guide RNA sequence of a CRISPR-Cas targeting system.

48. The method of paragraph 46 or 47, further comprising multiplex transformation of the population of cells with a plurality of genomic sequence-perturbation constructs.

49. The method of any one of paragraphs 46, 47 or 48, wherein the UPI sequence is attached to a perturbation-sequence-capture sequence, and the transfer particle comprise a perturbation-sequence-capture-binding-sequence having specific binding affinity for the perturbation-sequence-capture sequence attached to the UPI sequence.

50. The method of any one of paragraphs 46 or 47, wherein the UPI sequence is attached to a universal ligation handle sequence, whereby a USI is generated by split-pool ligation.

51. A method of determining any combination of protein detection, RNA detection, open chromatin detection, protein-protein interactions, protein-RNA interactions, or protein-DNA interactions comprising any of the preceding methods.

52. A method according to any of the preceding paragraphs including whole transcriptome amplification.

53. The method according to any one of the preceding paragraphs including comparing an RNA profile of the perturbed cell with any mutations in the cell to also correlate phenotypic or transcriptome profile and genotypic profile.

54. A method for determining genetic interactions by causing a set of P genetic perturbations in cells,

wherein the method comprises:

- (a) Determining, based upon random sampling, a subset of π genetic perturbations from the set of P genetic perturbations;

- (b) Performing said subset of π genetic perturbations in a population of cells;
- (c) Performing single-cell molecular profiling of the population of genetically perturbed cells of step (b);
- (d) Inferring, from the results of step (c) and based upon the random sampling of step (a), single-cell molecular profiles for the set of P genetic perturbations in cells.

55. The method according to paragraph 54,

wherein the method further comprises:

- (e) From the results of step (d), determining genetic interactions.

56. The method according to paragraph 54,

wherein the method further comprises:

- (f) Confirming genetic interactions determined at step (e) with additional genetic manipulations.

57. The method according to any one of paragraphs 54-56,

wherein said set of P genetic perturbations or said subset of π genetic perturbations comprises single-order genetic perturbations.

58. The method according to any one of paragraphs 54-57,

wherein said set of P genetic perturbations or said subset of π genetic perturbations comprises combinatorial genetic perturbations.

59. The method according to any one of paragraphs 54-58,

wherein said genetic perturbation comprises gene knock-down, gene knock-out, gene activation, gene insertion, or regulatory element deletion.

60. The method according to any one of paragraphs 54-59,

wherein said set of P genetic perturbations or said subset of π genetic perturbations comprises genome-wide perturbations.

61. The method according to any one of paragraphs 54-60,

wherein said set of P genetic perturbations or said subset of π genetic perturbations comprises k-order combinations of single genetic perturbations, wherein k is an integer ranging from 2 to 15, and

wherein step (e) comprises determining k-order genetic interactions.

62. The method according to any one of paragraphs 54-61,

wherein said set of P genetic perturbations comprises combinatorial genetic perturbations, such as k-order combinations of single-order genetic perturbations, wherein k is an integer ranging from 2 to 15, and

wherein step (e) comprises determining j-order genetic interactions, with $j < k$.

63. The method according to any one of paragraphs 54-62, wherein step (b) comprises performing RNAi- or CRISPR-Cas-based perturbation.

64. The method according to any one of paragraphs 54-63, wherein step (b) comprises array-format or pool-format perturbation.

65. The method according to any one of paragraphs 54-64, wherein step (b) comprises pooled single or combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

66. The method according to any one of paragraphs 54-65, wherein step (b) comprises pooled combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

67. The method of paragraph 65 or 66, wherein each sgRNA comprises a unique molecular identifier.

68. The method of paragraph 65, 66 or 67, wherein each sgRNA is co-delivered with a reporter mRNA.

69. The method according to any one of paragraphs 54-68, wherein random sampling comprises matrix completion, tensor completion, compressed sensing, or kernel learning.

70. The method according to any one of paragraphs 54-69, wherein random sampling comprises matrix completion, tensor completion, or compressed sensing, and wherein π is of the order of $\log P$.

71. The method of any one of the preceding paragraphs, wherein the cell comprises a eukaryotic cell.

72. The method of paragraph 71, wherein the eukaryotic cell comprises a mammalian cell.

73. The method of paragraph 72, wherein the mammalian cell comprises a human cell.

74. The method of any one of the preceding paragraphs, wherein the cell is from a population comprising 10^2 to 10^8 cells and/or DNA or RNA or protein or post translational modification measurements or variables per cell comprise 50 or more.

75. A method for screening a collection of compounds or agents comprising performing a method as paragraphed in any of the preceding paragraphs wherein perturbing comprises exposing the cell to each compound or agent.

76. A therapeutic identified by the method of paragraph 75.

77. The method of any one of the preceding paragraphs wherein the perturbing is across a library of cells to thereby obtain RNA level and optionally protein level cell-to-cell circuit data at genomic or transcript or expression level.

78. The method of paragraph 77, wherein the library of cells comprises or is from a tissue sample.

79. The method of paragraph 78, wherein the tissue sample comprises or is from a biopsy from a mammalian subject.

80. The method of paragraph 79, wherein the mammalian subject comprises a human subject.

81. The method of paragraph 80, wherein the biopsy is from a tumor.

82. The method of any one of paragraphs 77 to 81, further comprising reconstructing cell-to-cell circuits.

83. The method of any preceding paragraph wherein open chromatin detection includes fragmenting chromatin inside isolated intact nuclei from a cell, adding universal primers at cutting sites, and uniquely tagging DNA that originated from the cell.

[001073] The invention is further described by the following numbered statement:

[001074] A method of reconstructing a cellular network or circuit, comprising:

(a) introducing at least 1, 2, 3, 4 or more single-order or combinatorial perturbations to a plurality of cells in a population of cells, wherein each cell in the plurality of the cells receives at least 1 perturbation; and

(b) measuring comprising:

(i) detecting genomic, genetic, proteomic, epigenetic and/or phenotypic differences in single cells compared to one or more cells that did not receive any perturbation, and

(ii) detecting the perturbation(s) in single cells.

* * *

[001075] Having thus described in detail preferred embodiments of the present invention, it is to be understood that the invention defined by the above paragraphs is not to be limited to particular details set forth in the above description as many apparent variations thereof are possible without departing from the spirit or scope of the present invention.

WHAT IS CLAIMED IS:

1. A method of reconstructing a cellular network or circuit, comprising
 - (a) introducing at least 1, 2, 3, 4 or more single-order or combinatorial perturbations to a plurality of cells in a population of cells, wherein each cell in the plurality of the cells receives at least 1 perturbation;
 - (b) measuring comprising:
 - (i) detecting genomic, genetic, proteomic, epigenetic and/or phenotypic differences in single cells compared to one or more cells that did not receive any perturbation, and
 - (ii) detecting the perturbation(s) in single cells; and
 - (c) determining measured differences relevant to the perturbations by applying a model accounting for co-variates to the measured differences, whereby intercellular and/or intracellular networks or circuits are inferred.
2. The method of claim 1, wherein the measuring in single cells comprises single cell sequencing.
3. The method of claim 2, wherein single cell sequencing comprises cell barcodes, whereby the cell-of-origin of each RNA is recorded..
4. The method of claim 2, wherein the single cell sequencing comprises unique molecular identifiers (UMI), whereby the capture rate of the measured signals, such as transcript copy number or probe binding events, in a single cell is determined.
5. The method of any of claims 1 to 4, wherein the model comprises accounting for the capture rate of measured signals, whether the perturbation actually perturbed the cell (phenotypic impact), the presence of subpopulations of either different cells or cell states, and/or analysis of matched cells without any perturbation.
6. The method of any of claims 1 to 5, wherein the single-order or combinatorial perturbations comprise 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99 or 100 perturbations.

7. The method of any of claims 1 to 6, wherein the perturbation(s) target genes in a pathway or intracellular network.
8. The method of any of claims 1 to 7, wherein the measuring comprises detecting the transcriptome of each of the single cells.
9. The method of any of claims 1 to 8, wherein the perturbation(s) comprise one or more genetic perturbation.
10. The method of any of claims 1 to 9, wherein the perturbation(s) comprise one or more epigenetic or epigenomic perturbation.
11. The method of any of the preceding claims wherein at least one perturbation is introduced with RNAi- or a CRISPR-Cas system.
12. The method of any of the preceding claims wherein at least one perturbation is introduced via a chemical agent, biological agent, an intracellular spatial relationship between two or more cells, an increase or decrease of temperature, addition or subtraction of energy, electromagnetic energy, or ultrasound.
13. The method of any one of the preceding claims wherein the cell comprises a cell in a model non-human organism, a model non-human mammal that expresses a Cas protein, a mouse that expresses a Cas protein, a cell *in vivo* or a cell *ex vivo* or a cell *in vitro*.
14. The method of any one of the preceding claims wherein measuring or measured differences comprises measuring or measured differences of DNA, RNA, protein or post translational modification; or measuring or measured differences of protein or post translational modification correlated to RNA and/or DNA level(s).
15. The method of any preceding claim wherein perturbing or perturbation(s) comprise(s) genetic perturbing.
16. The method of any preceding claim wherein perturbing or perturbation(s) comprise(s) single-order perturbations.
17. The method of any preceding claim wherein perturbing or perturbation(s) comprise(s) combinatorial perturbations.
18. The method of any preceding claim wherein perturbing or perturbation(s) comprises gene knock-down, gene knock-out, gene activation, gene insertion, or regulatory element deletion.

19. The method of any preceding claim wherein perturbing or perturbation(s) comprises genome-wide perturbation.

20. The method of any preceding claim wherein perturbing or perturbation(s) comprises performing CRISPR-Cas-based perturbation.

21. The method of any preceding claim wherein perturbing or perturbation(s) comprises performing pooled single or combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

22. The method of any preceding claim wherein perturbing or perturbation(s) comprises performing pooled combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

23. The method of any of claims 20 to 22, wherein each sgRNA is associated with a unique perturbation barcode.

24. The method of claims 23, wherein each sgRNA is co-delivered with a reporter mRNA comprising the unique perturbation barcode.

25. The method of any preceding claim wherein perturbing or perturbation(s) comprises subjecting the cell to an increase or decrease in temperature.

26. The method of any preceding claim wherein perturbing or perturbation(s) comprises subjecting the cell to a chemical agent.

27. The method of any preceding claim wherein perturbing or perturbation(s) comprises subjecting the cell to a biological agent.

28. The method of claim 27, wherein the biological agent is a toll like receptor agonist or cytokine.

29. The method of any preceding claim wherein perturbing or perturbation(s) comprises subjecting the cell to a chemical agent, biological agent and/or temperature increase or decrease across a gradient.

30. The method of any preceding claim wherein the cell is in a microfluidic system.

31. The method of any preceding claim wherein the cell is in a droplet.

32. The method of any preceding claim wherein perturbing or perturbation(s) comprises transforming or transducing the cell or a population that includes and from which the cell is isolated with one or more genomic sequence-perturbation constructs that perturbs a genomic sequence in the cell.

33. The method of any preceding claim wherein perturbing or perturbation(s) comprises multiplex transformation with a plurality of genomic sequence-perturbation constructs.

34. The method of any one of the preceding claims, wherein proteins or transcripts expressed in single cells are determined in response to a perturbation, wherein the proteins or transcripts are detected in the single cells by binding of more than one labeling ligand comprising an oligonucleotide tag, wherein the oligonucleotide tag comprises a unique constituent identifier (UCI) specific for a target protein or transcript.

35. The method of claim 34, wherein, single cells are fixed in discrete particles.

36. The method of claim 34, wherein the oligonucleotide tag and sgRNA perturbation barcode comprises a universal ligation handle sequence, whereby a unique cell barcode may be generated by split-pool ligation.

37. The method of any of claims 34 to 36, wherein the labeling ligand comprises an oligonucleotide label comprising a regulatory sequence configured for amplification by T7 polymerase.

38. The method of any of claims 34 to 36, wherein the labeling ligands comprise oligonucleotide sequences configured to hybridize to a transcript specific region.

39. The method of any of claims 34 to 38, wherein the oligonucleotide label further comprises a photocleavable linker.

40. The method of any of claims 34 to 38, wherein the oligonucleotide label further comprises a restriction enzyme site between the labeling ligand and unique constituent identifier (UCI).

41. The method of any of claims 36 to 38, wherein the ligation handle comprises a restriction site for producing an overhang complementary with a first index sequence overhang, and wherein the method further comprises digestion with a restriction enzyme.

42. The method of method of any of claims 36 to 38, wherein the ligation handle comprises a nucleotide sequence complementary with a ligation primer sequence and wherein the overhang complementary with a first index sequence overhang is produced by hybridization of the ligation primer to the ligation handle.

43. The method of any one of claims 34 to 42, further comprising quantitating relative amount of UCI sequence associated with a first cell to the amount of the same UCI sequence

associated with a second cell, whereby the relative differences of a cellular constituent between cell(s) are determined.

44. The method of any one of claims 34 to 43, wherein the labeling ligand comprises an antibody or an antibody fragment.

45. The method of claim 44, wherein the antibody fragment is a nanobody, Fab, Fab', (Fab')₂, Fv, ScFv, diabody, triabody, tetrabody, Bis-scFv, minibody, Fab2, or Fab3 fragment.

46. The method of any of claims 34 to 43, wherein the labeling ligand comprises an aptamer.

47. The method of any of claims 34 to 43, wherein the labeling ligand is a nucleotide sequence complementary to a target sequence.

48. A method according to any of the preceding claims, wherein single cell sequencing comprises whole transcriptome amplification.

49. The method according to any one of the preceding claims including comparing an RNA profile of the perturbed cell with any mutations in the cell to also correlate phenotypic or transcriptome profile and genotypic profile.

50. The method according to any one of the preceding claims further comprising determining genetic interactions by causing a set of P genetic perturbations in single cells of the population of cells,

wherein the method comprises:

(a) determining, based upon random sampling, a subset of π genetic perturbations from the set of P genetic perturbations;

(b) performing said subset of π genetic perturbations in a population of cells;

(c) performing single-cell molecular profiling of the population of genetically perturbed cells of step (b);

(d) inferring, from the results of step (c) and based upon the random sampling of step (a), single-cell molecular profiles for the set of P genetic perturbations in cells.

51. The method according to claim 50,

wherein the method further comprises:

(e) from the results of step (d), determining genetic interactions.

52. The method according to claim 50,

wherein the method further comprises:

(f) confirming genetic interactions determined at step (e) with additional genetic manipulations.

53. The method according to any one of claims 50-52, wherein said set of P genetic perturbations or said subset of π genetic perturbations comprises single-order genetic perturbations.

54. The method according to any one of claims 50-53, wherein said set of P genetic perturbations or said subset of π genetic perturbations comprises combinatorial genetic perturbations.

55. The method according to any one of claims 50-54, wherein said genetic perturbation comprises gene knock-down, gene knock-out, gene activation, gene insertion, or regulatory element deletion.

56. The method according to any one of claims 50-55, wherein said set of P genetic perturbations or said subset of π genetic perturbations comprises genome-wide perturbations.

57. The method according to any one of claims 50-56, wherein said set of P genetic perturbations or said subset of π genetic perturbations comprises k-order combinations of single genetic perturbations, wherein k is an integer ranging from 2 to 15, and wherein step (e) comprises determining k-order genetic interactions.

58. The method according to any one of claims 50-57, wherein said set of P genetic perturbations comprises combinatorial genetic perturbations, such as k-order combinations of single-order genetic perturbations, wherein k is an integer ranging from 2 to 15, and wherein step (e) comprises determining j-order genetic interactions, with $j < k$.

59. The method according to any one of claims 50-58, wherein step (b) comprises performing RNAi- or CRIPSR-Cas-based perturbation.

60. The method according to any one of claims 50-59, wherein step (b) comprises array-format or pool-format perturbation.

61. The method according to any one of claims 50-60, wherein step (b) comprises pooled single or combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.

62. The method according to any one of claims 50-61, wherein step (b) comprises pooled combinatorial CRISPR-Cas-based perturbation with a genome-wide library of sgRNAs.
63. The method according to any one of claims 50-62, wherein random sampling comprises matrix completion, tensor completion, compressed sensing, or kernel learning.
64. The method according to any one of claims 50-63, wherein random sampling comprises matrix completion, tensor completion, or compressed sensing, and wherein π is of the order of $\log P$.
65. The method of any one of the preceding claims, wherein the cell comprises a eukaryotic cell.
66. The method of claim 65, wherein the eukaryotic cell comprises a mammalian cell.
67. The method of claim 66, wherein the mammalian cell comprises a human cell.
68. The method of any one of the preceding claims, wherein the cell is from a population comprising 10^2 to 10^8 cells and DNA or RNA or protein or post translational modification measurements or variables per cell comprise 50 or more.
69. The method of any one of the preceding claims, wherein perturbation of the population of cells is performed *in vivo*.
70. The method of any one of the preceding claims, wherein perturbation of the population of cells is performed *ex vivo* and the population of cells is adoptively transferred to a subject.
71. The method of claims 69 or 70, wherein the population of cells comprise tumor cells.
72. The method of any of claims 69 to 71, wherein a lineage barcode is associated with single cells, whereby the lineage or clonality of single cells may be determined.
73. The method of any one of the preceding claims wherein the perturbing is across a library of cells to thereby obtain RNA level and optionally protein level, whereby cell-to-cell circuit data at genomic or transcript or expression level is determined.
74. The method of claim 73, wherein the library of cells comprises or is from a tissue sample.

75. The method of claim 74, wherein the tissue sample comprises or is from a biopsy from a mammalian subject.

76. The method of claim 75, wherein the mammalian subject comprises a human subject.

77. The method of claim 76, wherein the biopsy is from a tumor.

78. The method of any one of claims 73 to 77, further comprising reconstructing cell-to-cell circuits.

79. The method of any preceding claims, wherein open chromatin is measured and comprises fragmenting chromatin inside isolated intact nuclei from a cell, adding universal primers at cutting sites, and uniquely tagging DNA that originated from the cell.

80. The method of any preceding claims, wherein measuring protein and RNA levels comprises CyTOF.

81. A method of determining any combination of protein detection, RNA detection, open chromatin detection, protein-protein interactions, protein-RNA interactions, or protein-DNA interactions comprising any of the preceding methods.

82. A method for screening compounds or agents capable of modifying a cellular network or circuit comprising performing a method as claimed in any of the preceding claims, wherein perturbing further comprises exposing the cell to each compound or agent.

83. A therapeutic identified by the method of claim 82.

Single-cell Clarity Drop-seq for NGS multiplex protein quantification

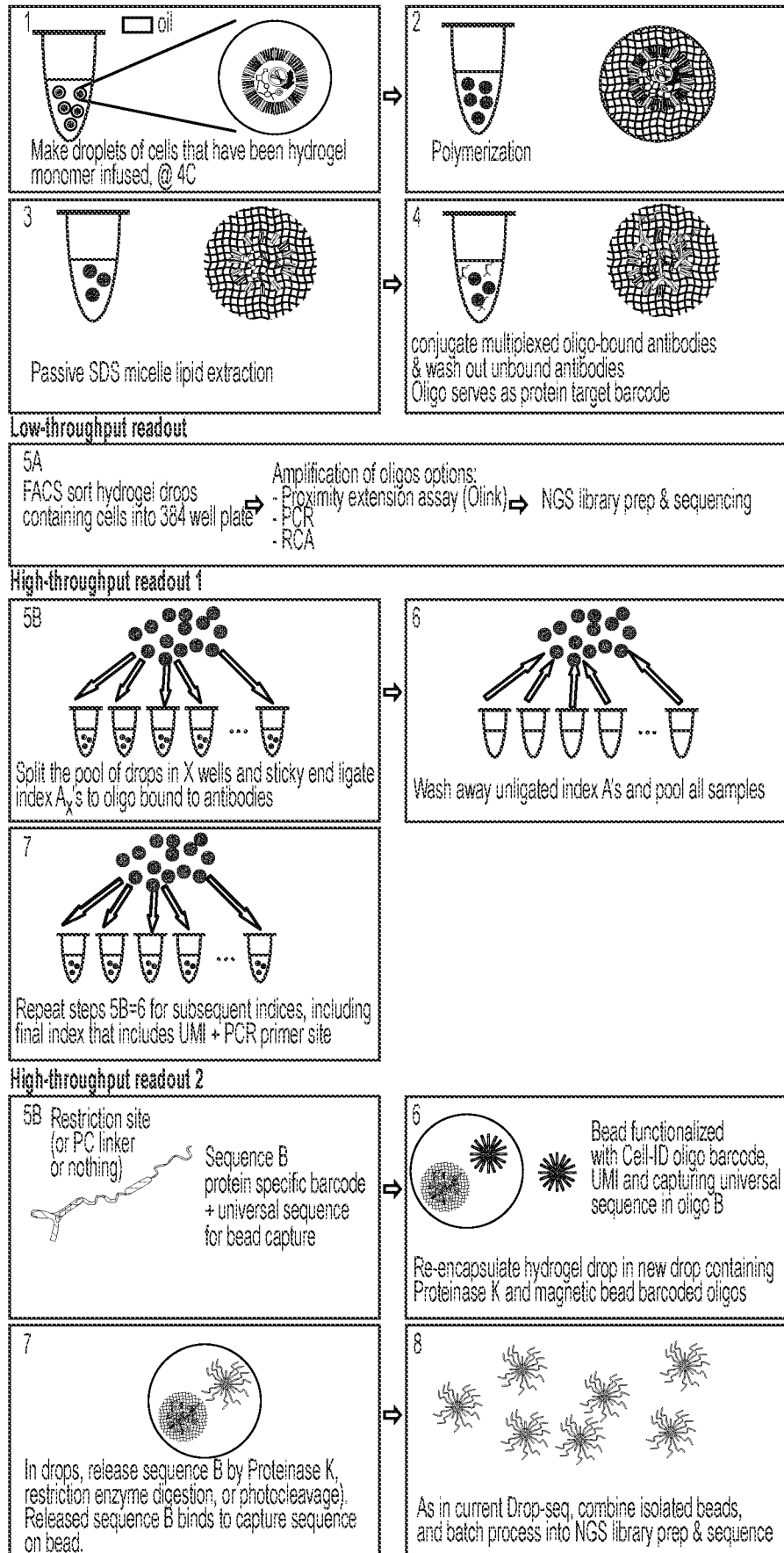


FIG. 1

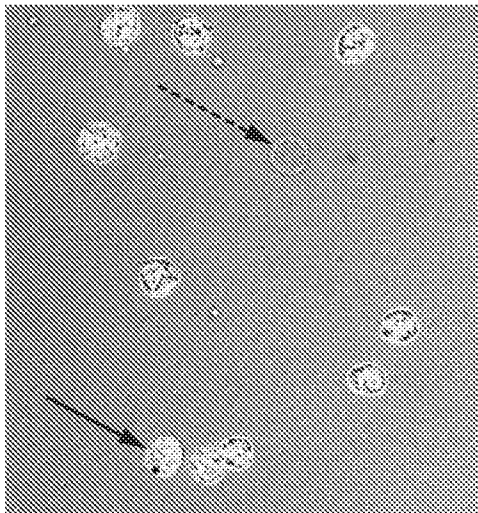


FIG. 2A

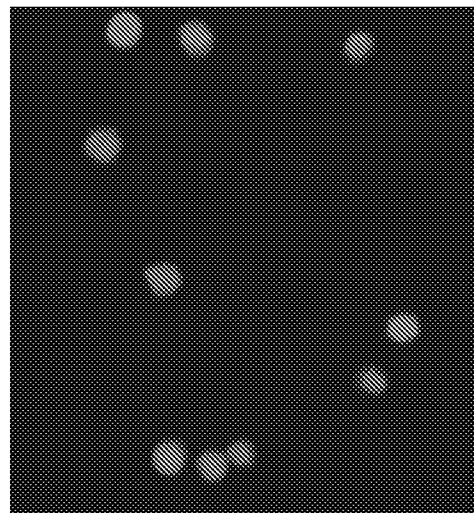


FIG. 2B

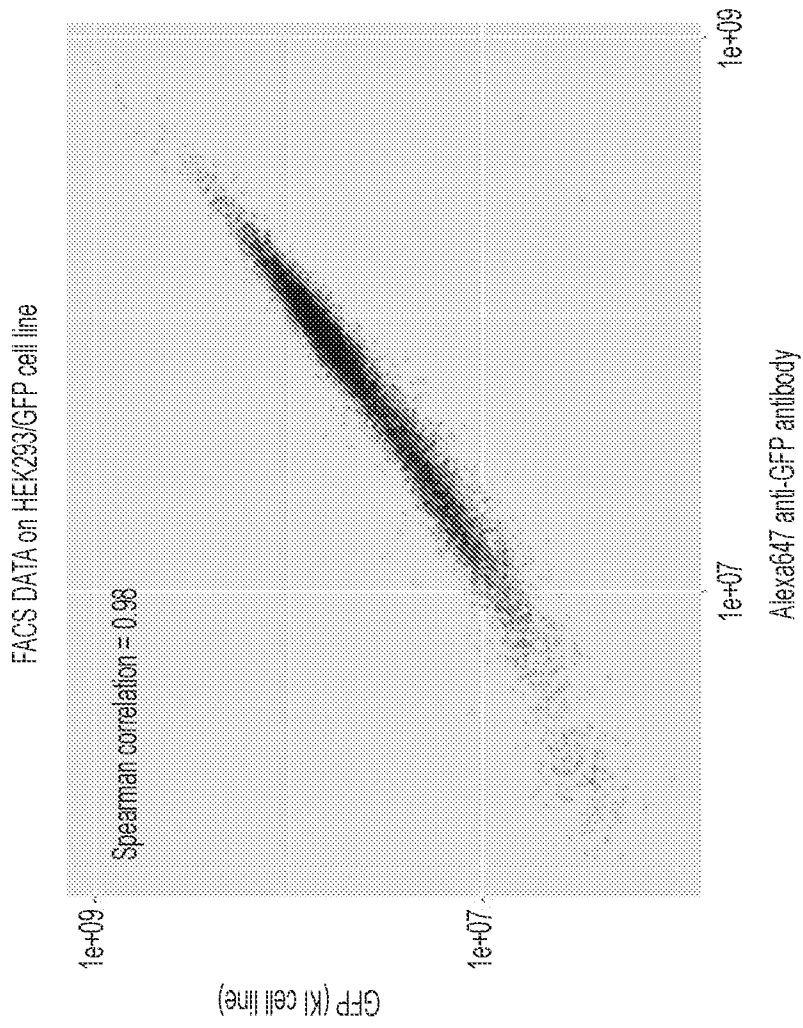


FIG. 2C

FIG. 3

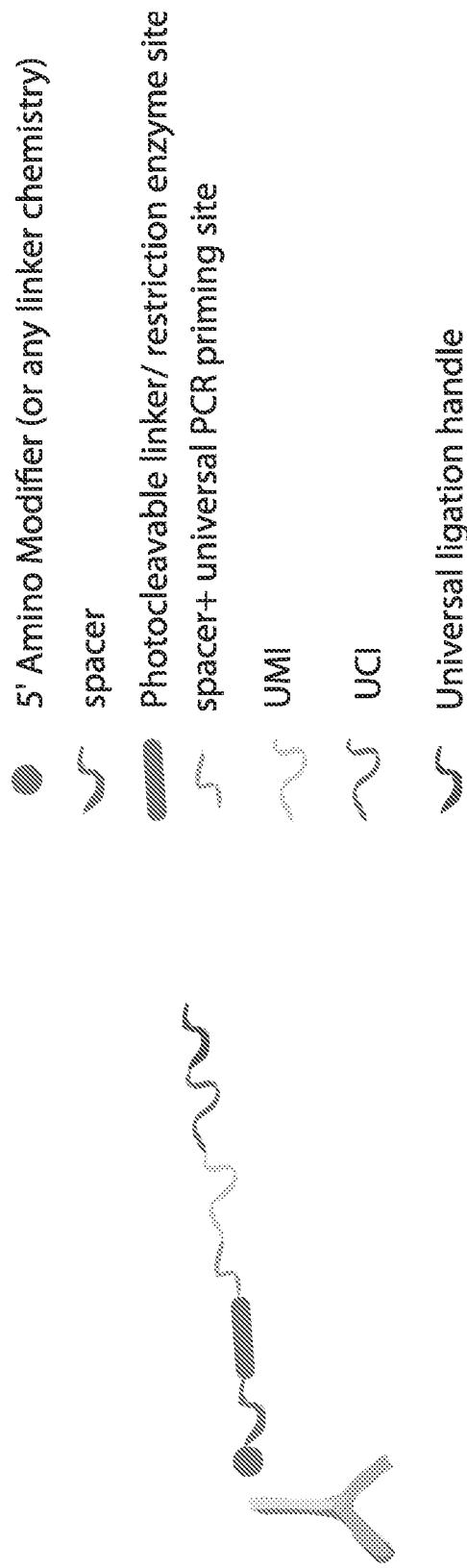


FIG. 4

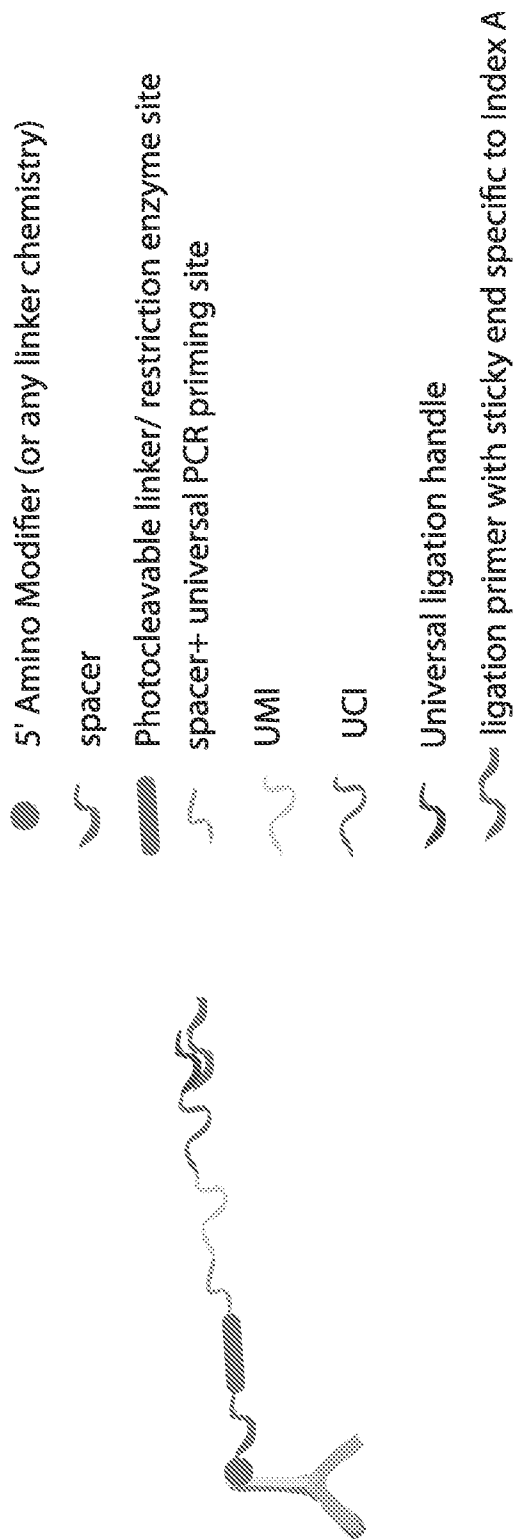


FIG. 5

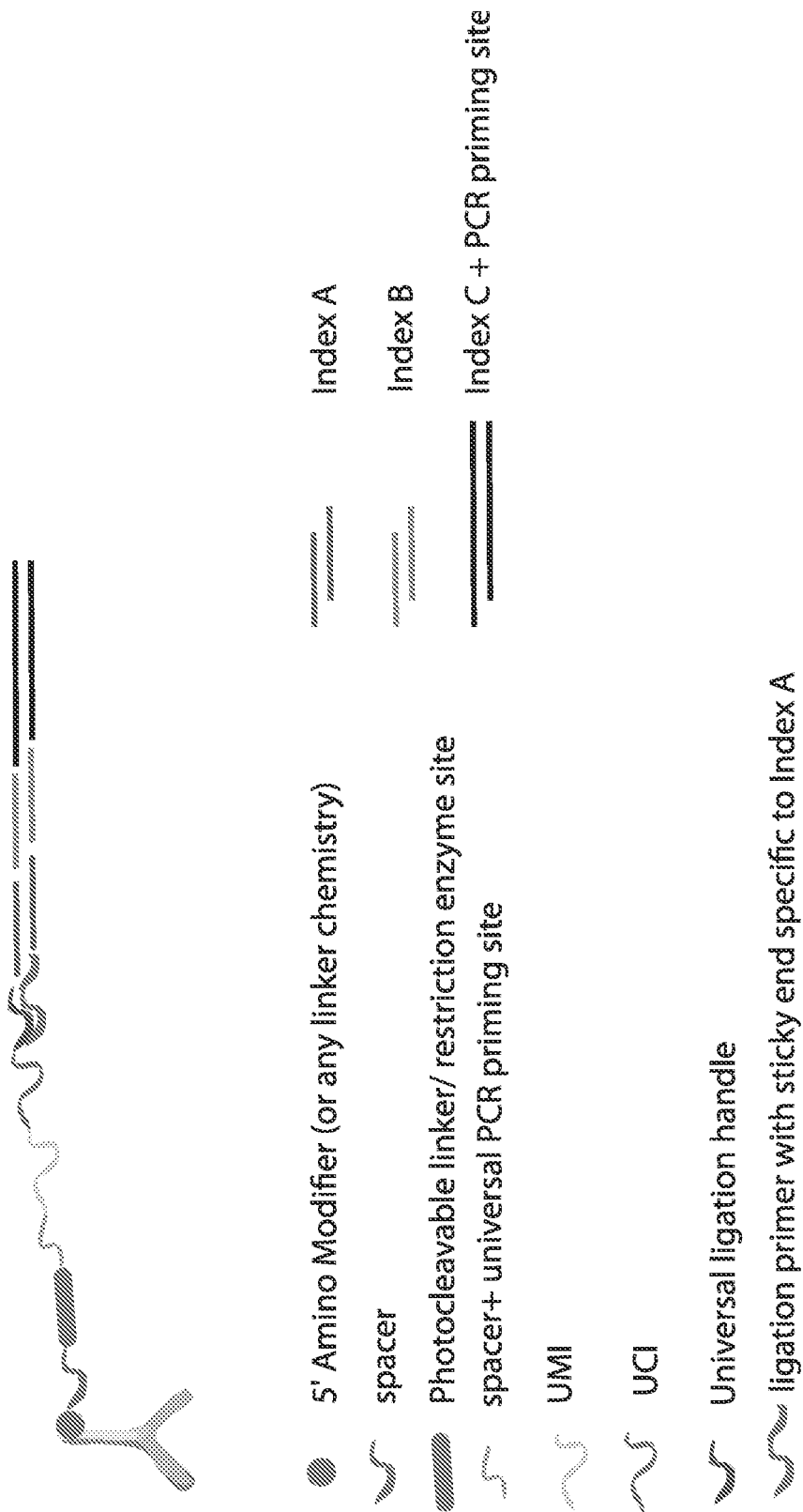


FIG. 6

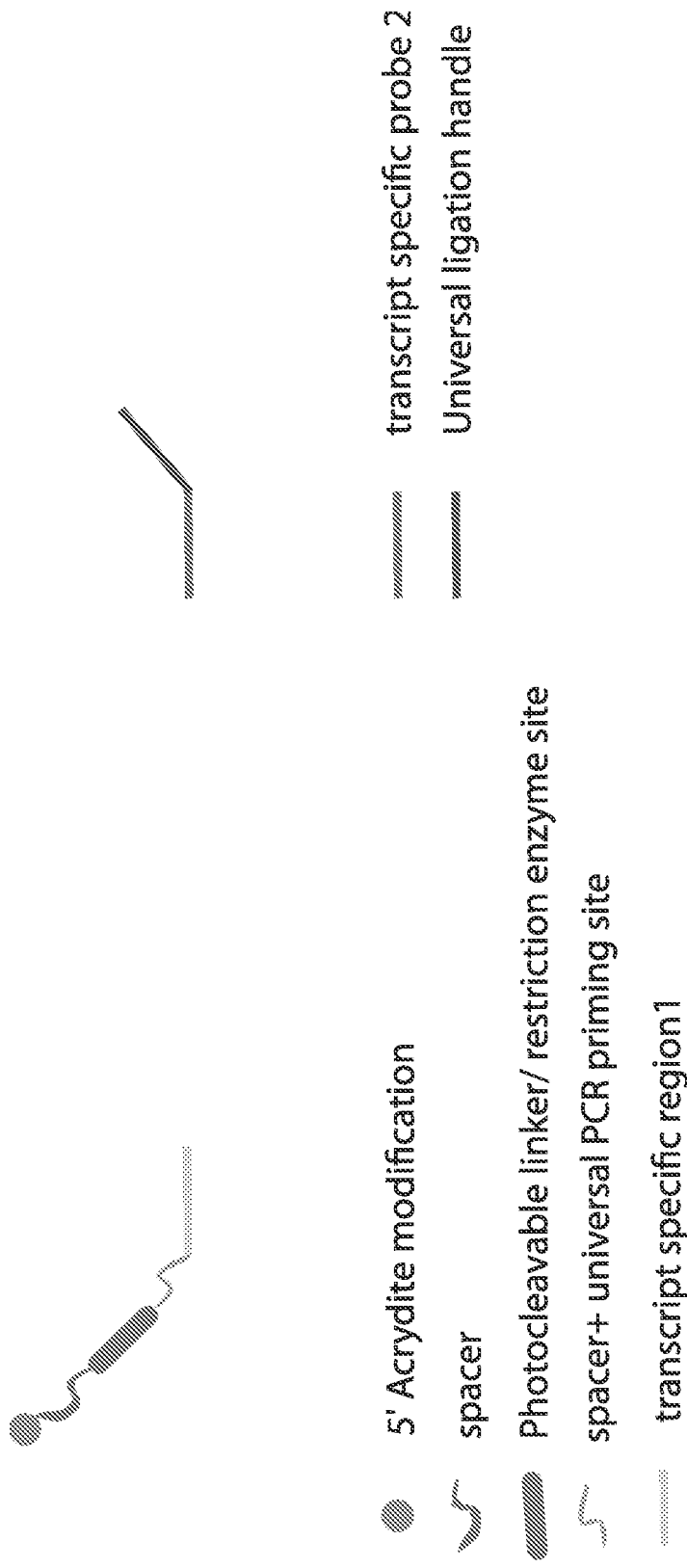
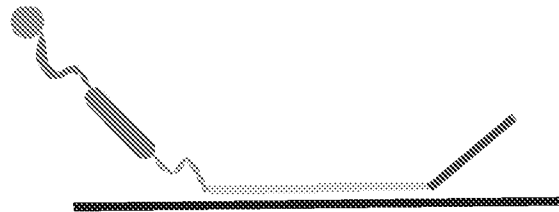


FIG. 7



- | | | | |
|---|--|---|---------------------------|
| ● | 5' Acrydite modification | | |
| ~ | spacer | ▨ | Universal ligation handle |
| ▨ | Photocleavable linker/ restriction enzyme site | | |
| ~ | spacer+ universal PCR priming site | — | RNA transcript |
| ⋯ | transcript specific region | | |

FIG. 8

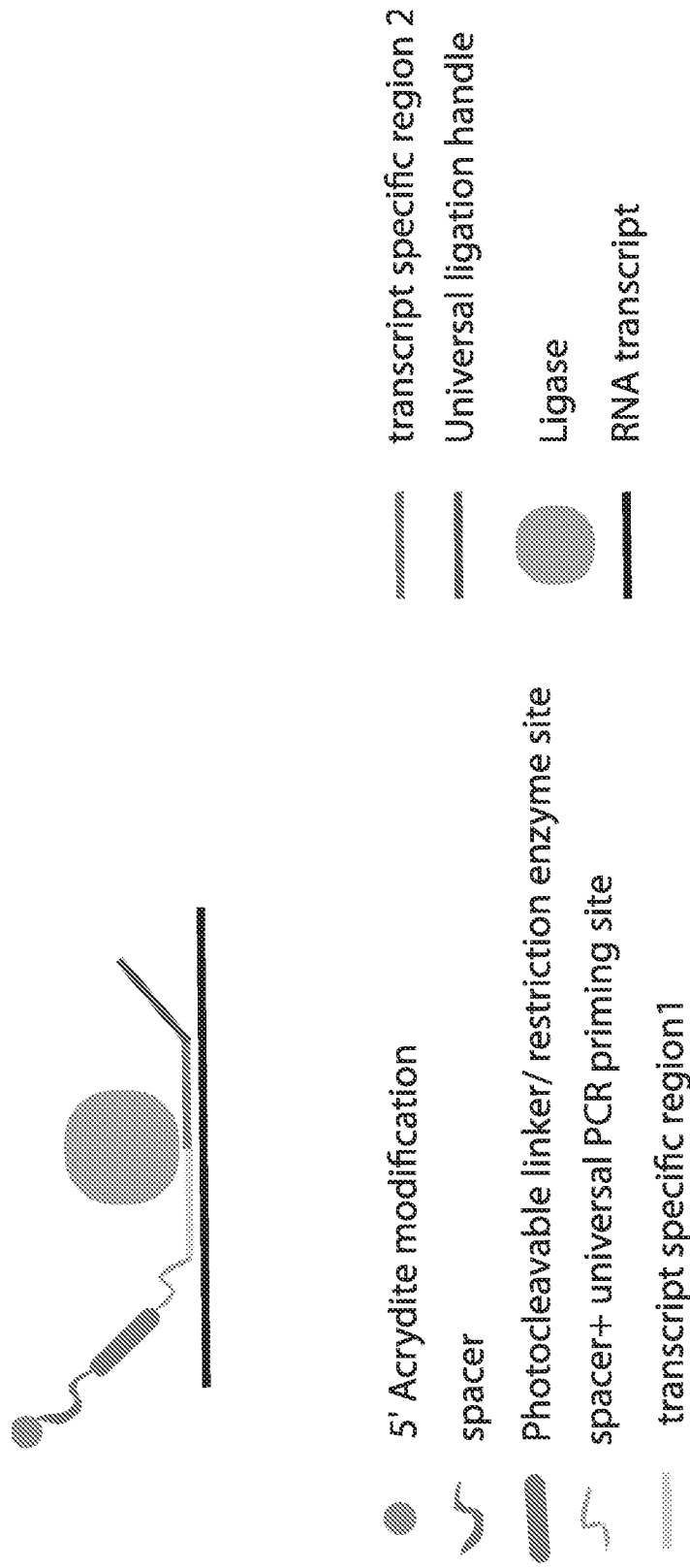


FIG. 9

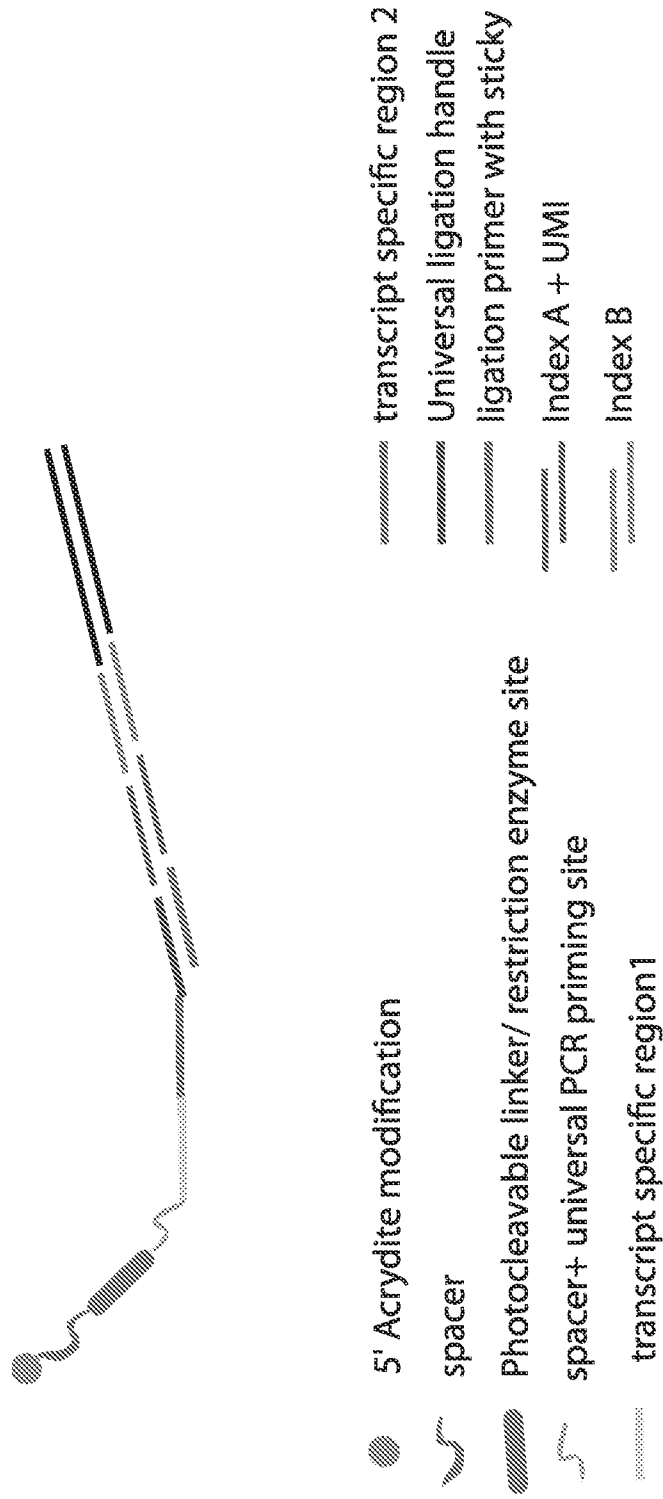
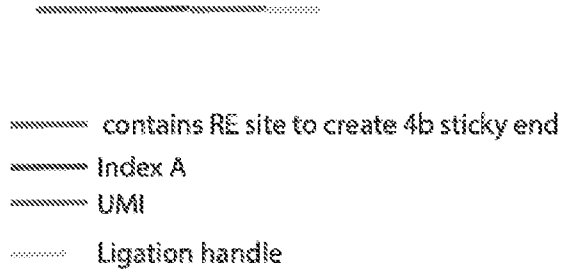
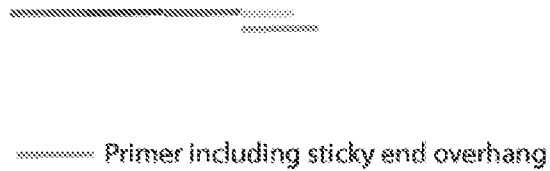


FIG. 10

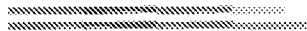
1. Synthesize ssDNA with the following structure:



2. hybridize primer providing sticky end at 3' end



3. DNA polymerase for second strand synthesis



4. Restriction enzyme digestion to generate 5' sticky end

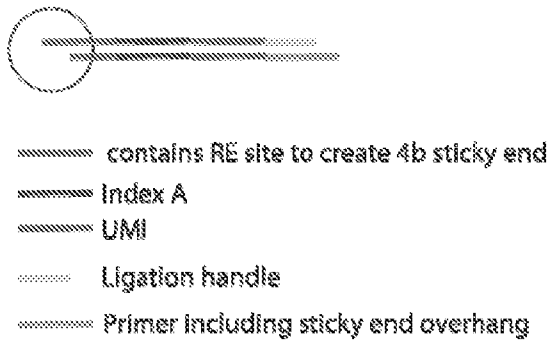


FIG. 11

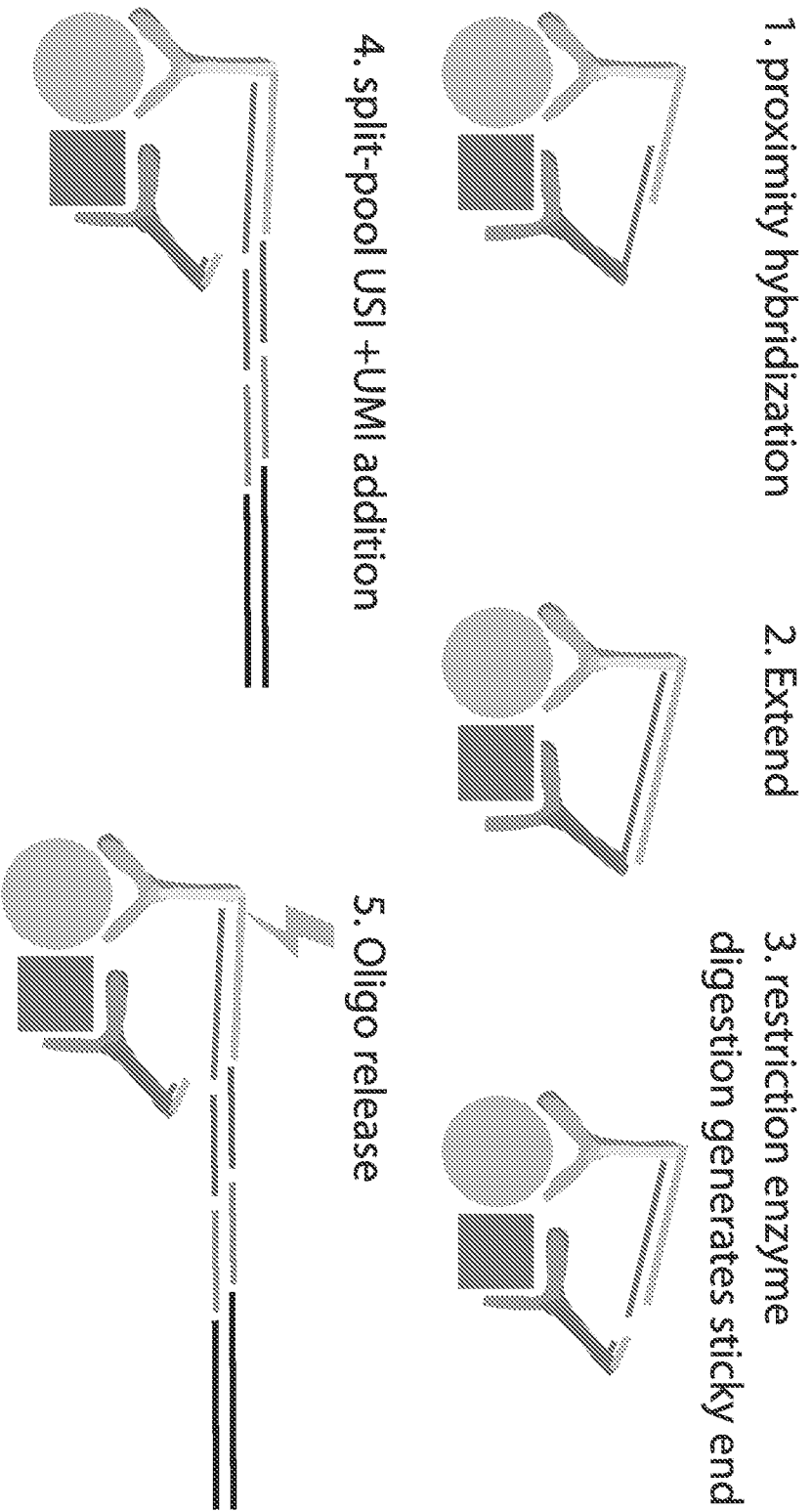


FIG. 12

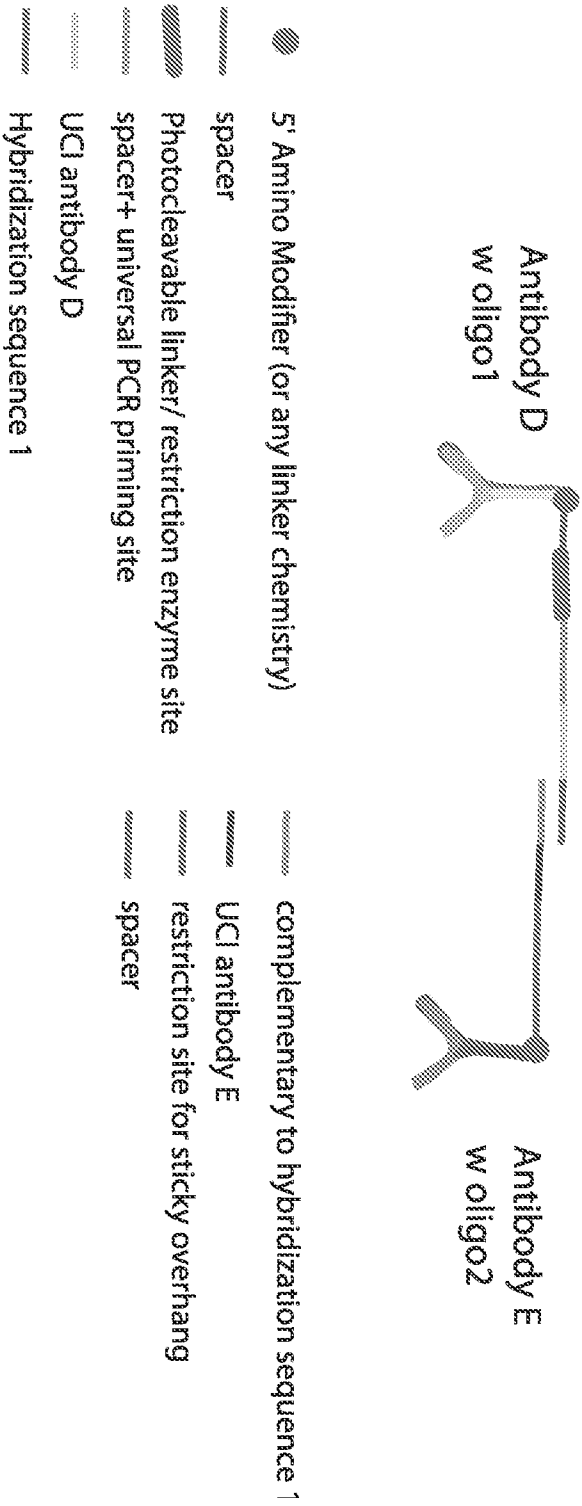
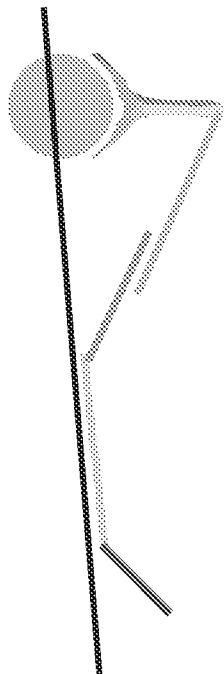
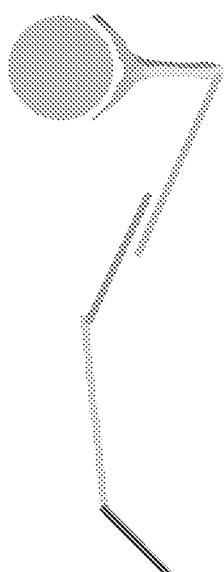


FIG. 13

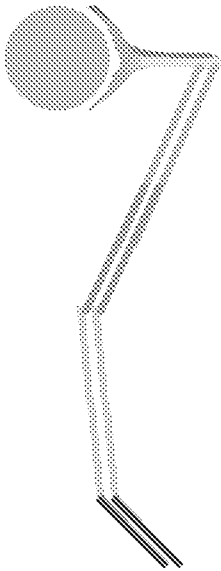
1. proximity hybridization



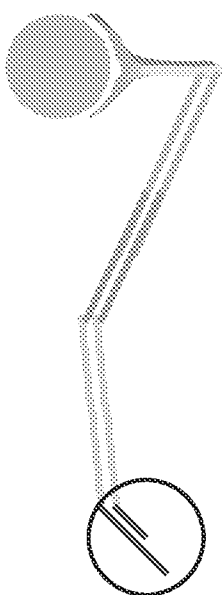
2. RNase



3. Extend



4. restriction enzyme digestion

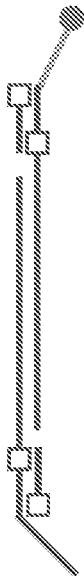


As previously continue to:

- 5. split-pool USI + UMI addition
- 6. Oligo release
- 7. Amplification + Sequencing

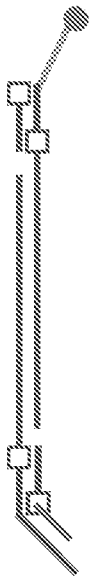
FIG. 14

1. tagmentation of genomic DNA (Tn5)



2. polymerize hydrogel drop, incorporates genomic DNA into mesh.

3. anneal ligation primer, generating sticky end



4. extend



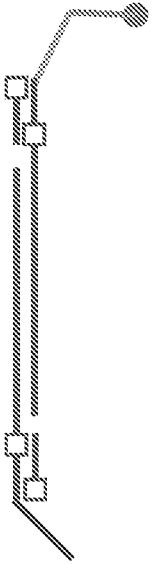
As previously continue to:

5. split-pool USI + UMI addition

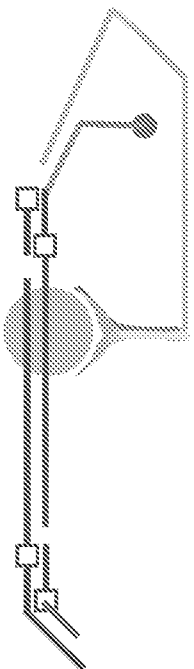
6. Amplification (T7/PCR) + Sequencing

FIG. 15

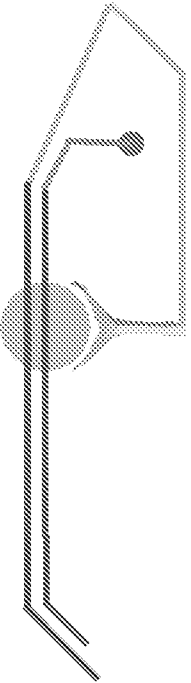
1. tagmentation of genomic DNA (Tn5)
+ polymerize in hydrogel drop



2. hybridize proximity probes and
ligation handle primer (→)



3. Extend



As previously continue to:

- 4. split-pool USI + UMI addition
- 5. Amplification + Sequencing

FIG. 16

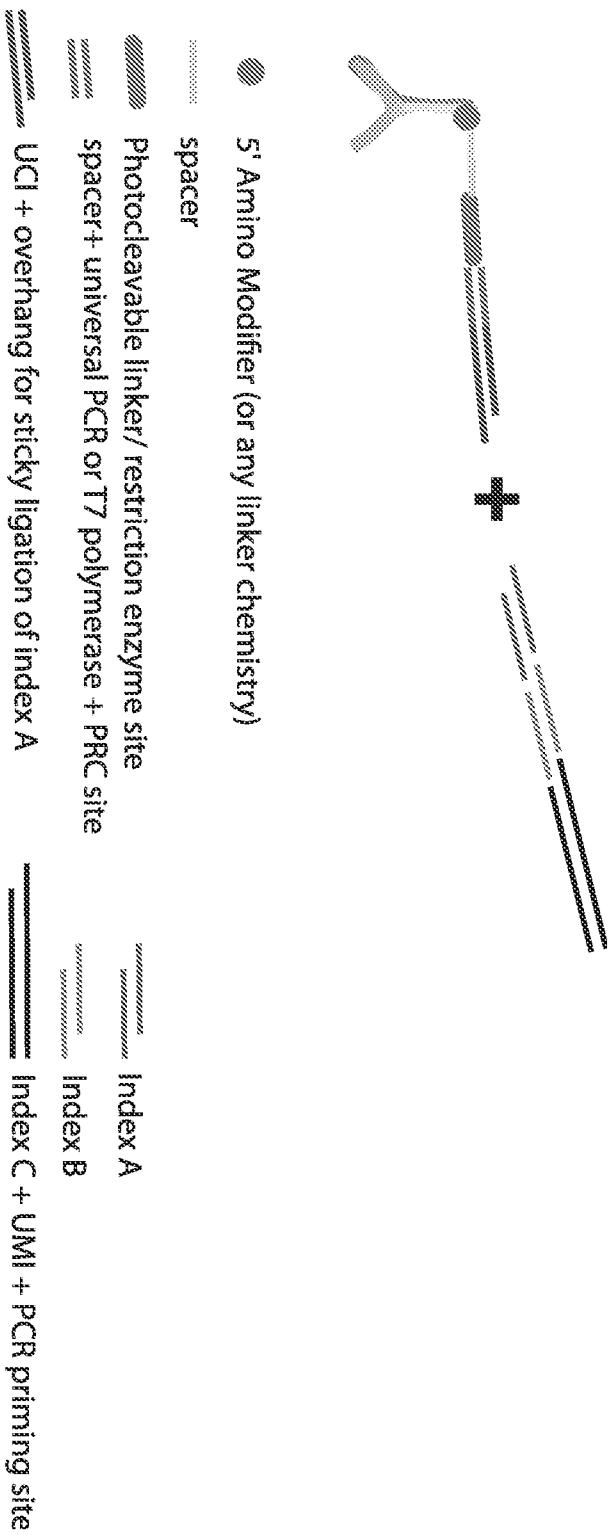


FIG. 17

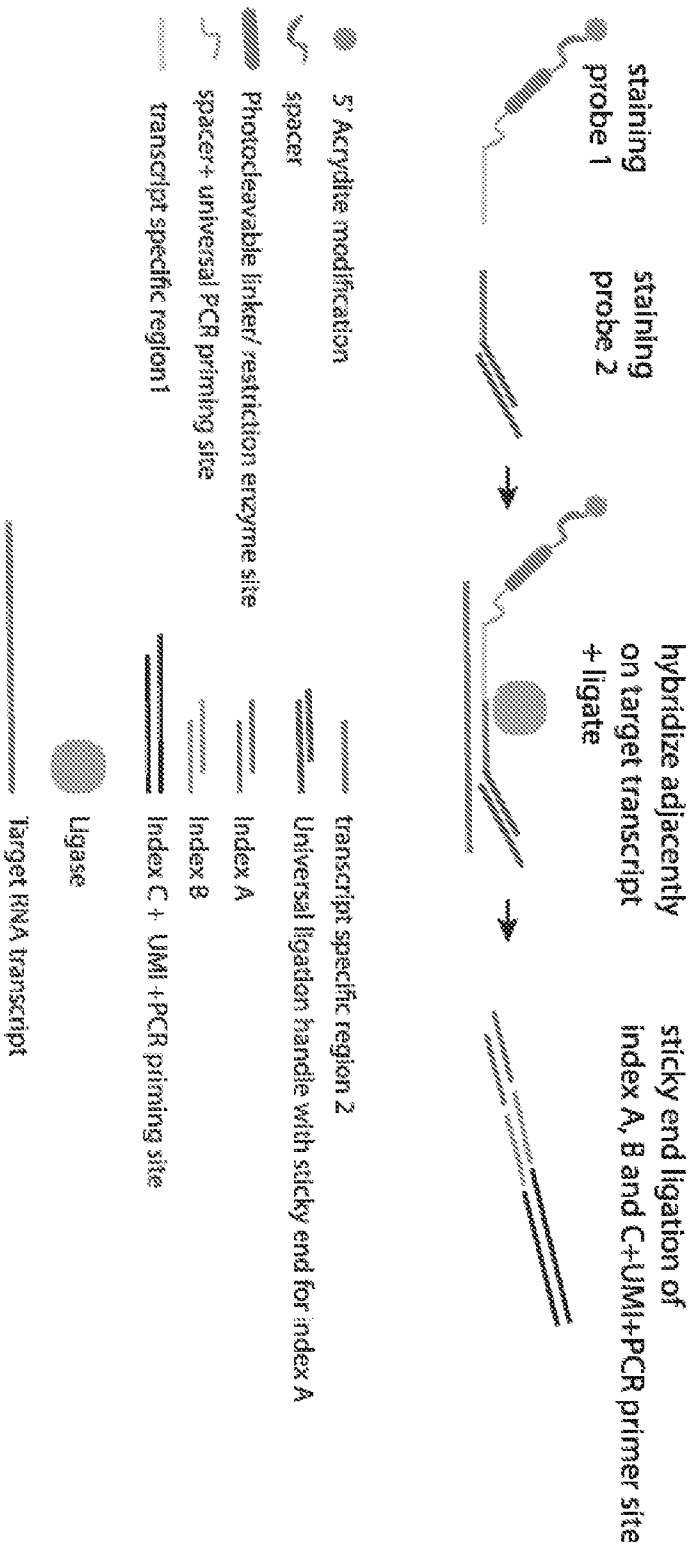
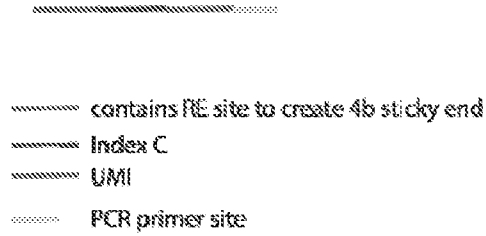
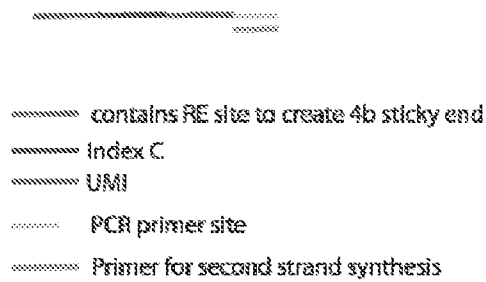


FIG. 18

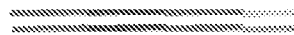
1. Synthesize ssDNA with the following structure:



2. Hybridize universal primer



3. DNA polymerase for second strand synthesis



4. Restriction enzyme digest

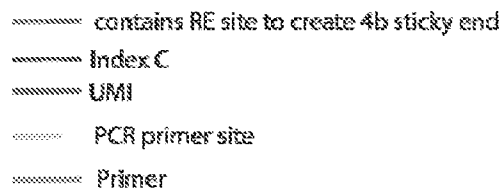
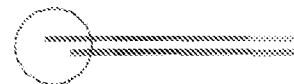


FIG. 19

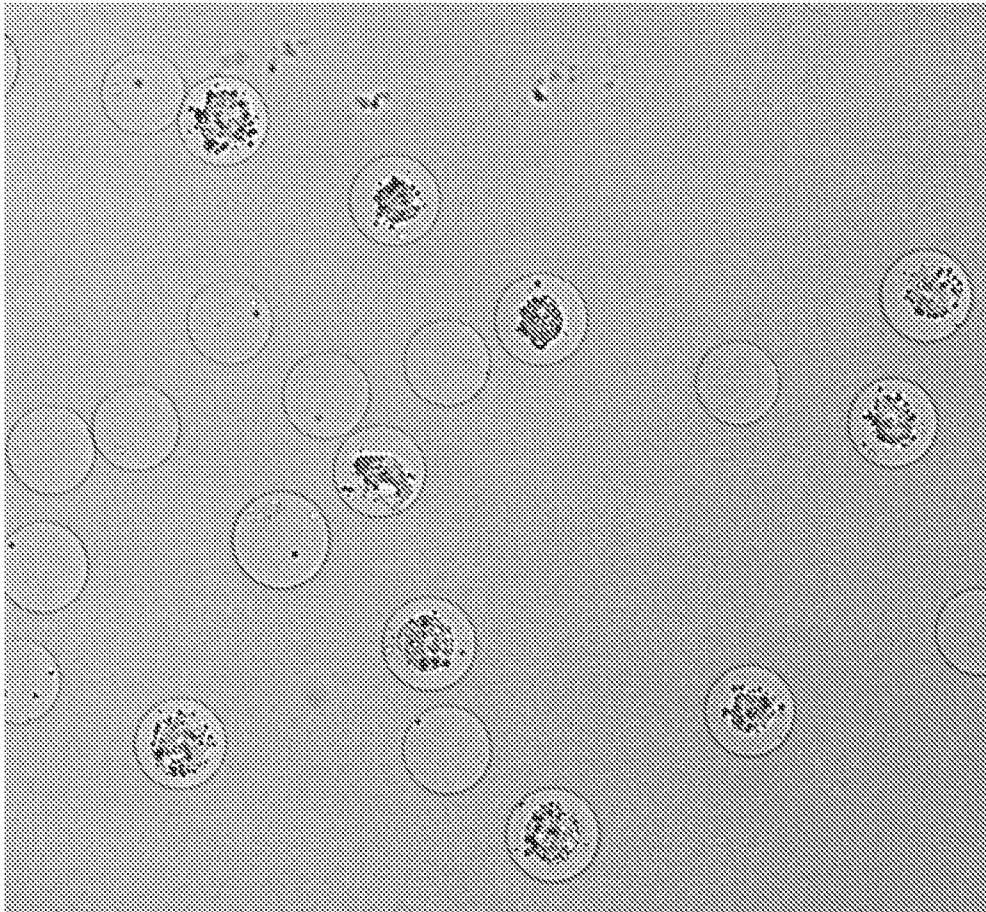
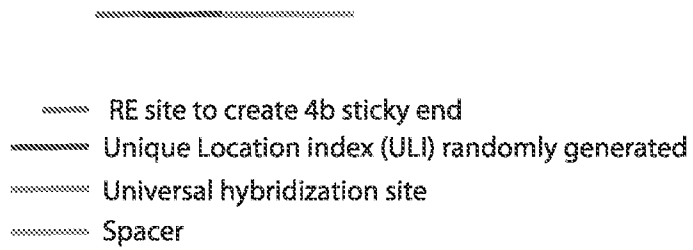
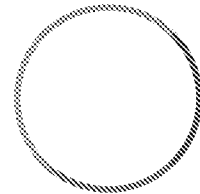


FIG. 20

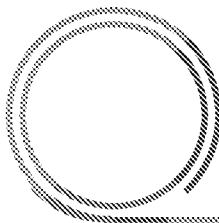
1) Synthesize ssDNA containing:



2) Circularize using CirLigase



3) Perform Rolling circle amplification



4) isolate linear amplicon and use as staining probe


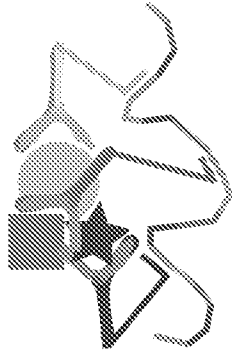
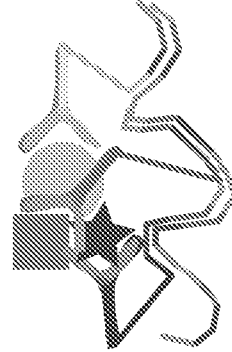


FIG. 21

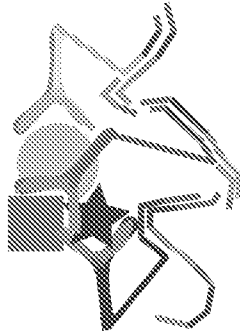
1) stain with ULI-probe



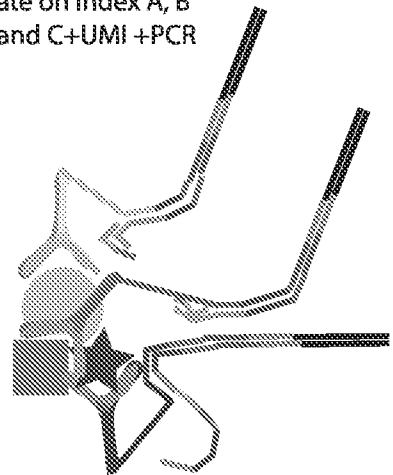
2) extend



3) Restriction enzyme digest generates 4bp overhang for sticky end ligation of USI + UMI



4) ligate on index A, B and C+UMI +PCR



5) Release oligo's from ligand and PCR or T7 amplify, each sequence to amplify (top strand) has:
 PCR FWD + UCI + universal hybridization site + ULI + USI + UMI + PCR Rev

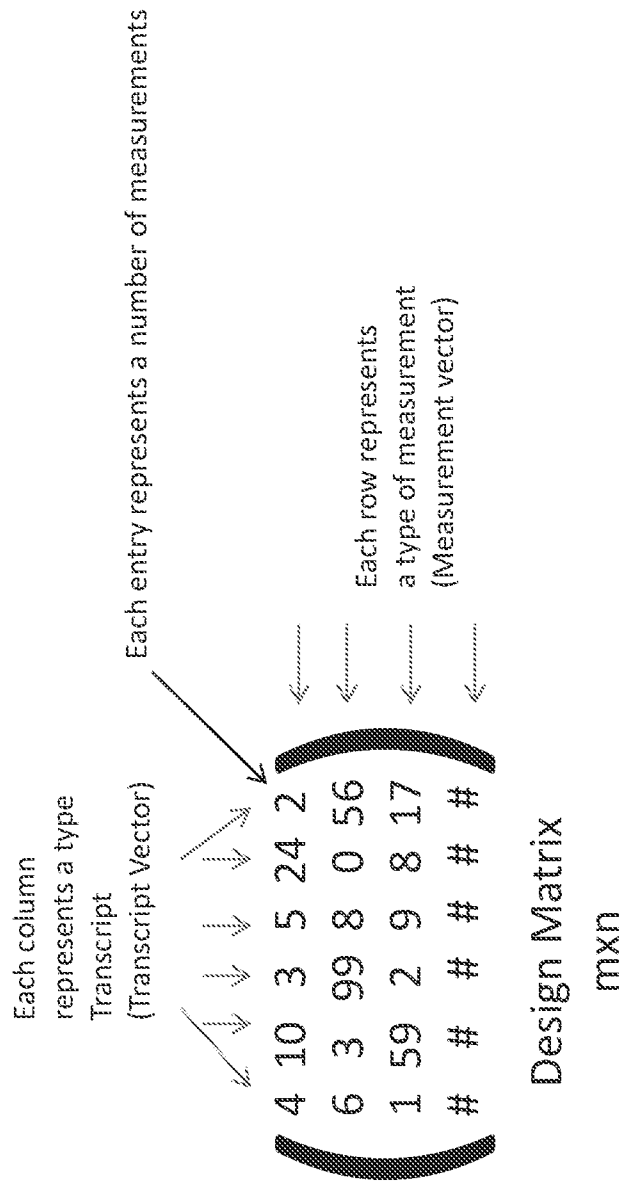


FIG. 22

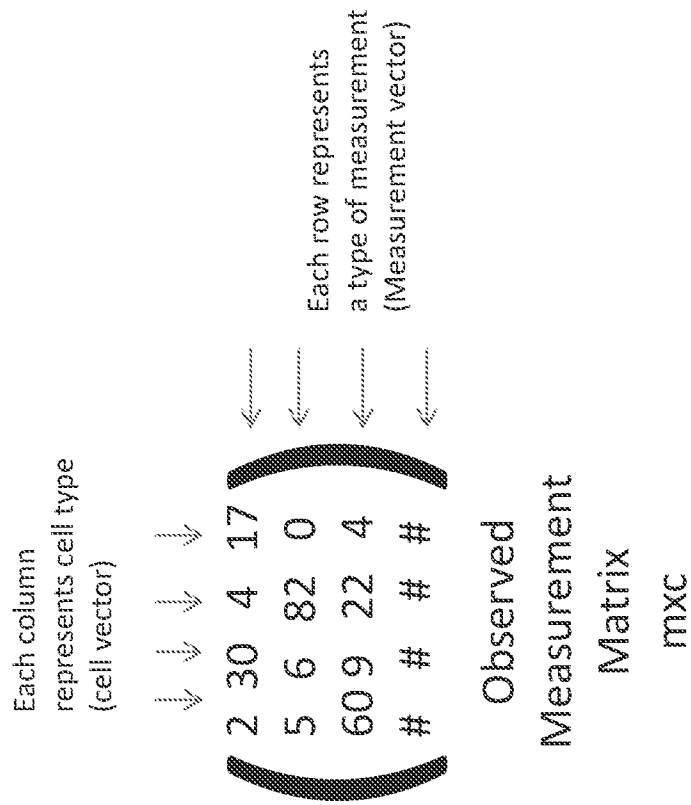


FIG. 23

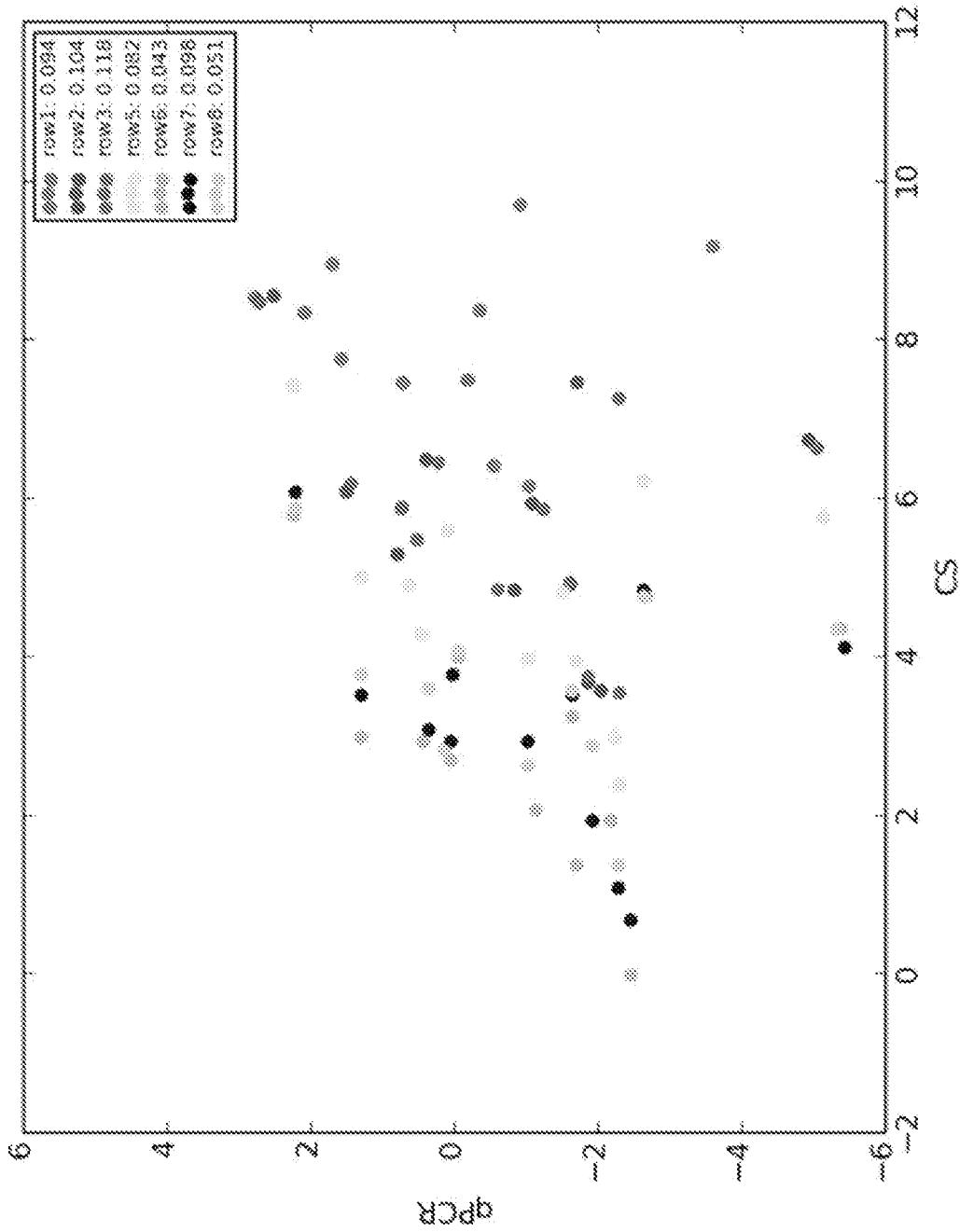


FIG. 24

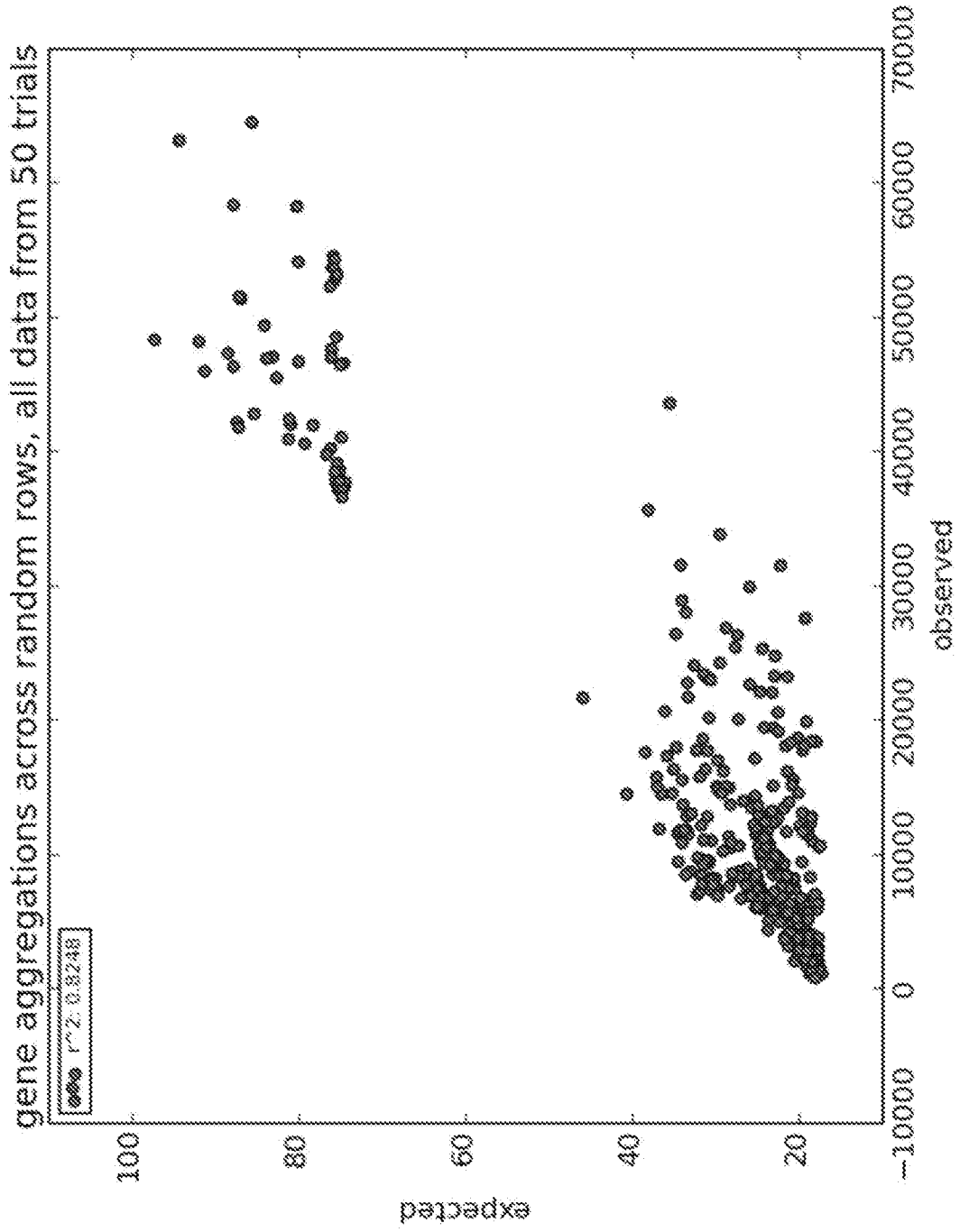


FIG. 25

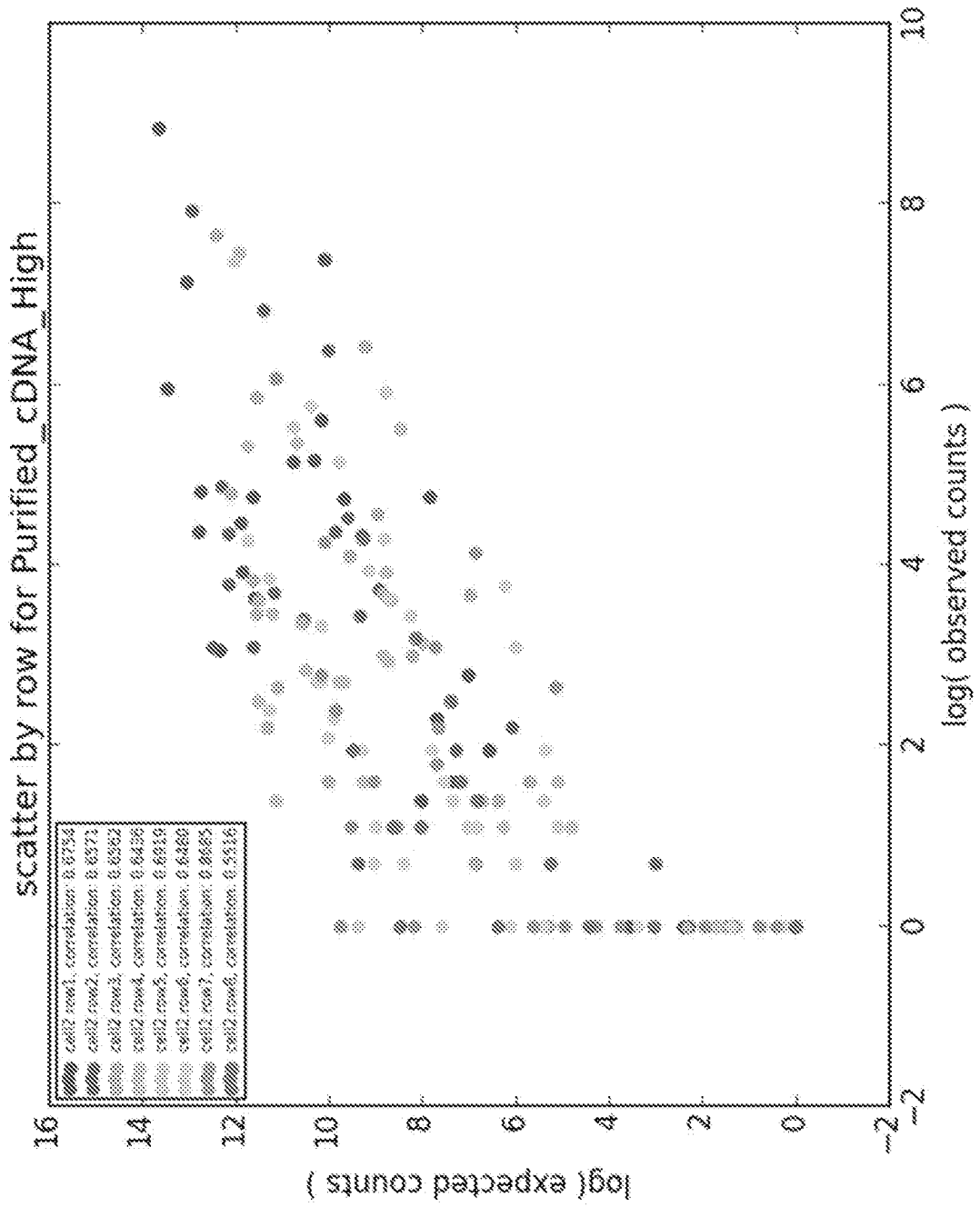


FIG. 26

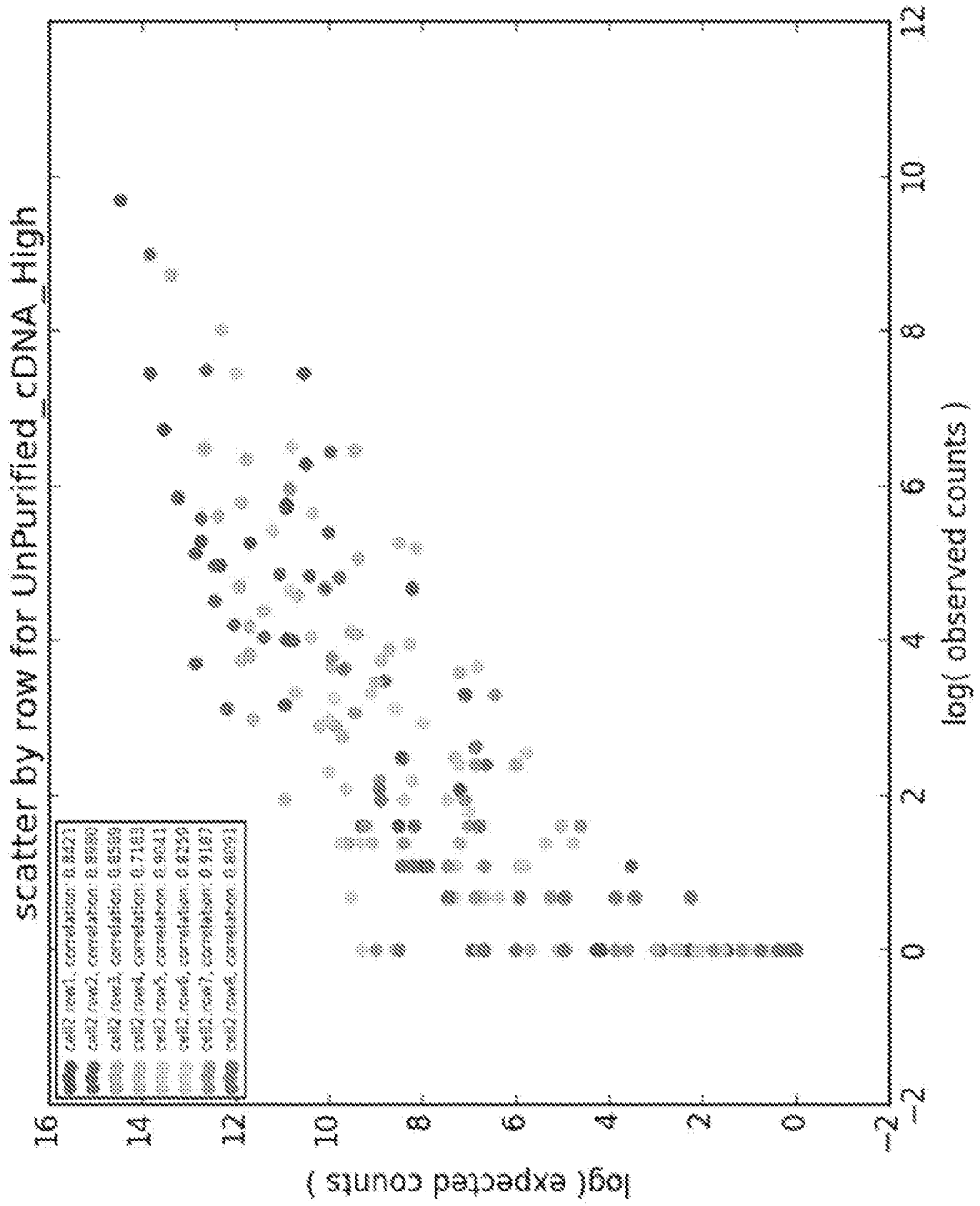


FIG. 27

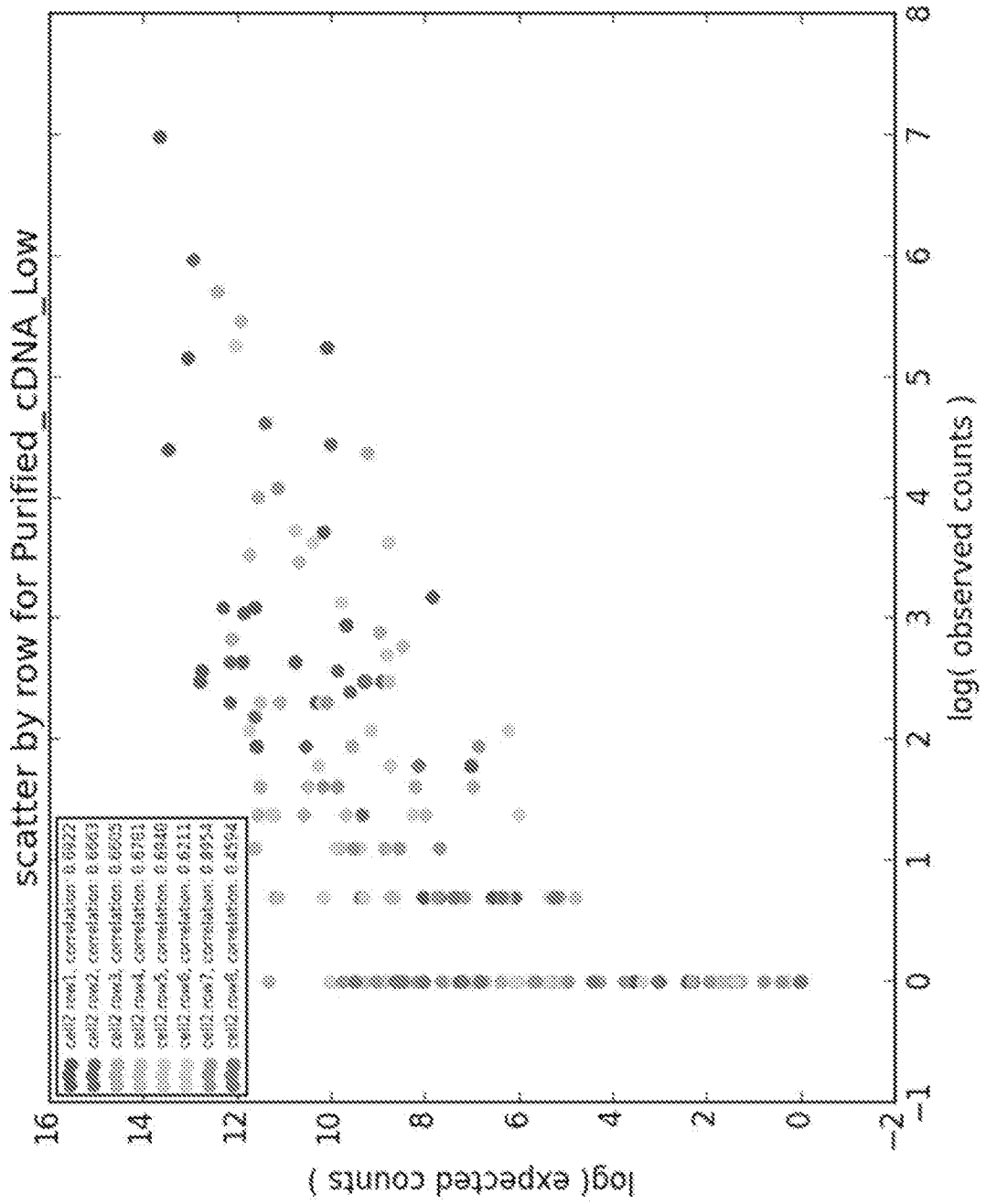


FIG. 28

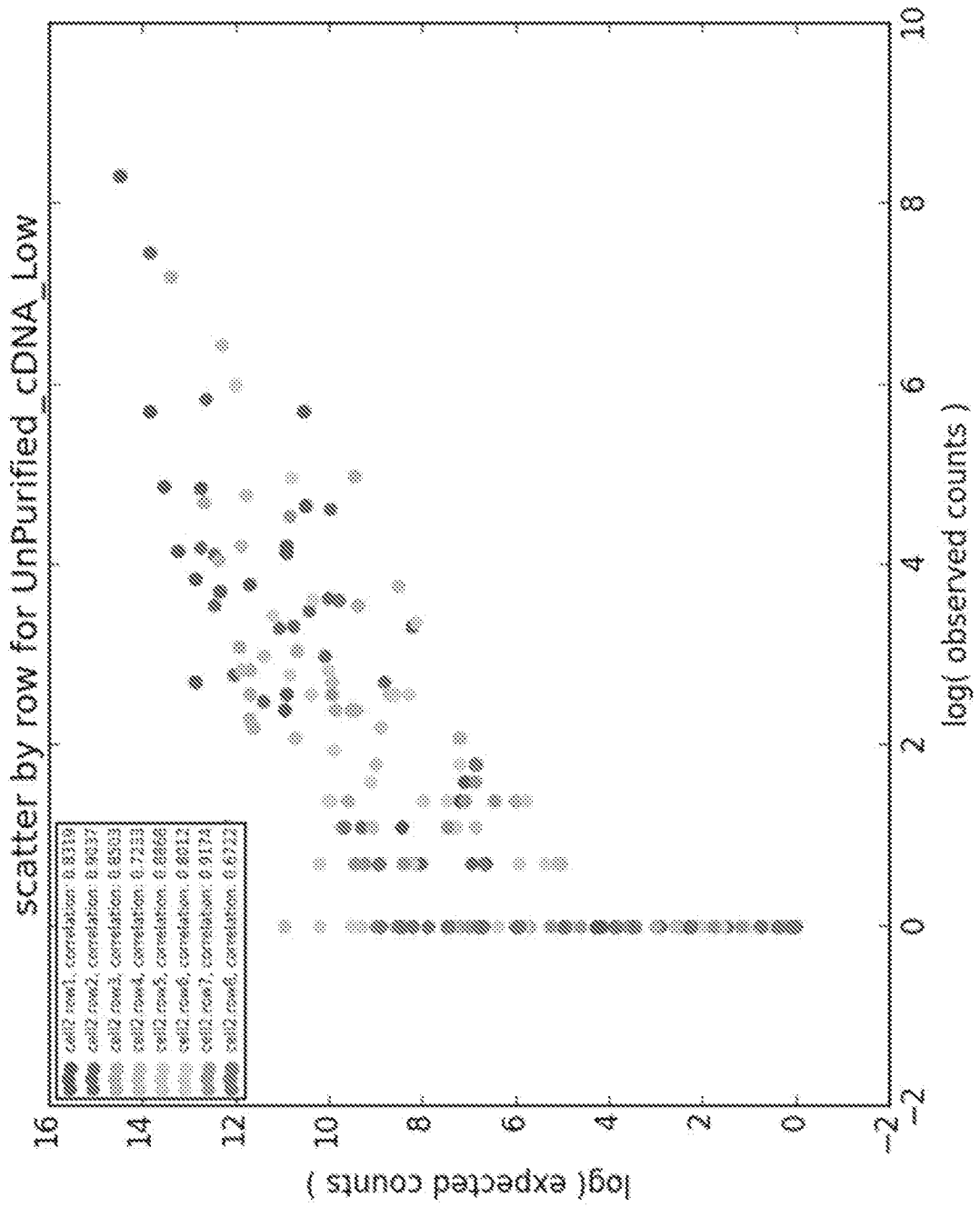


FIG. 29

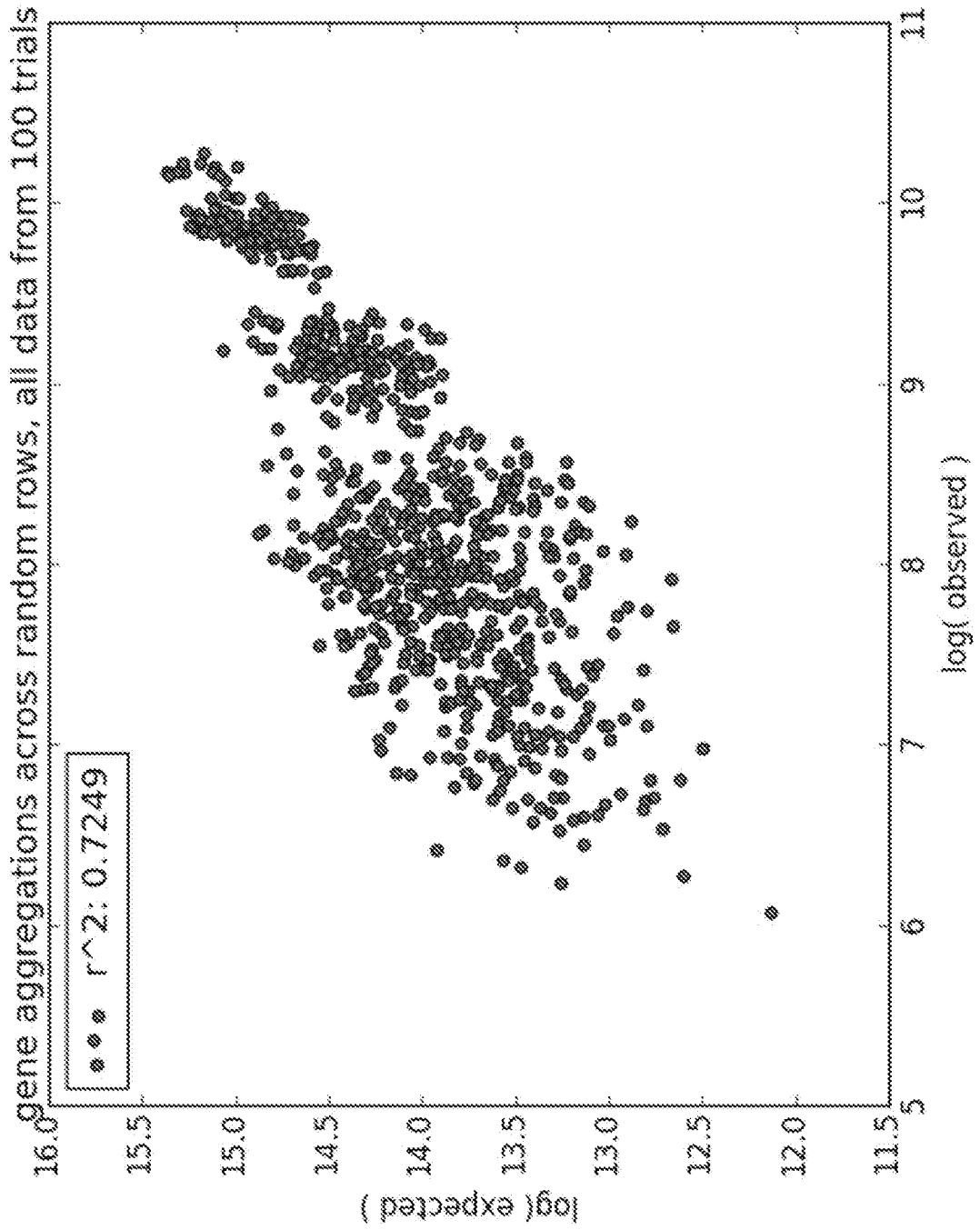


FIG. 30

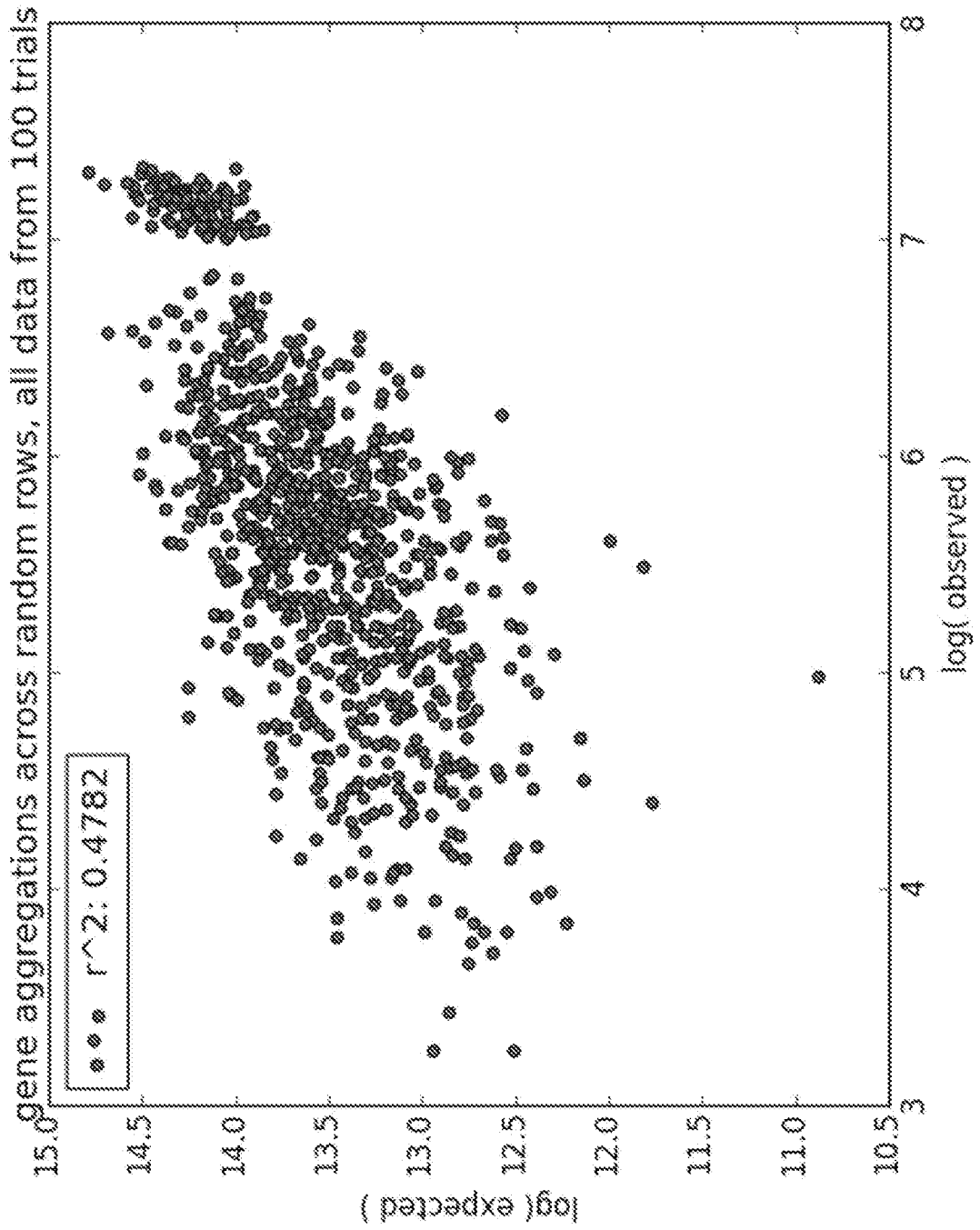


FIG. 31

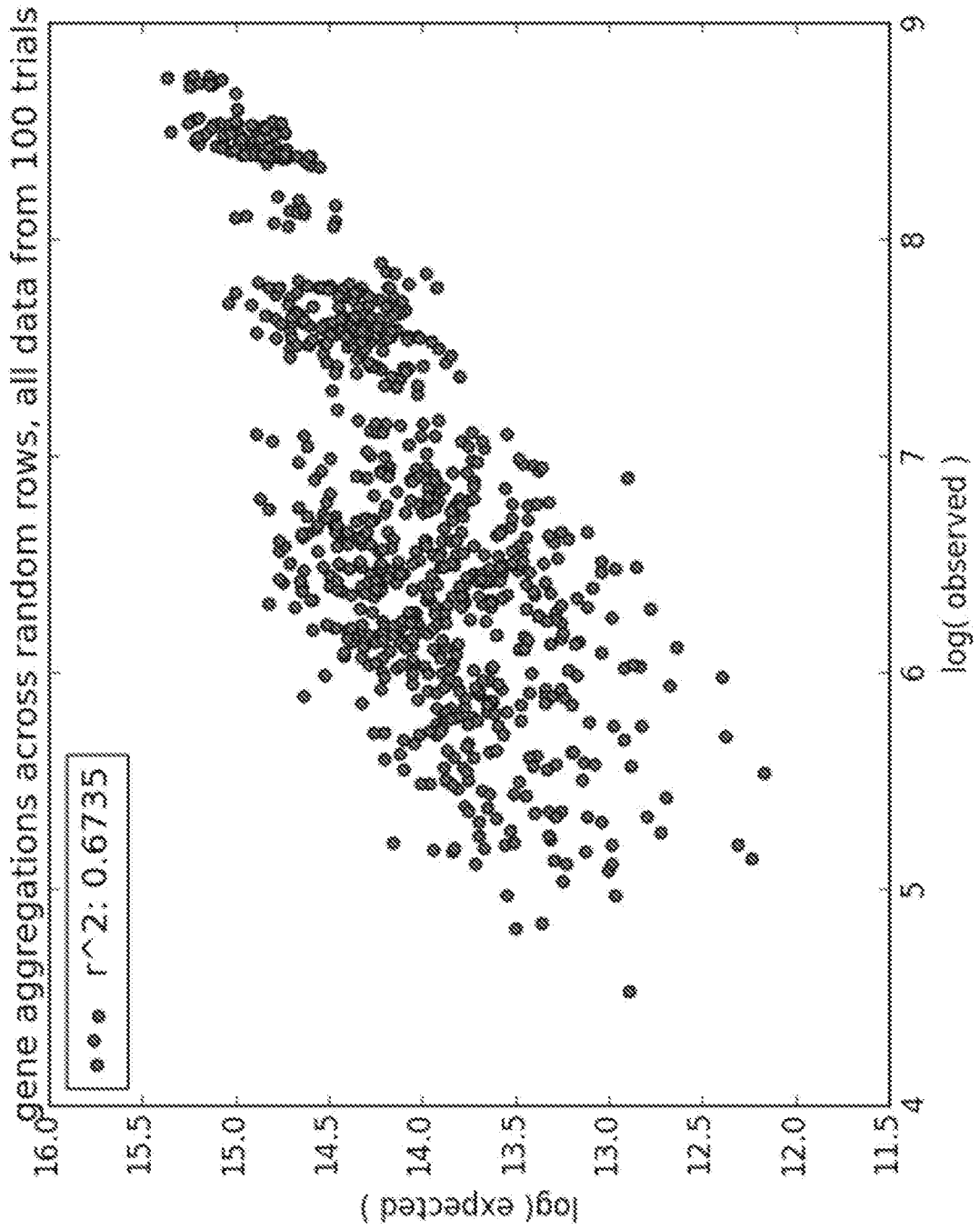


FIG. 32

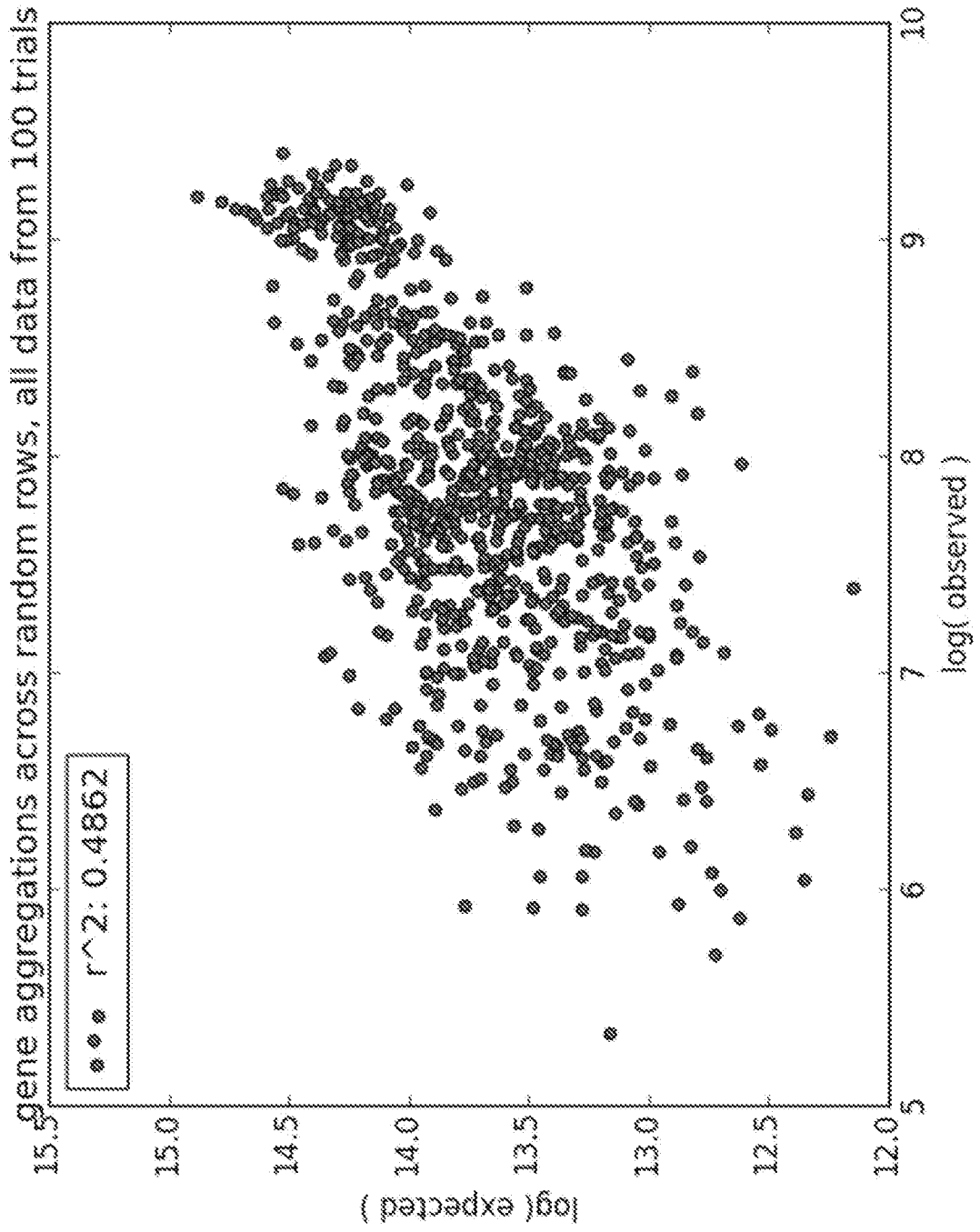


FIG. 33

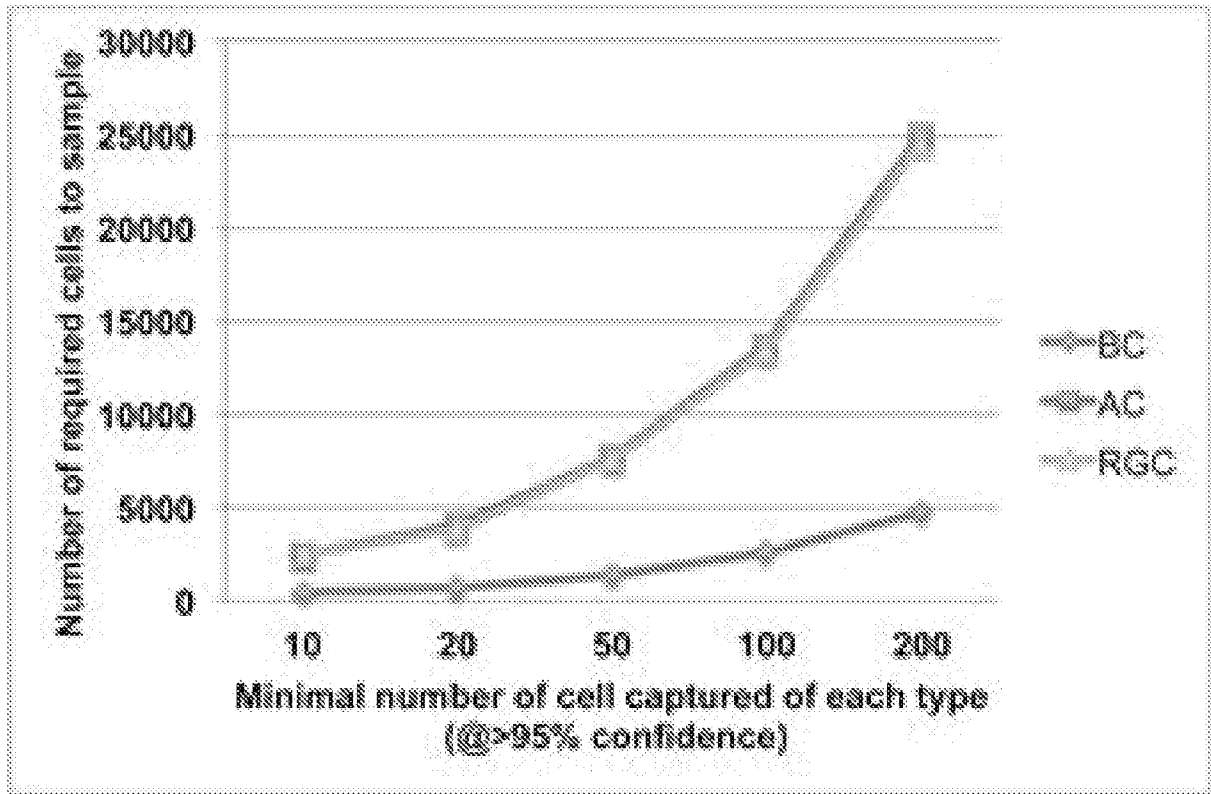


FIG. 34

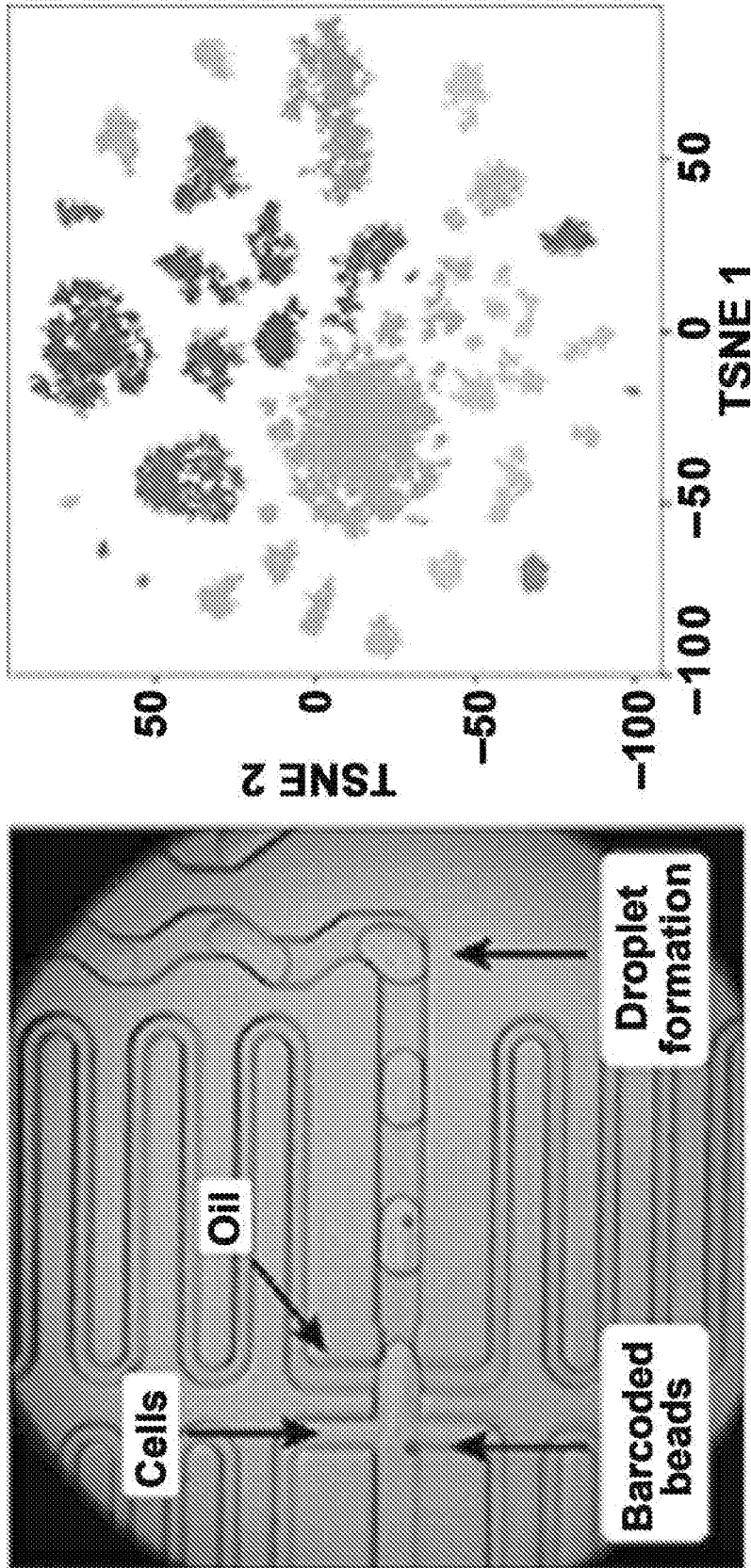


FIG. 35

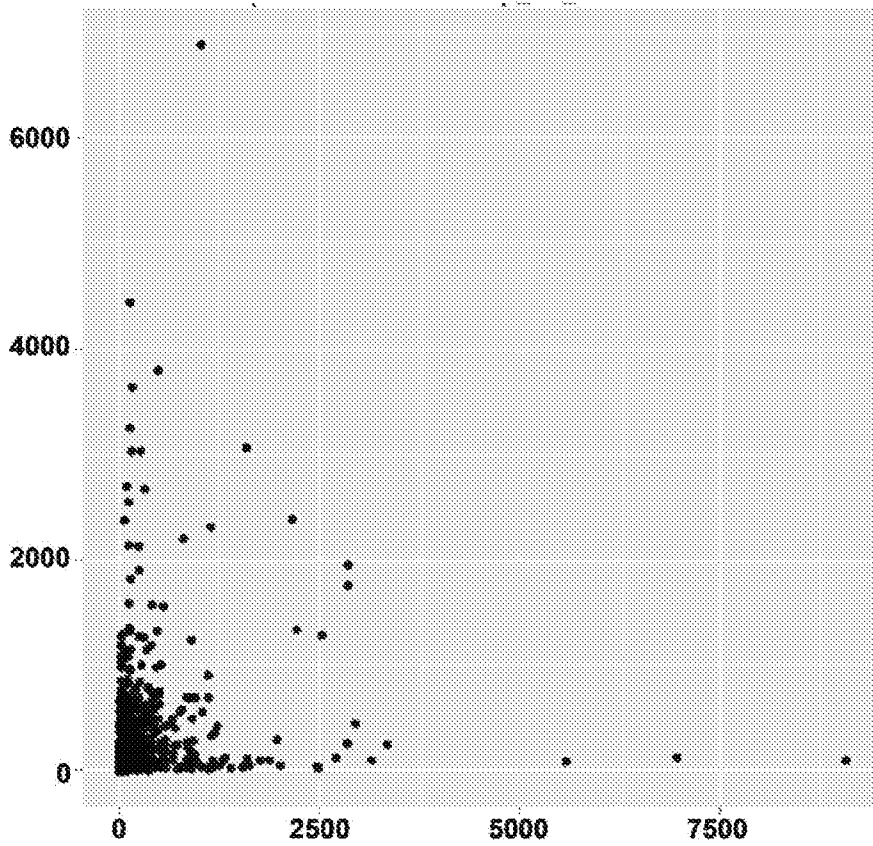
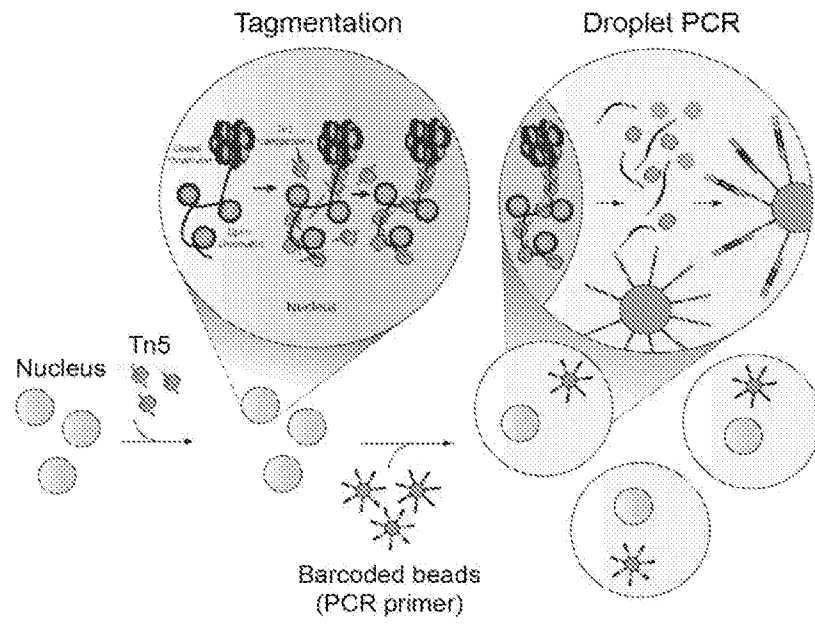


FIG. 36

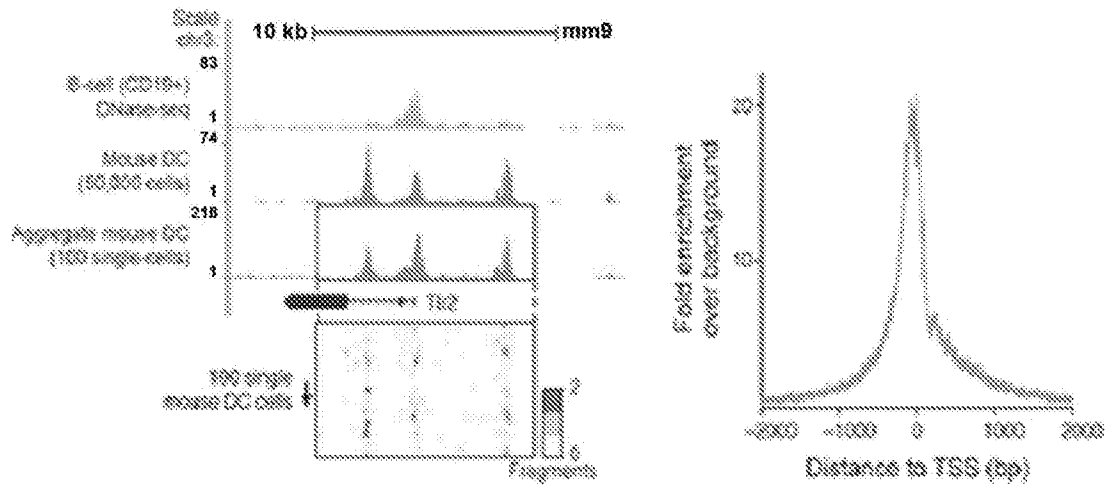


FIG. 37

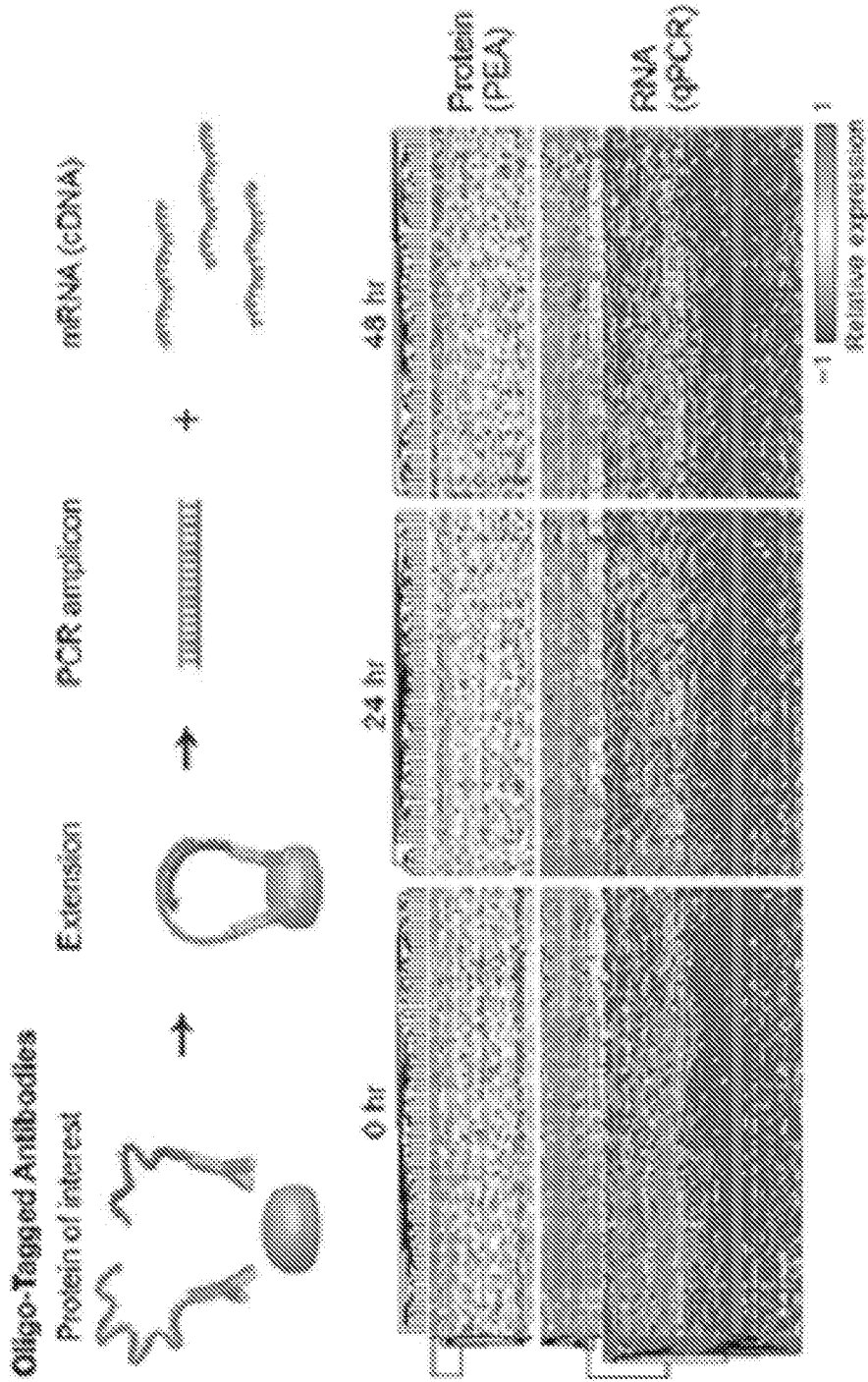


FIG. 38

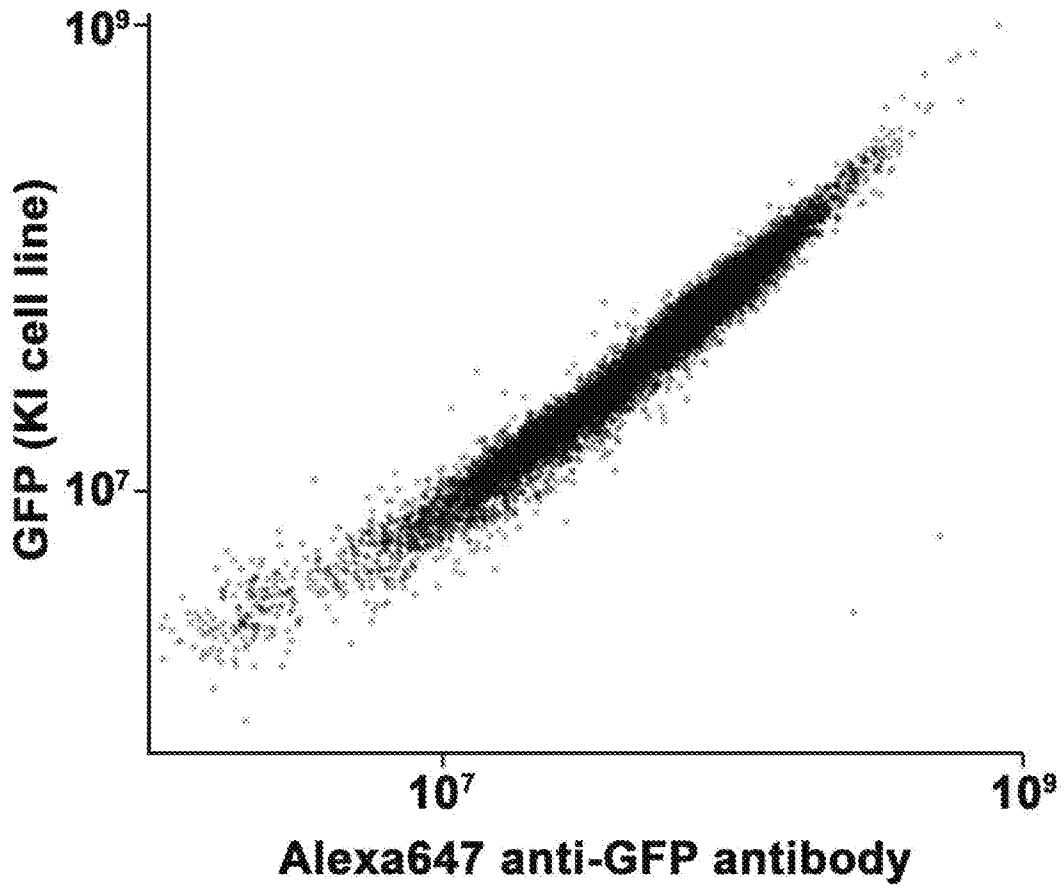
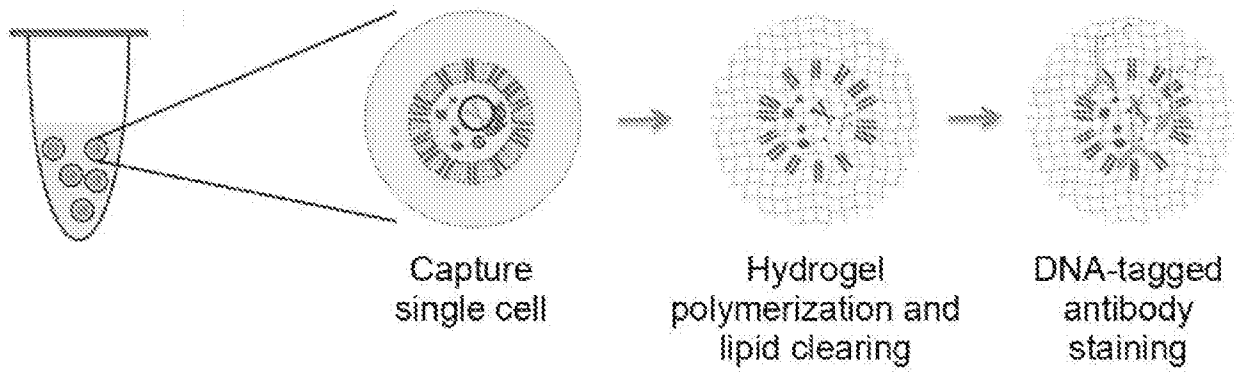


FIG. 39

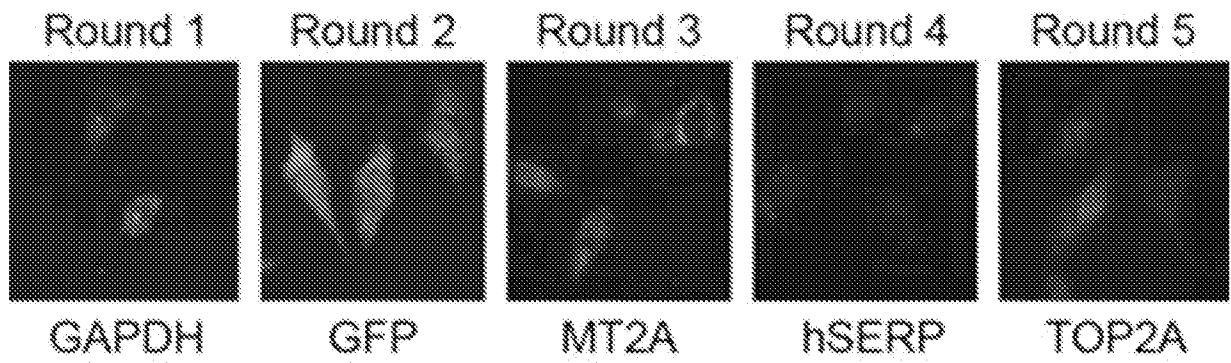


FIG. 40

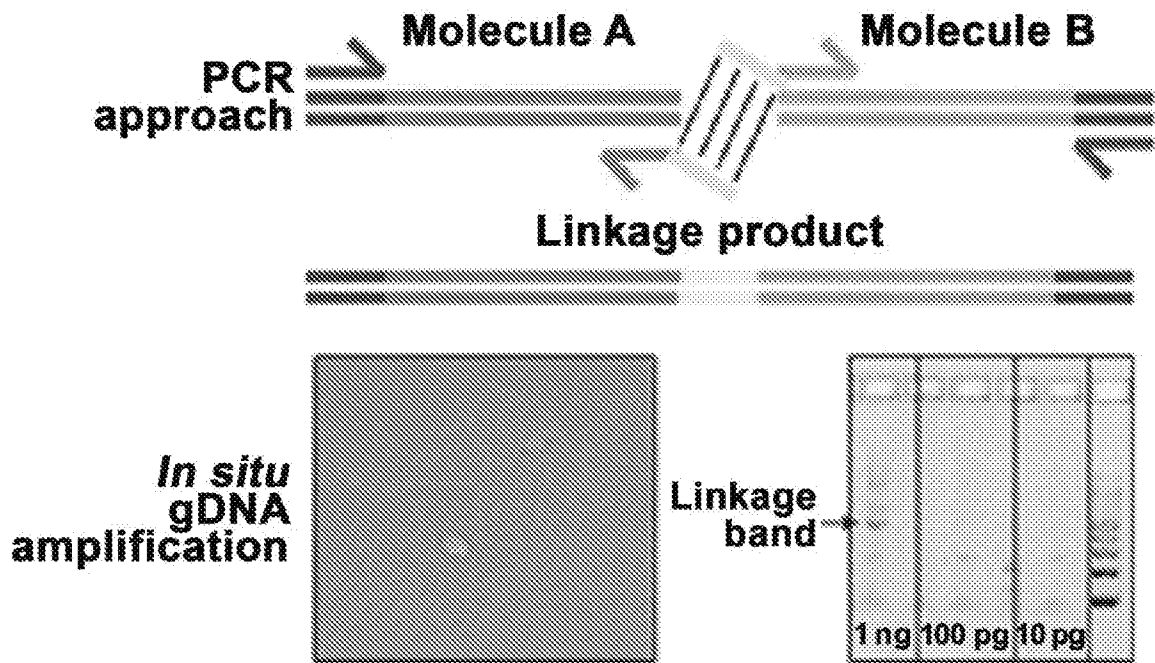


FIG. 41

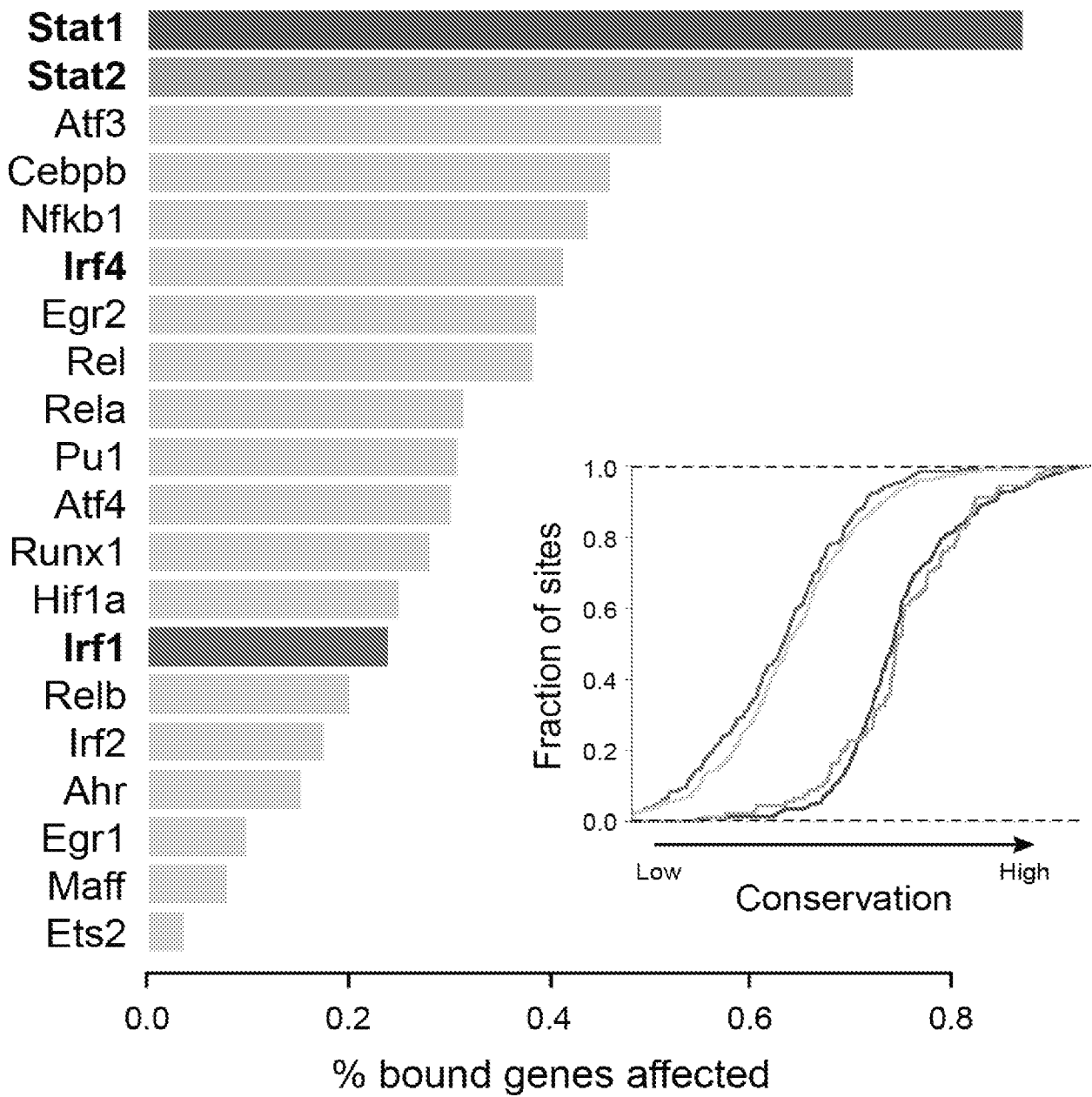


FIG. 42

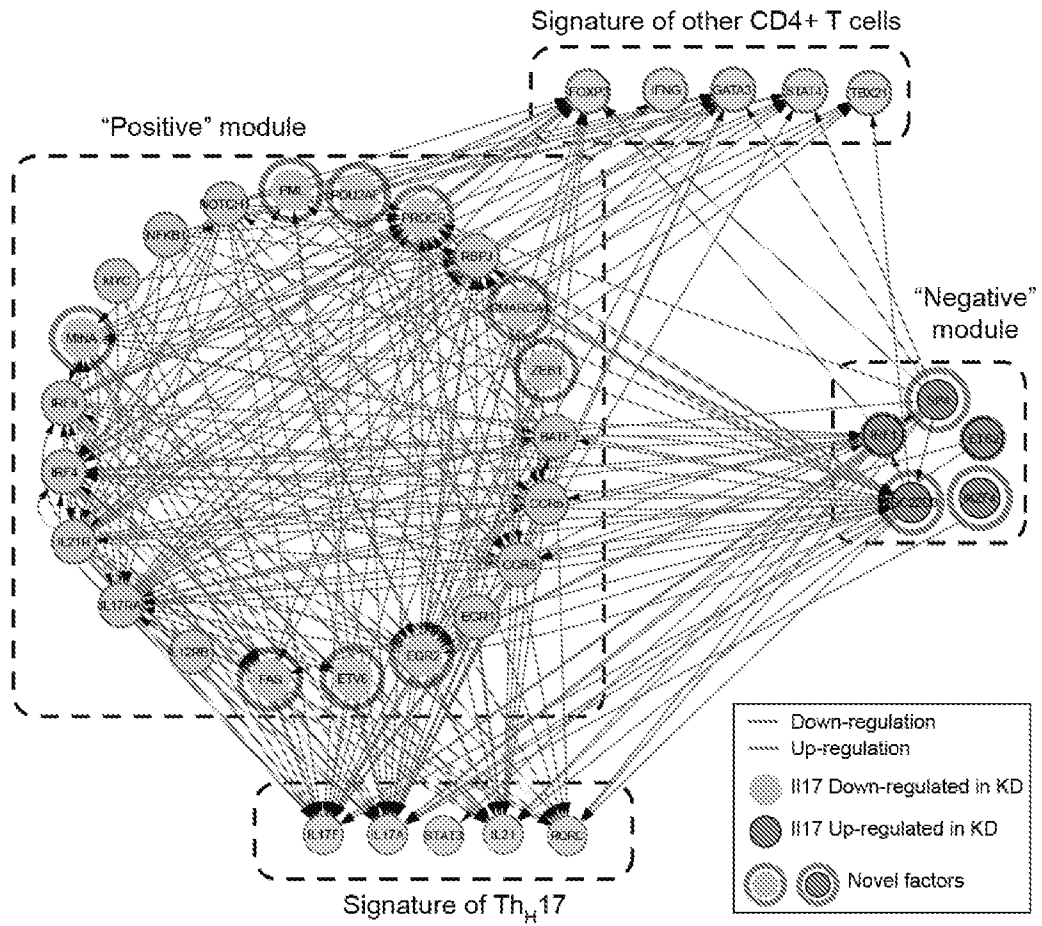


FIG. 43

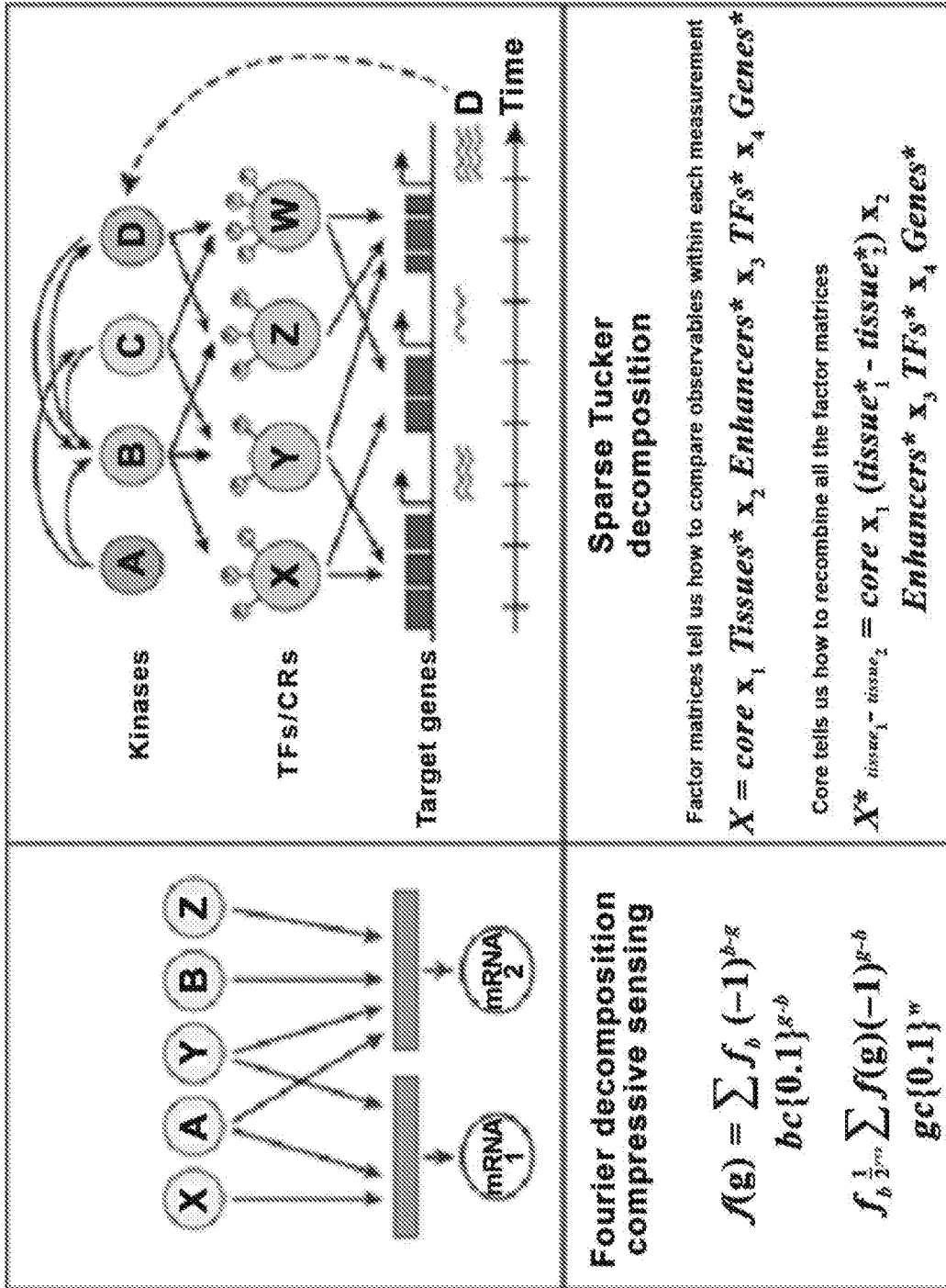


FIG. 44

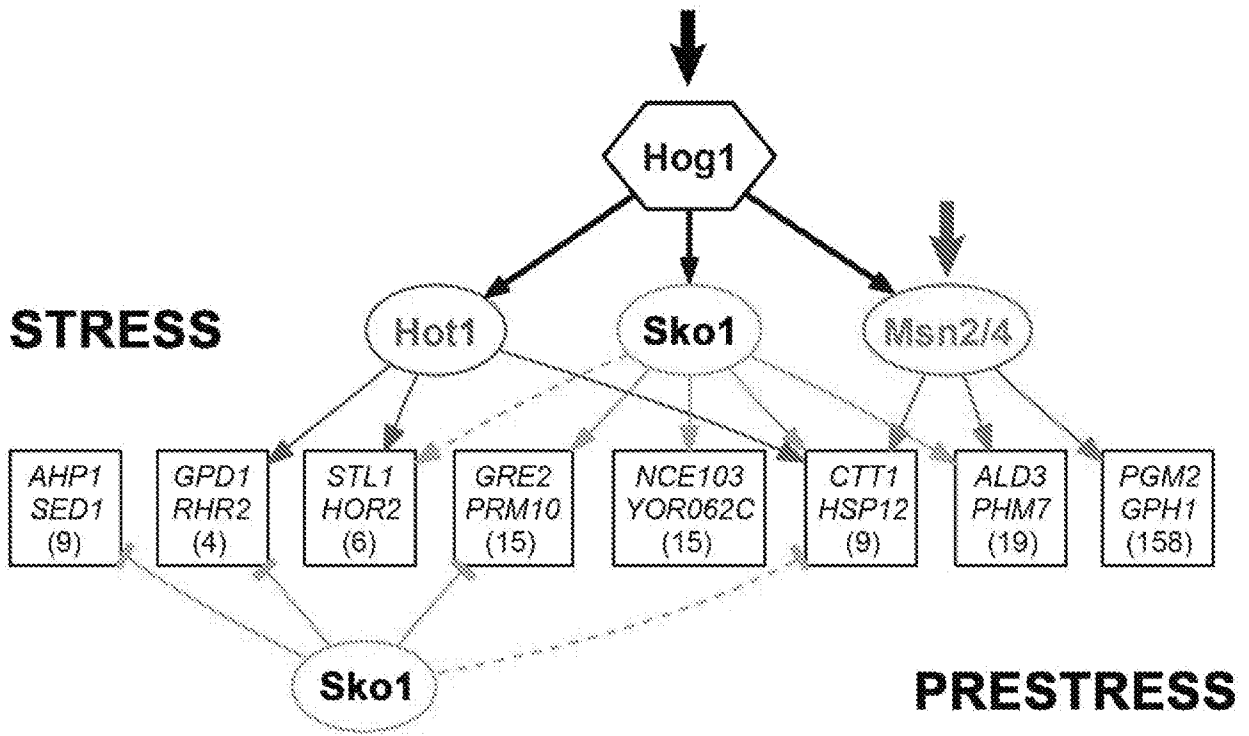


FIG. 45

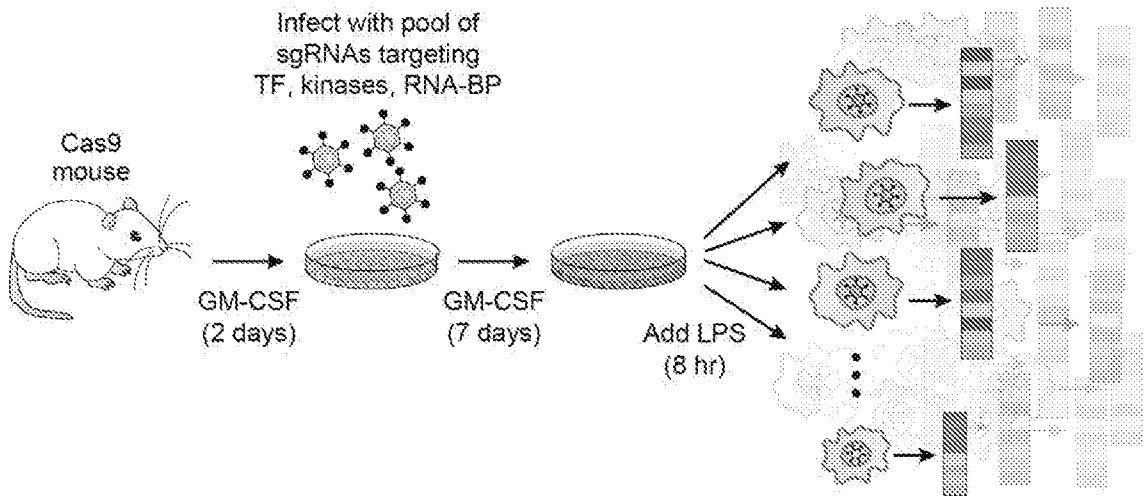


FIG. 46

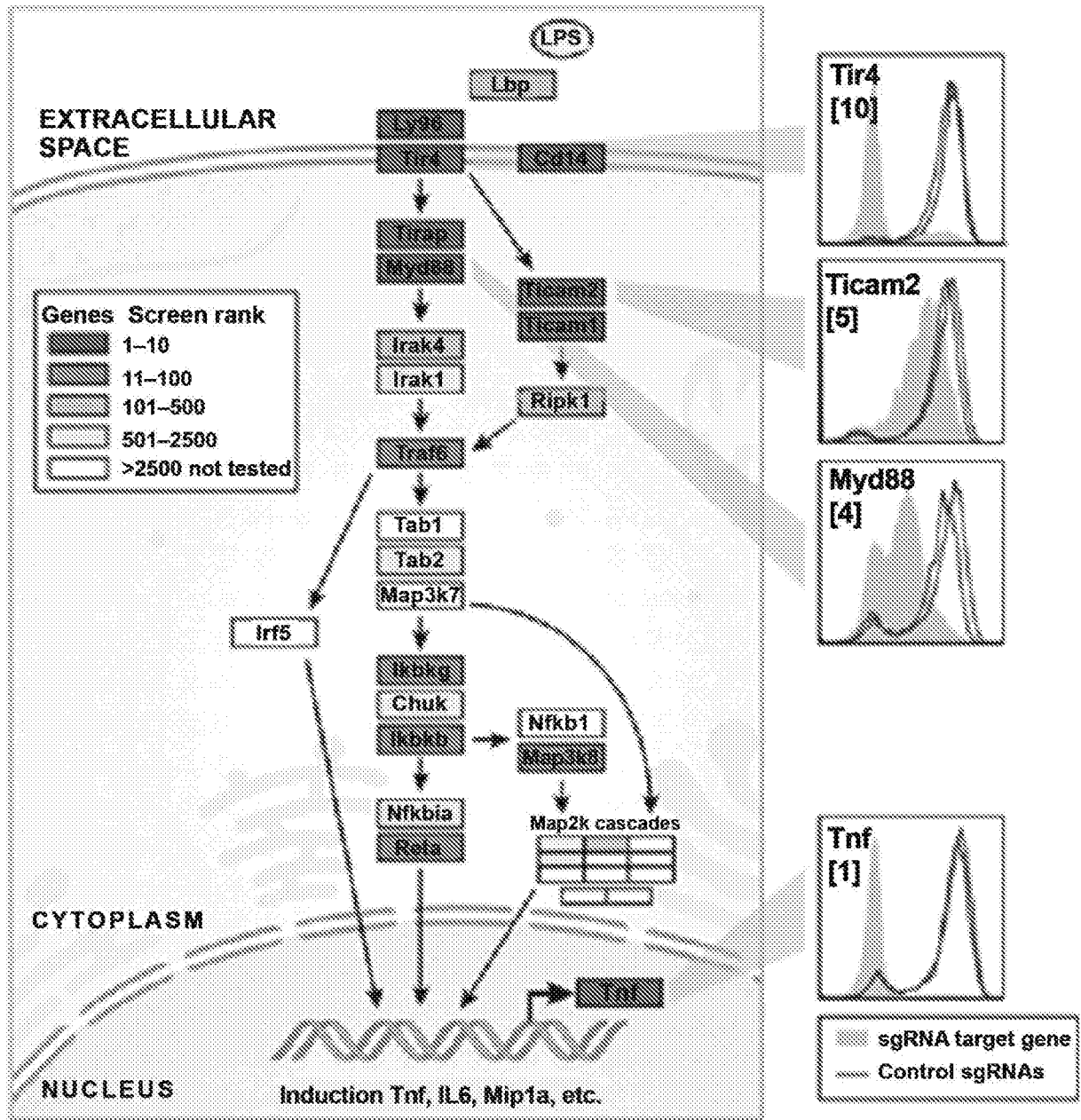


FIG. 47

FIG. 48A-B

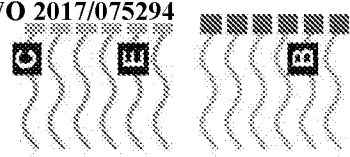
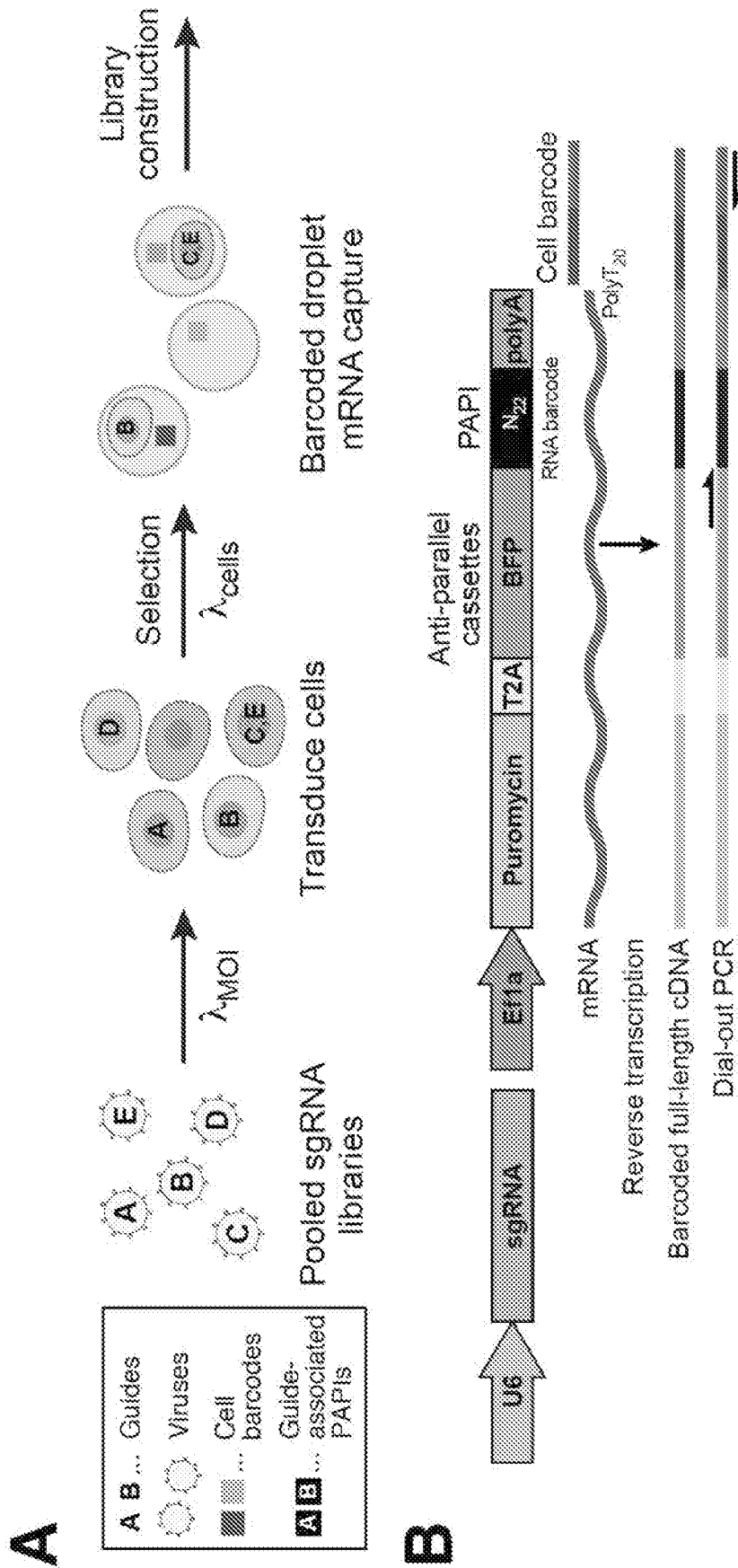


FIG. 48C-G

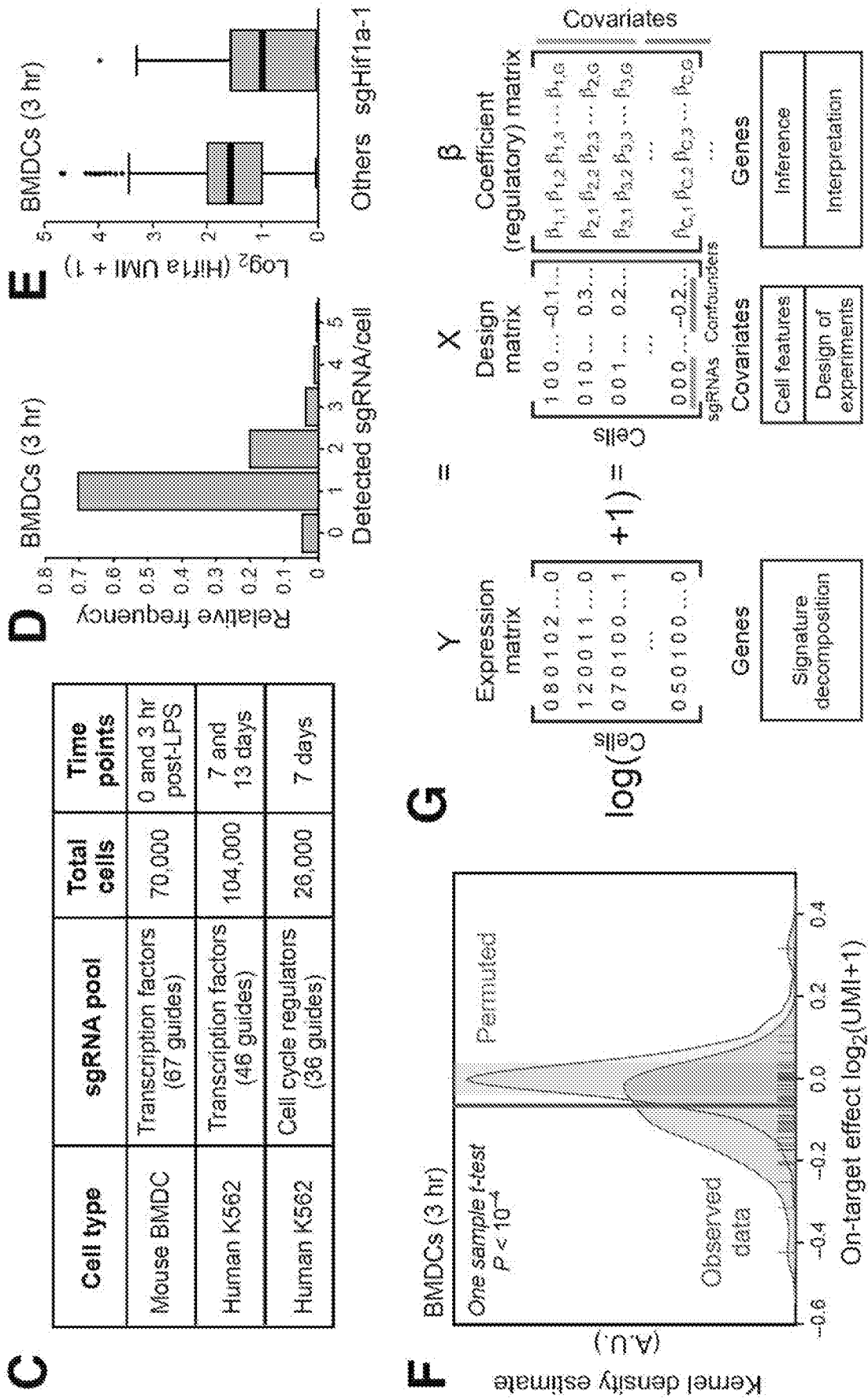


FIG. 49A-G

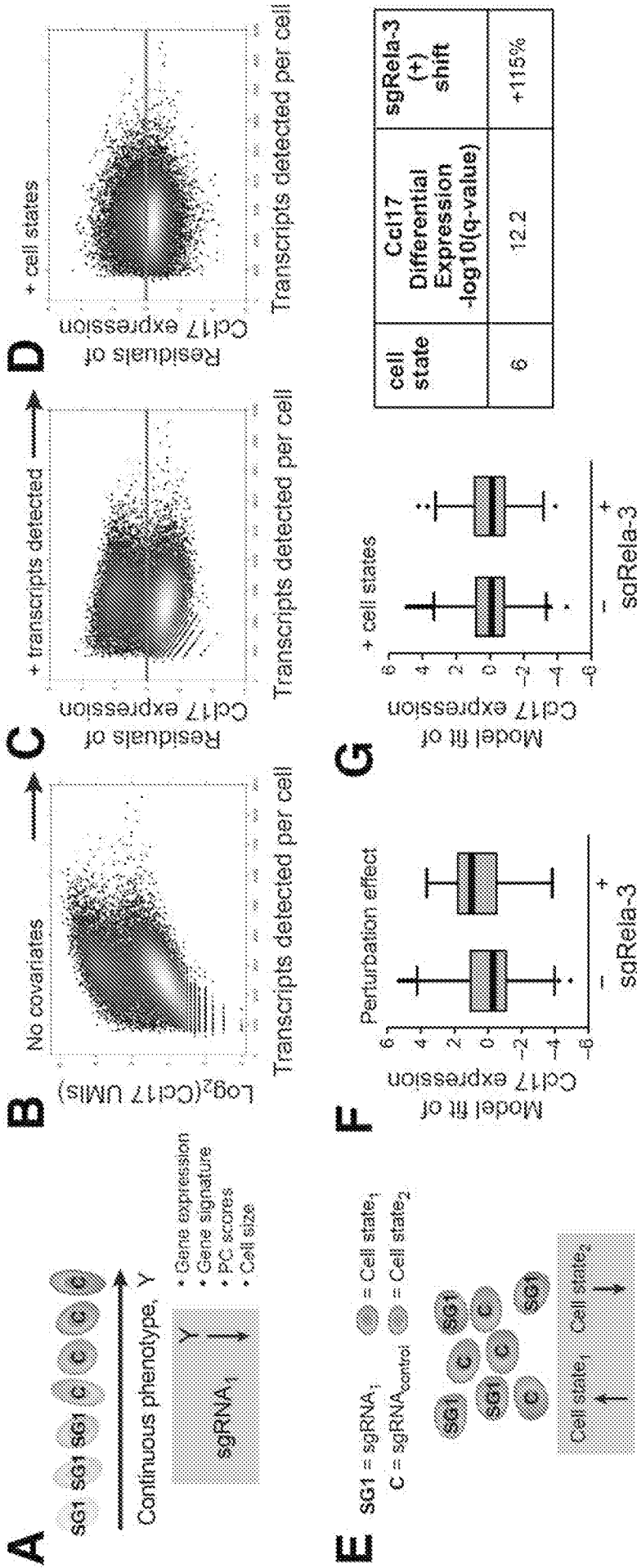


FIG. 49H-J

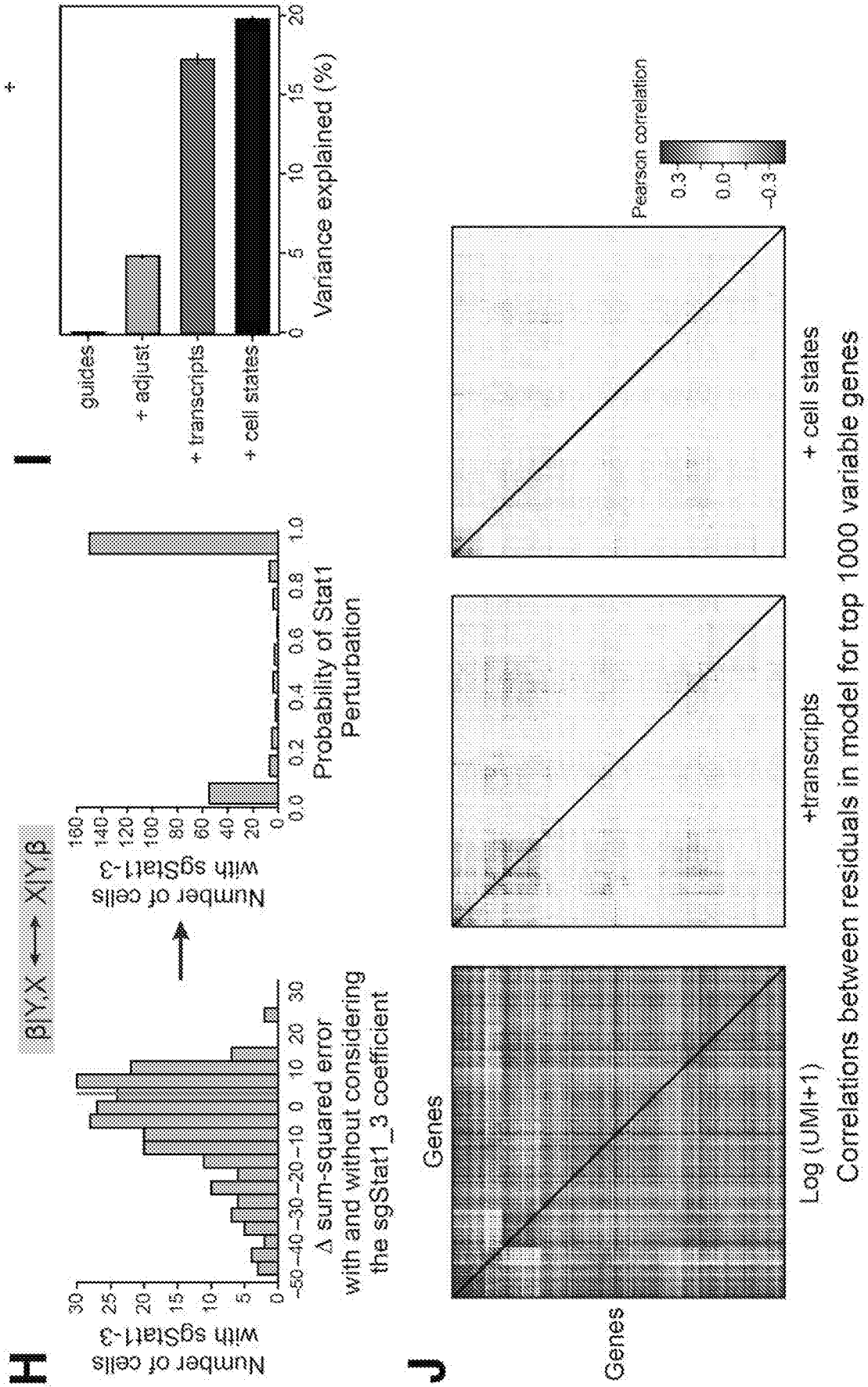


FIG. 50A-B

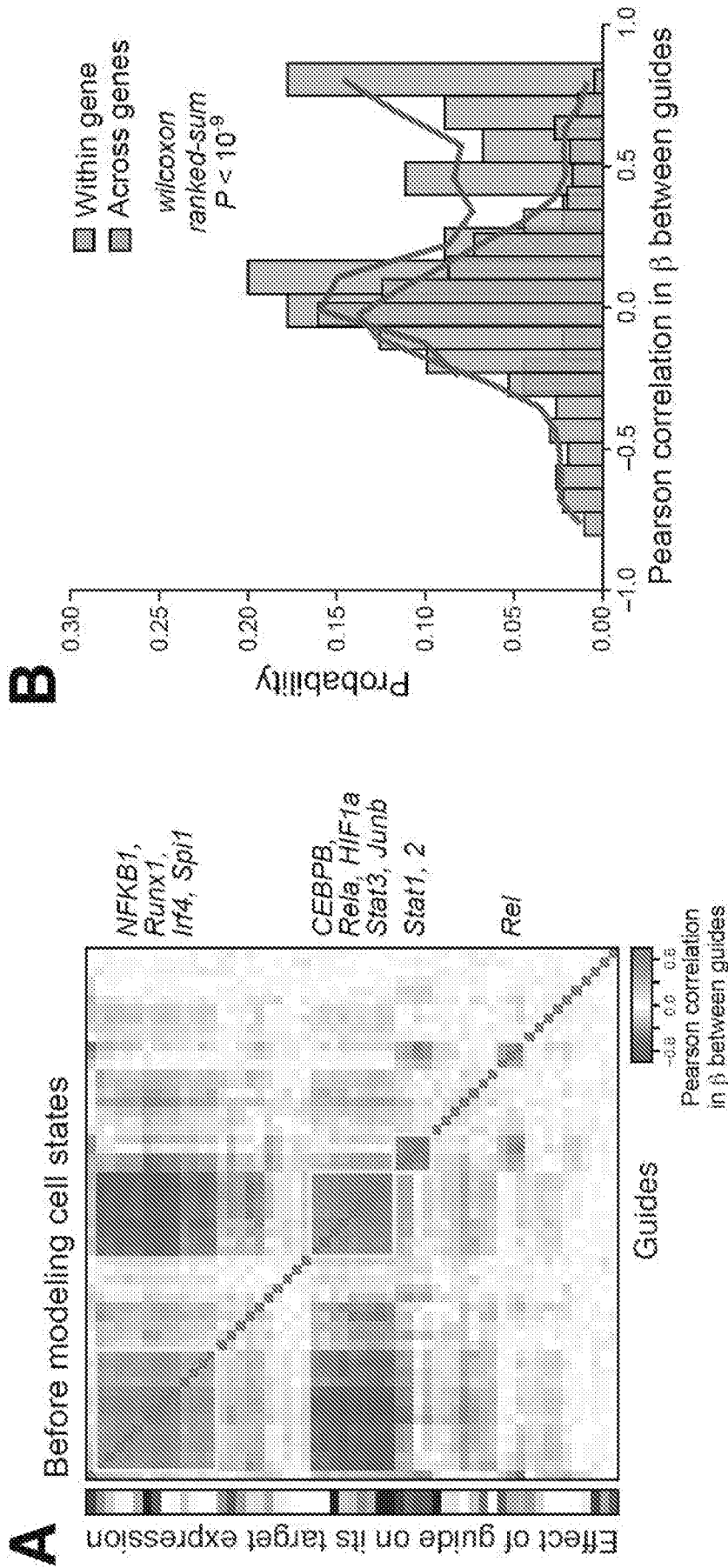


FIG. 50C-D

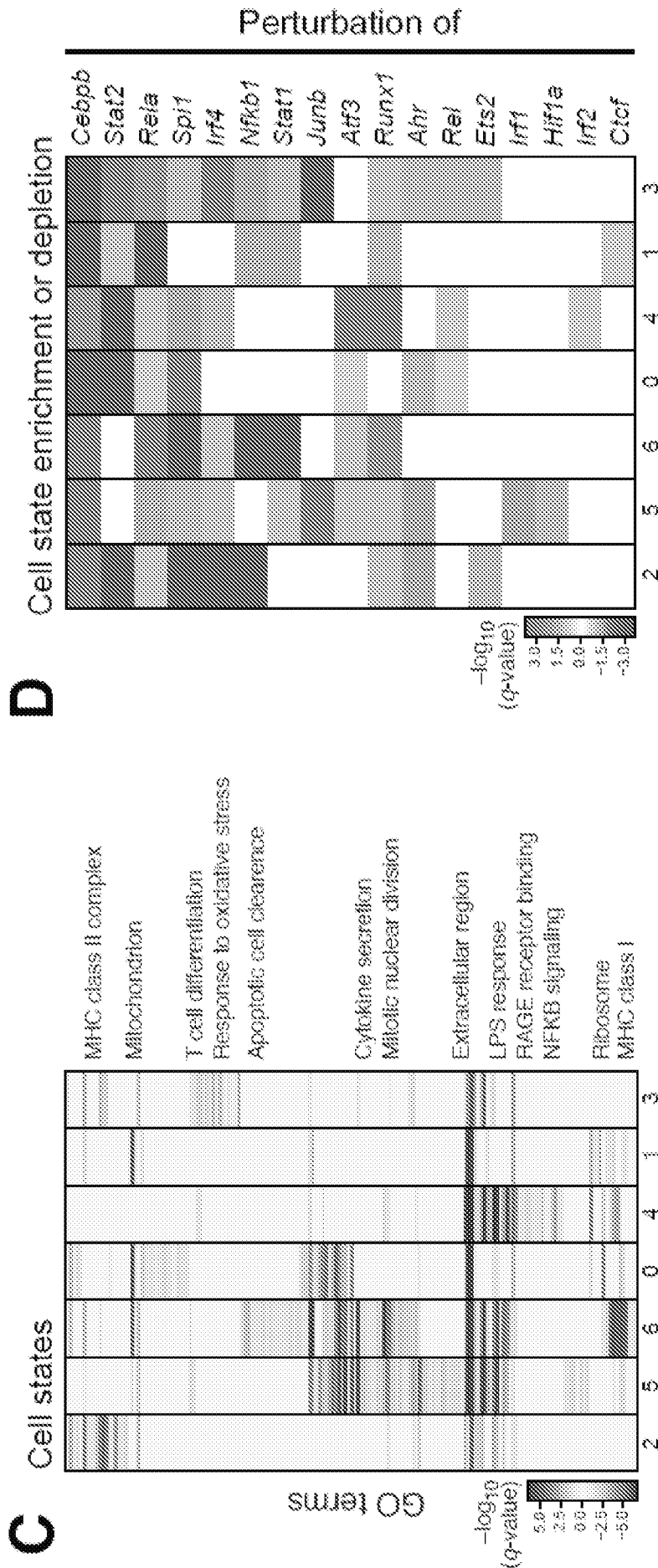


FIG. 50E-F

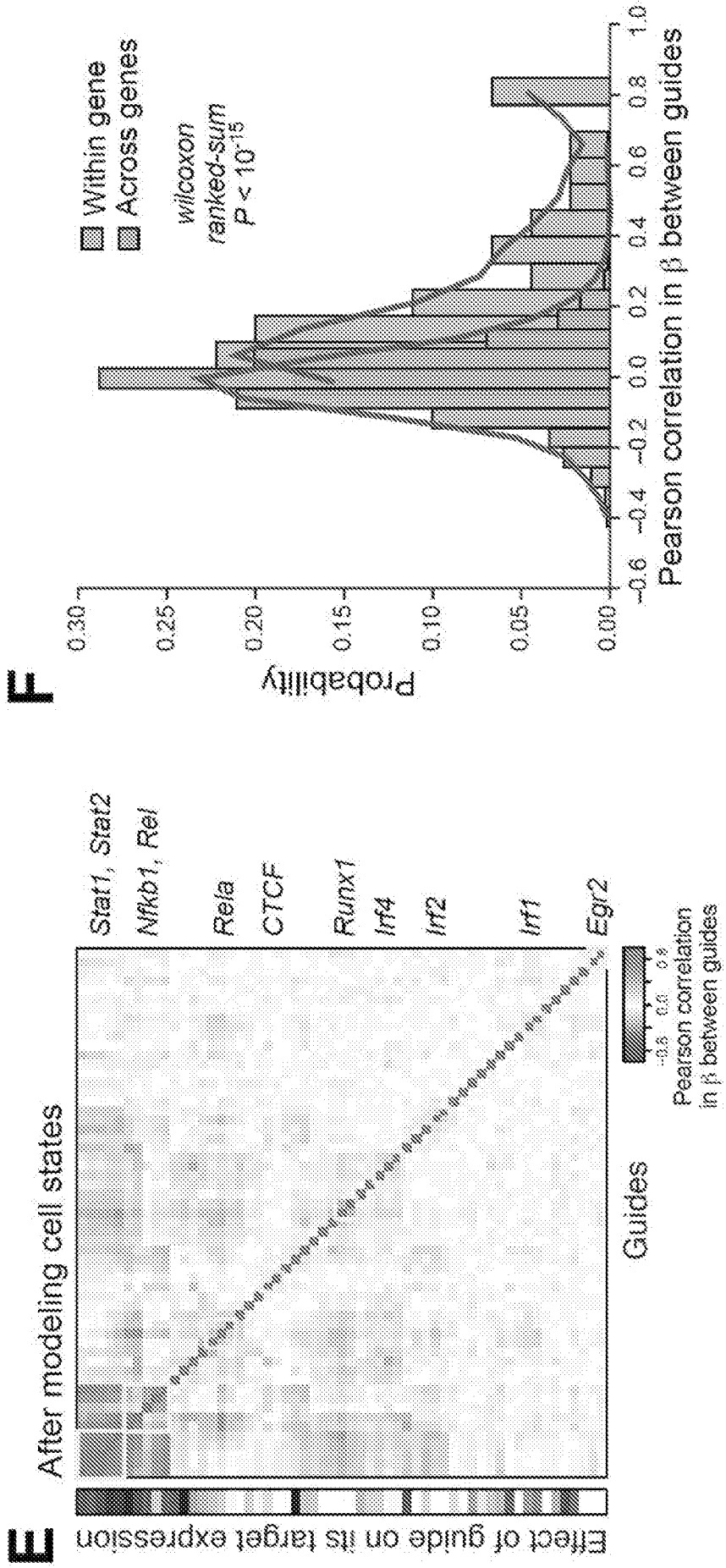


FIG. 51A

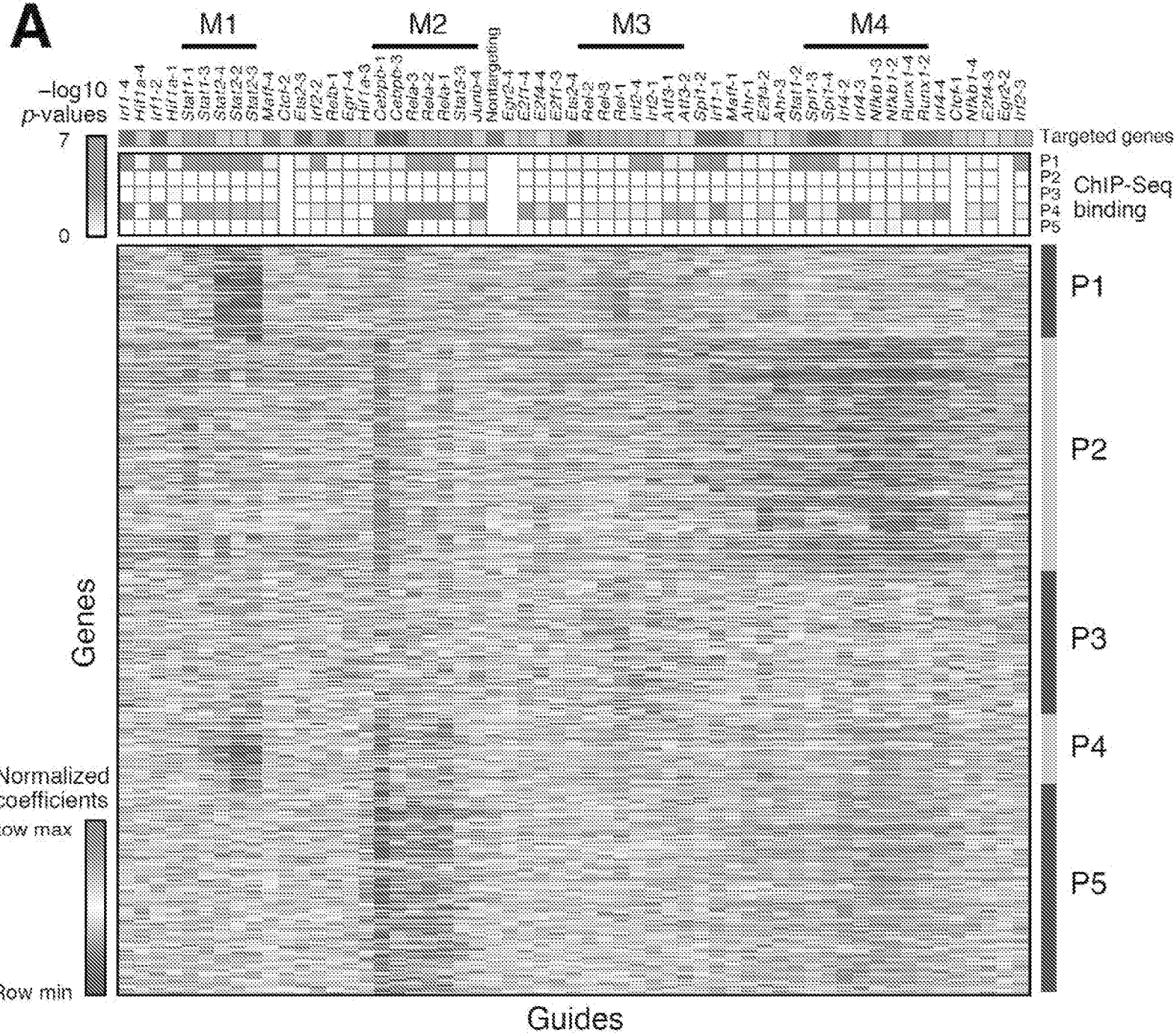


FIG. 51B

B

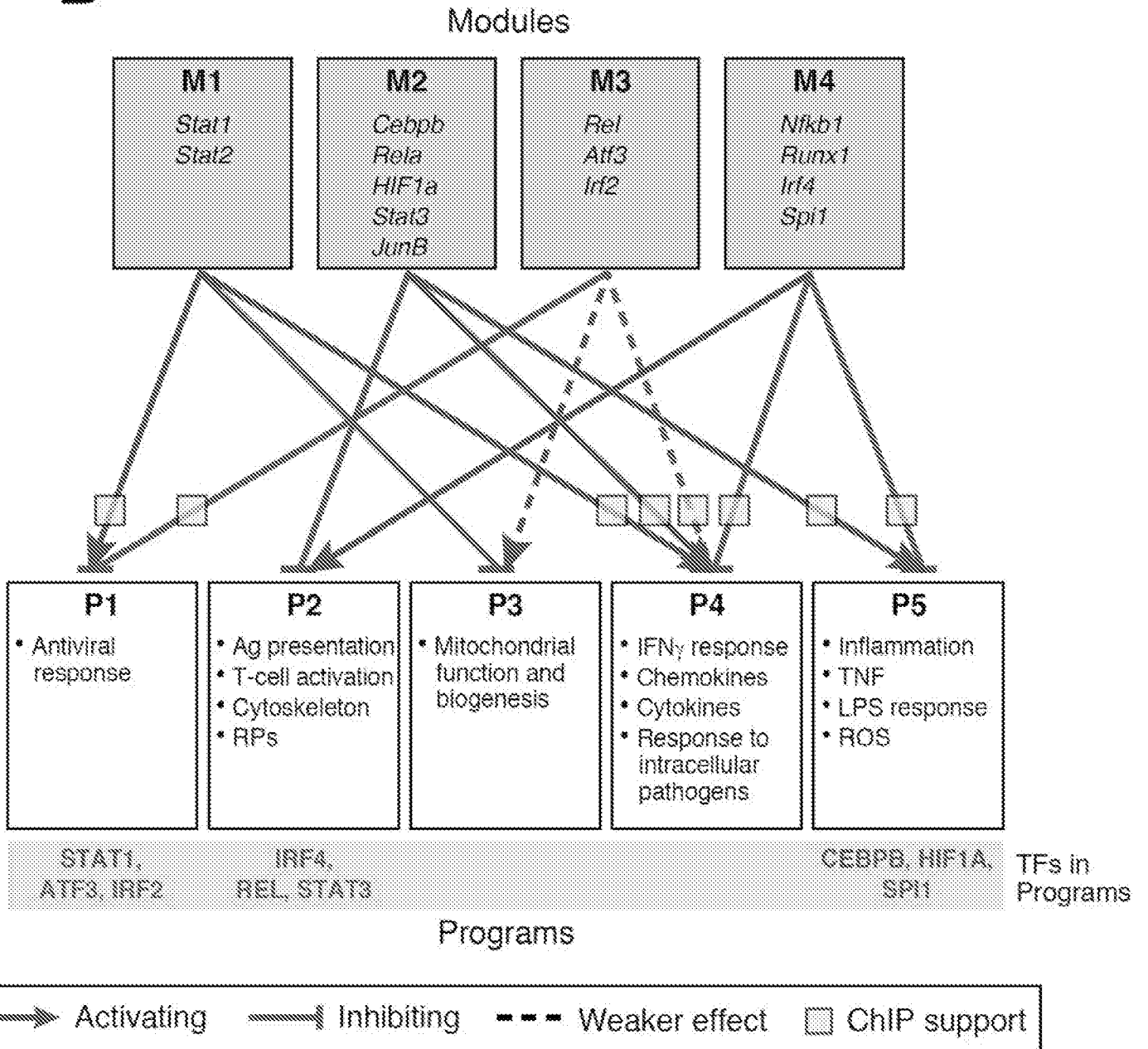


FIG. 51C

C

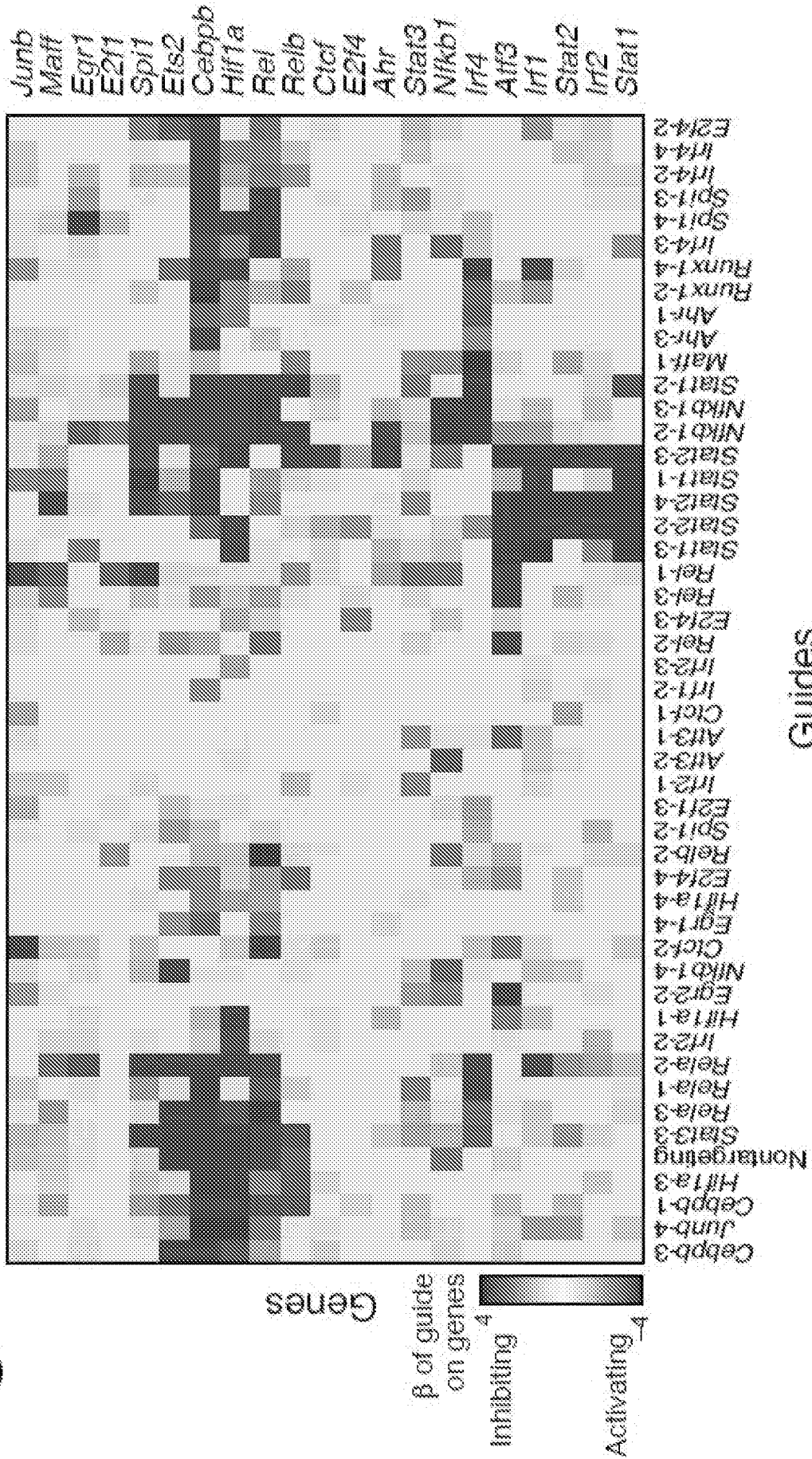
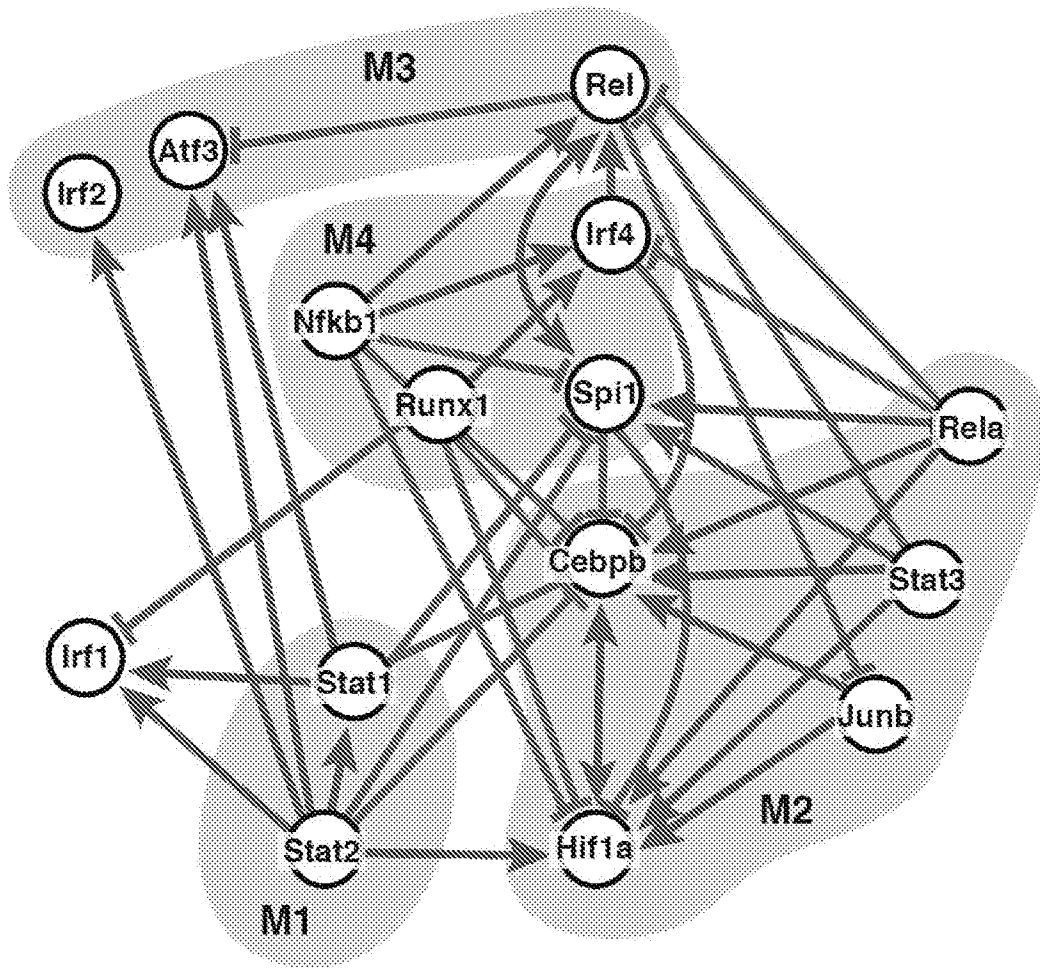


FIG. 51D

D



E

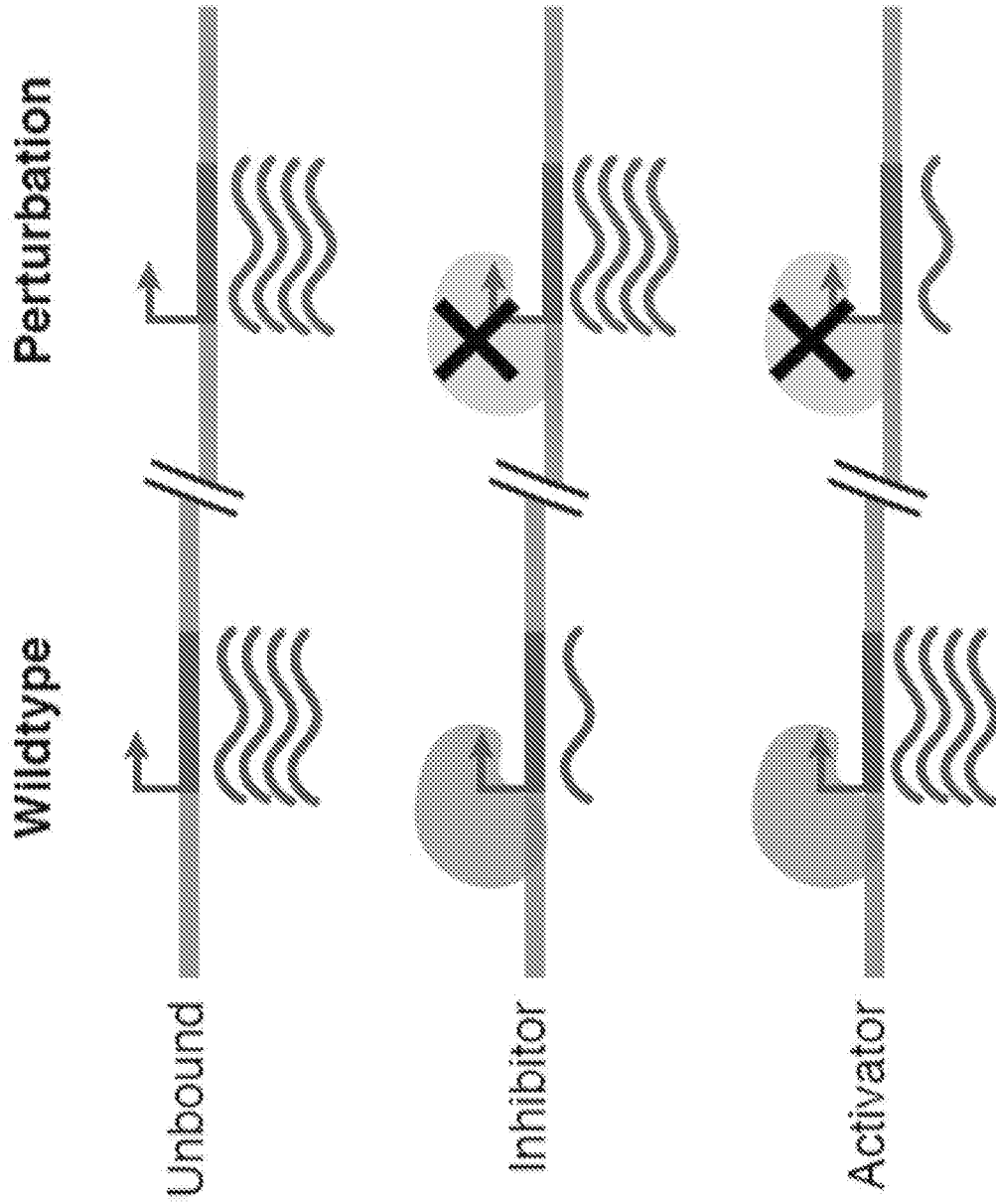


FIG. 51E

FIG. 51F

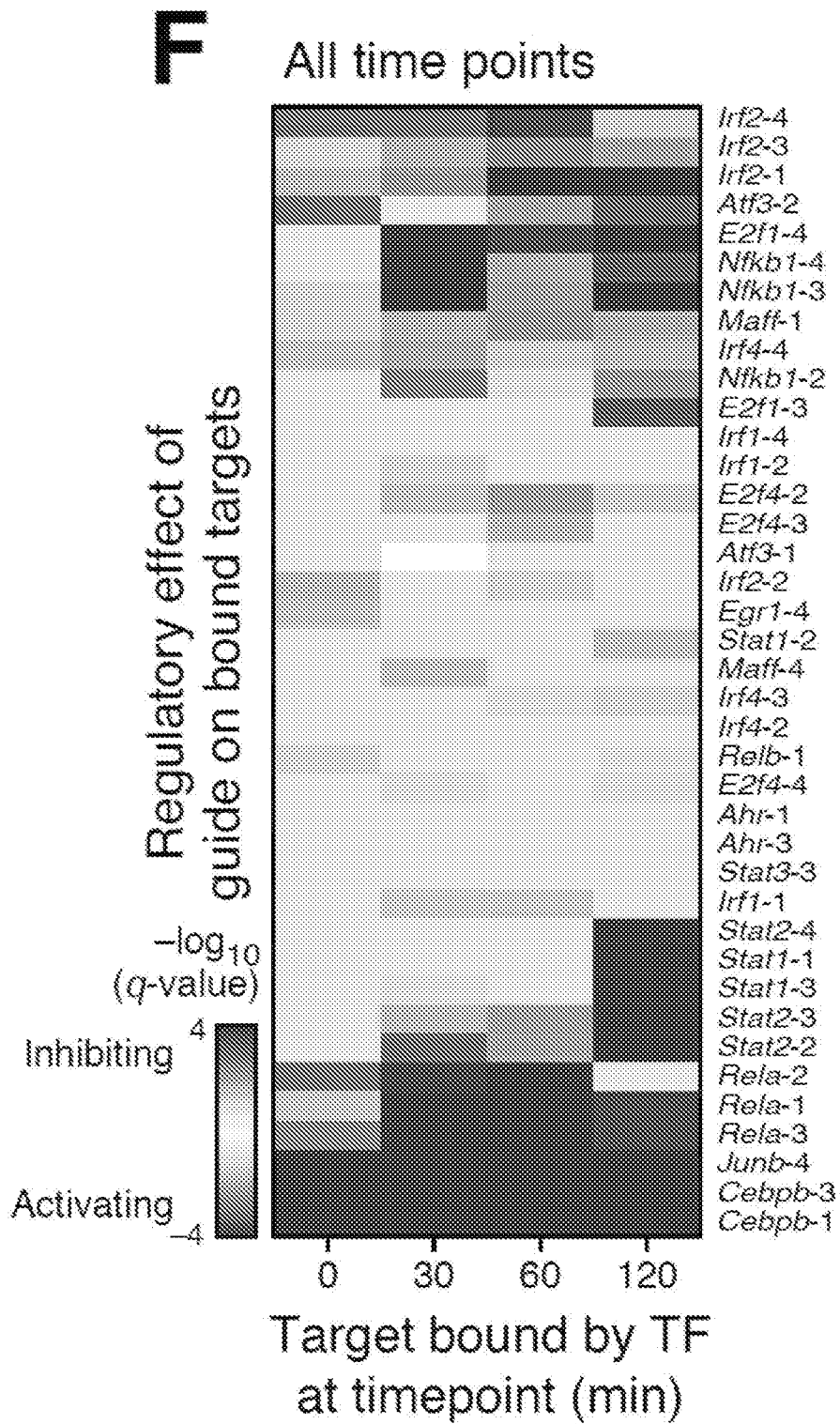


FIG. 51G

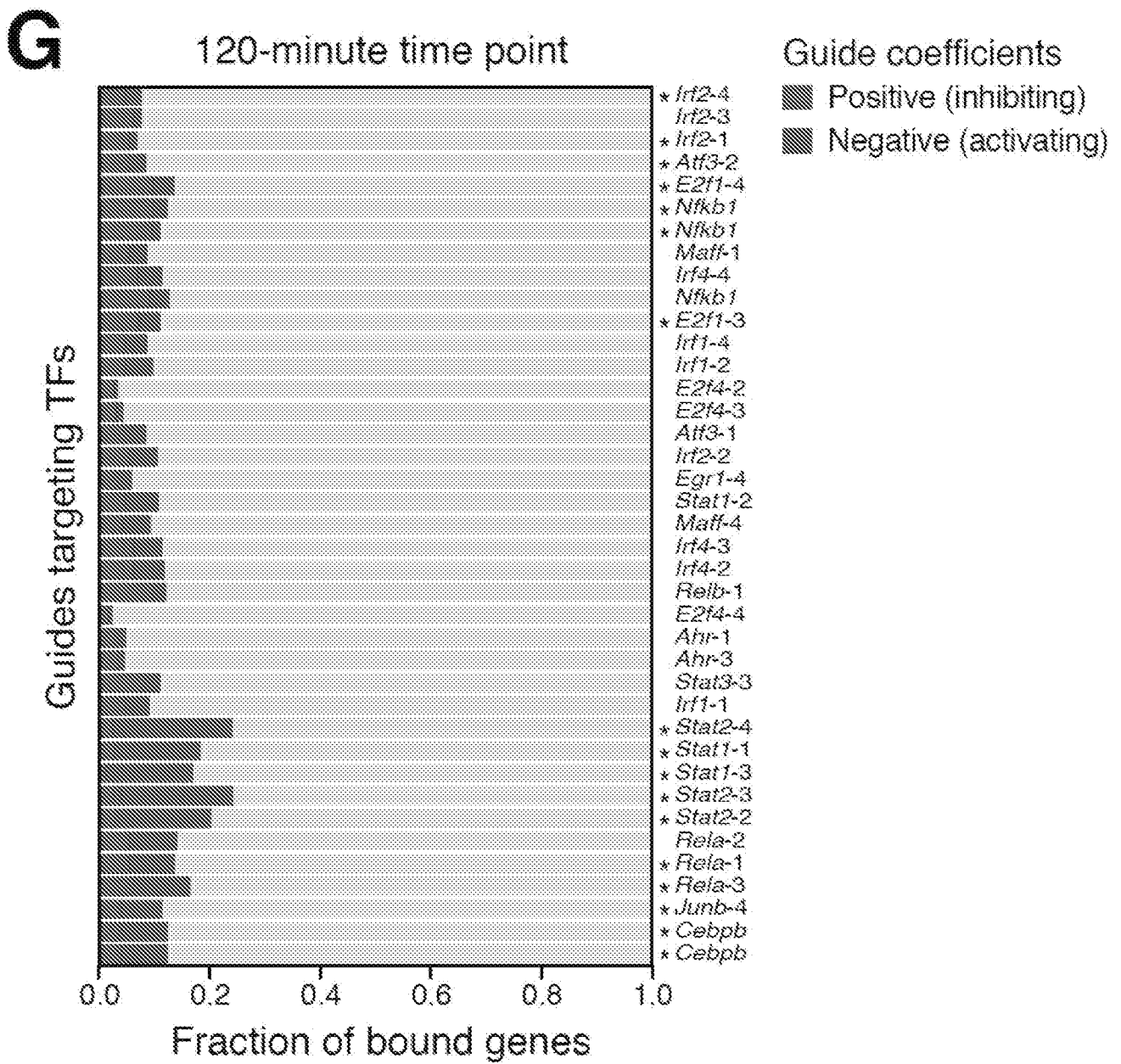


FIG. 52A

$$\begin{array}{c}
 \mathbf{A} \\
 \text{Expression matrix or} \\
 \text{cell state probabilities} \\
 \begin{bmatrix}
 0 & 8 & 0 & 1 & 0 & 2 & \dots & 0 \\
 1 & 2 & 0 & 0 & 1 & 1 & \dots & 0 \\
 0 & 7 & 0 & 1 & 0 & 0 & \dots & 1 \\
 \dots & & & & & & & \\
 1 & 0 & 0 & 0 & 1 & 0 & \dots & 0
 \end{bmatrix} \\
 \frac{C}{\Theta}
 \end{array}
 =
 \begin{array}{c}
 \mathbf{Y} \\
 \text{Perturbation matrix} \\
 \text{with interactions} \\
 \begin{bmatrix}
 1 & 0 & 0 \\
 0 & 1 & 0 \\
 1 & 1 & 1 \\
 \dots \\
 0 & 1 & 0
 \end{bmatrix} \\
 \frac{C}{\Theta}
 \end{array}
 =
 \begin{array}{c}
 \mathbf{X} \\
 \text{Regulatory} \\
 \text{matrix} \\
 \begin{bmatrix}
 \beta_{A,1} & \beta_{A,2} & \dots & \beta_{A,G} \\
 \beta_{B,1} & \beta_{B,2} & \dots & \beta_{B,G} \\
 \beta_{AB,1} & \beta_{AB,2} & \dots & \beta_{AB,G}
 \end{bmatrix} \\
 \text{Perturbations}
 \end{array}
 \end{array}$$

$\mathbf{A} = \mathbf{A}\beta_{\mathbf{A}} + \mathbf{B}\beta_{\mathbf{B}} + \mathbf{AB}\beta_{\mathbf{AB}}$

FIG. 52B

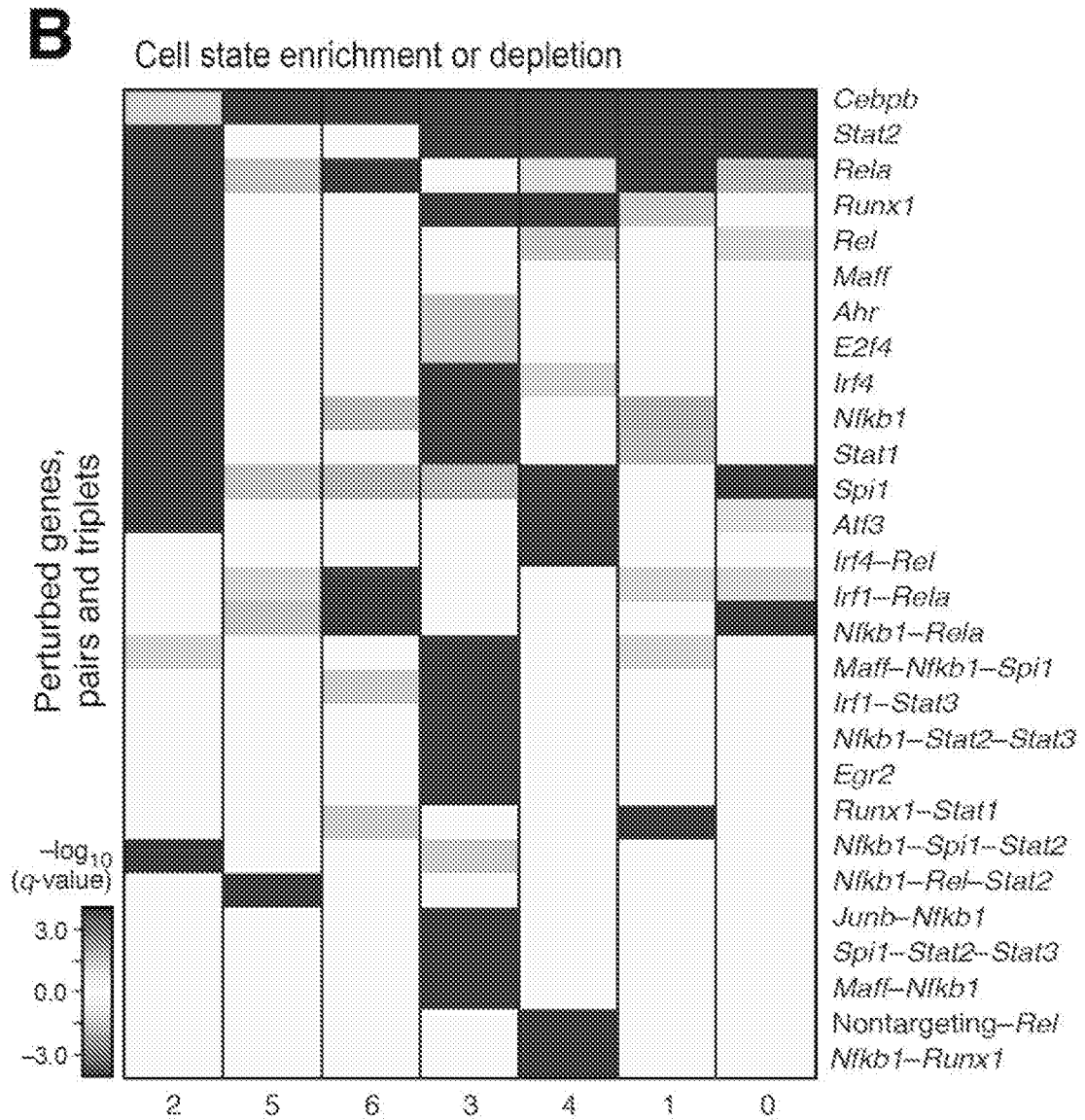


FIG. 52C

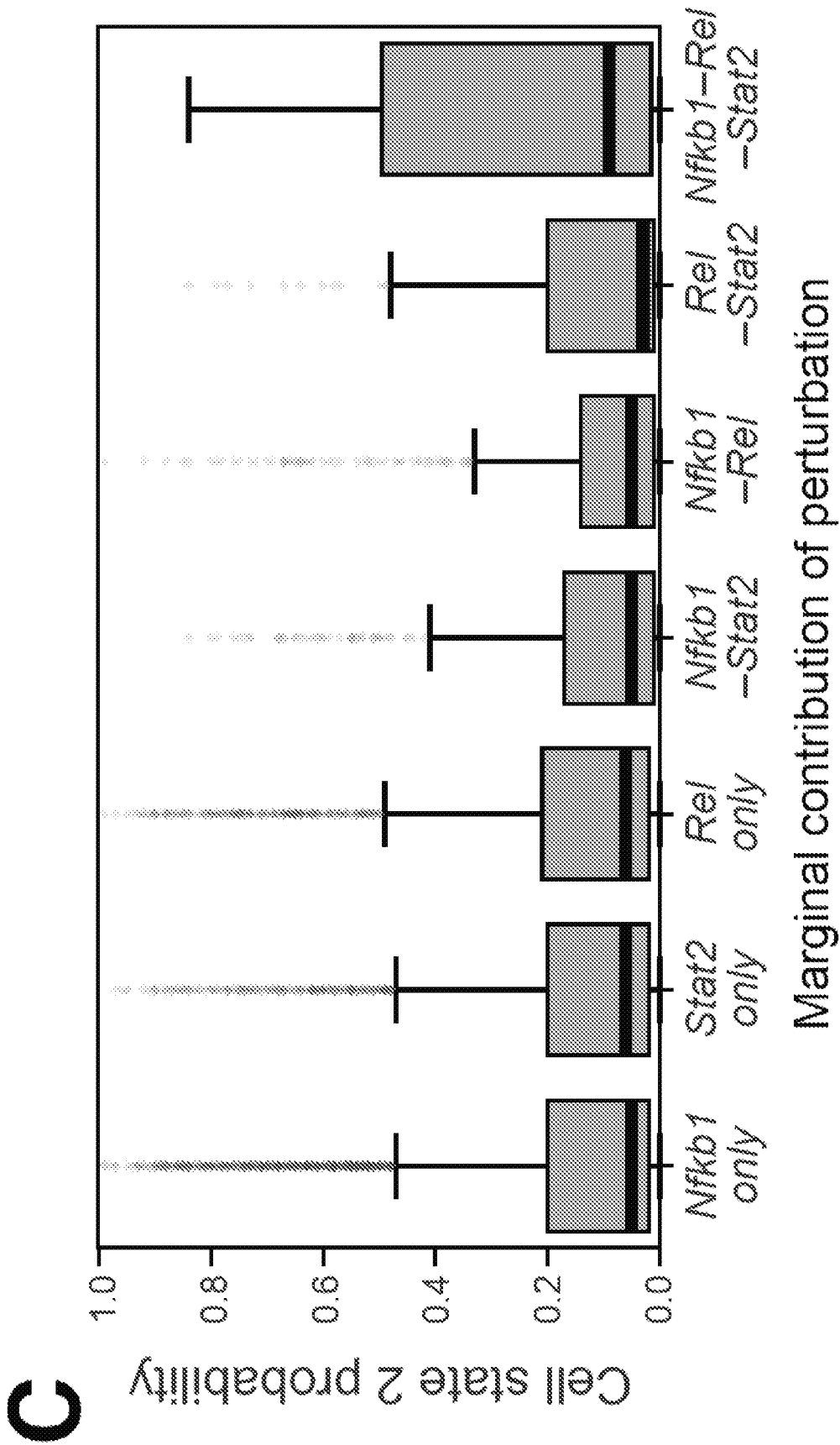
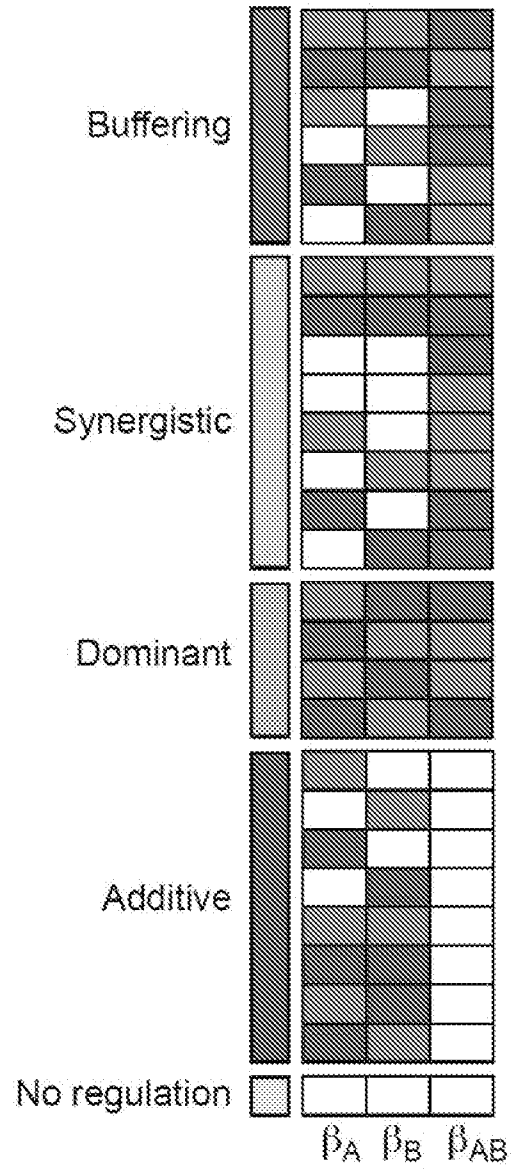


FIG. 52D

D

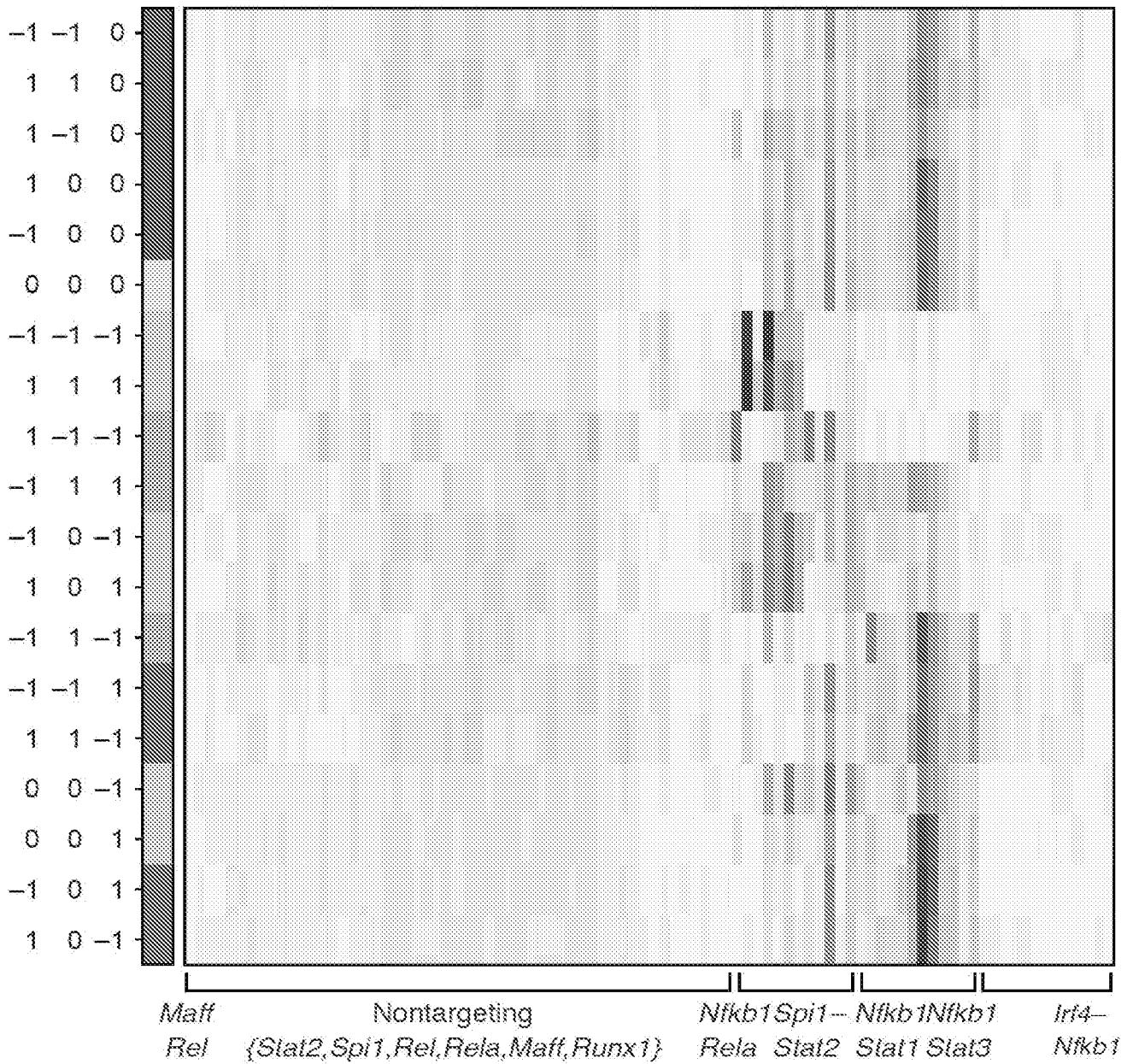


$3^3 = 27$ possibilities,
4 major sets

-
- 0
- +

FIG. 52E

E



Buffering Additive
 Synergistic No regulation
 Dominant

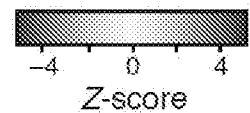


FIG. 52F

F

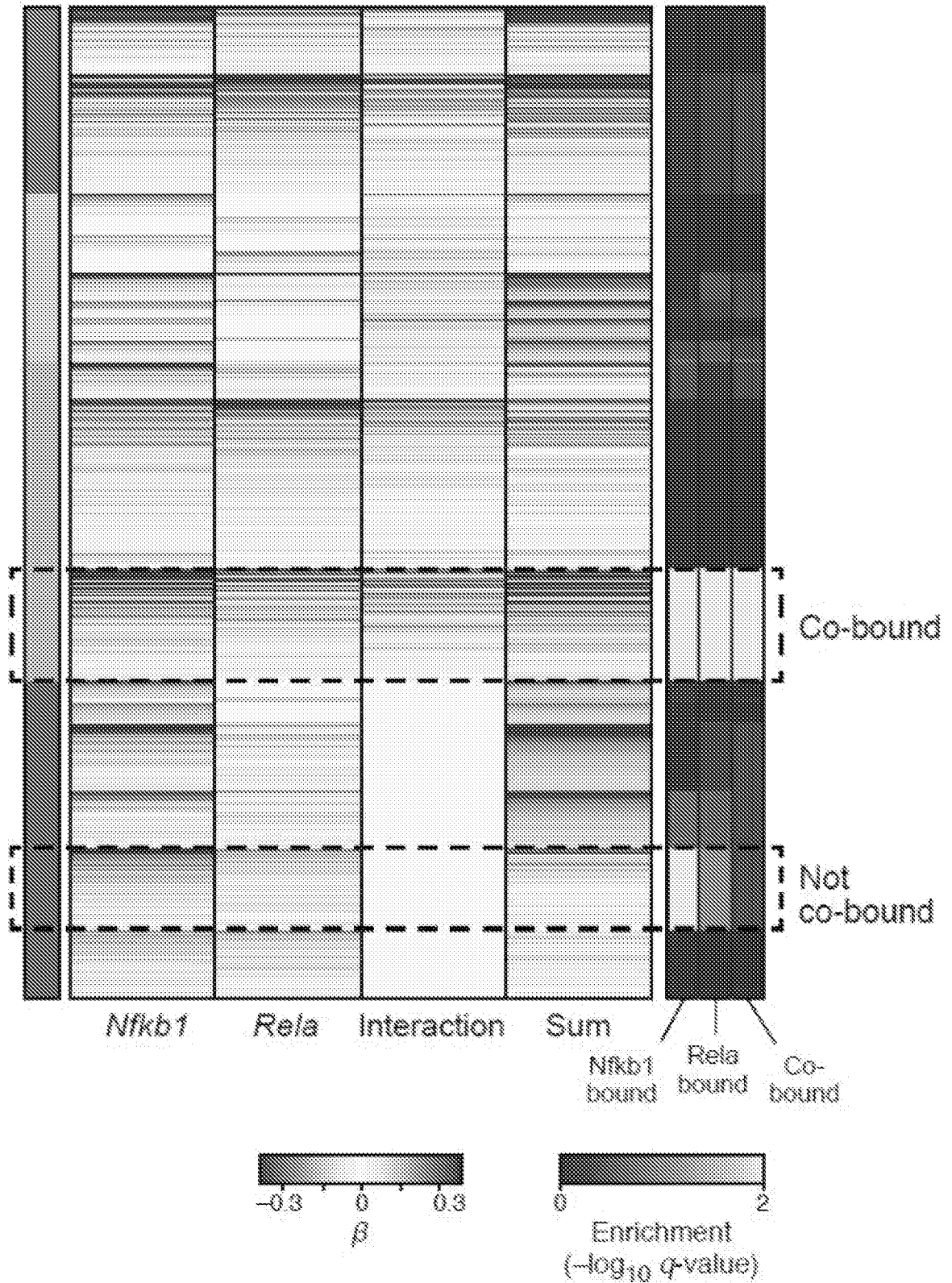


FIG. 53A-B

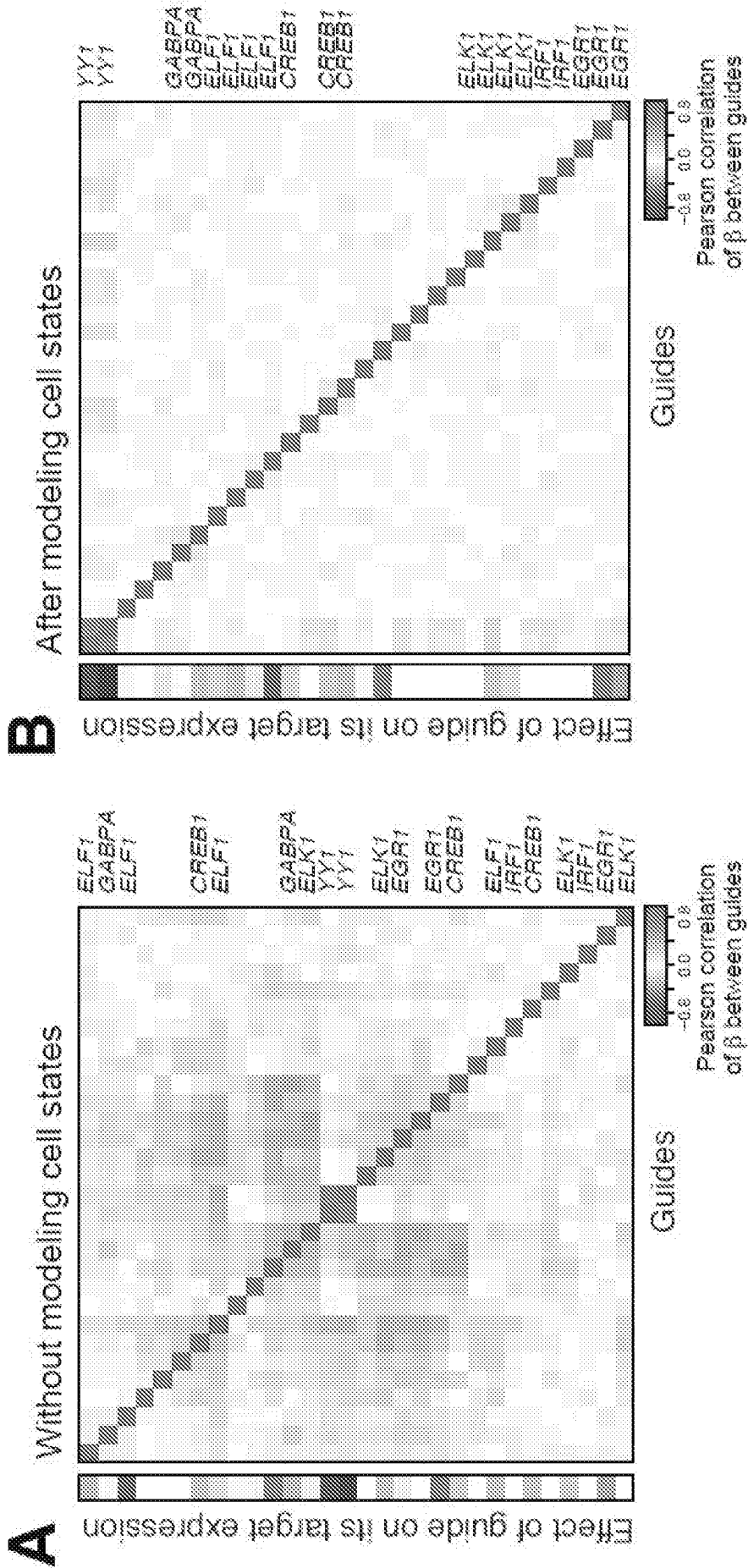
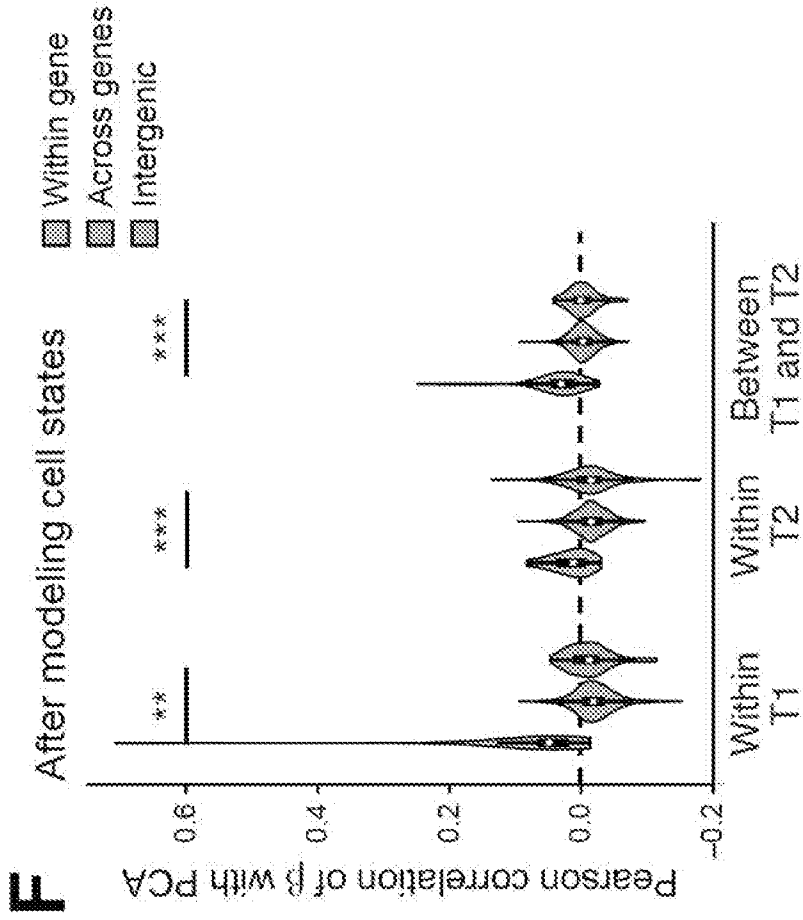
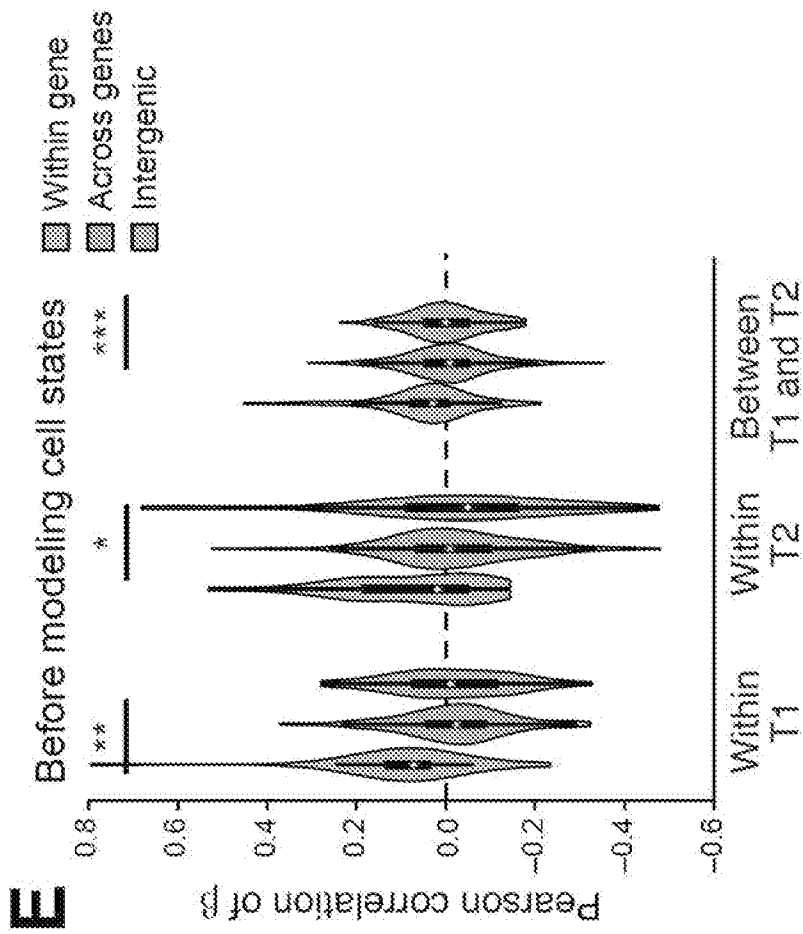


FIG. 53E-F



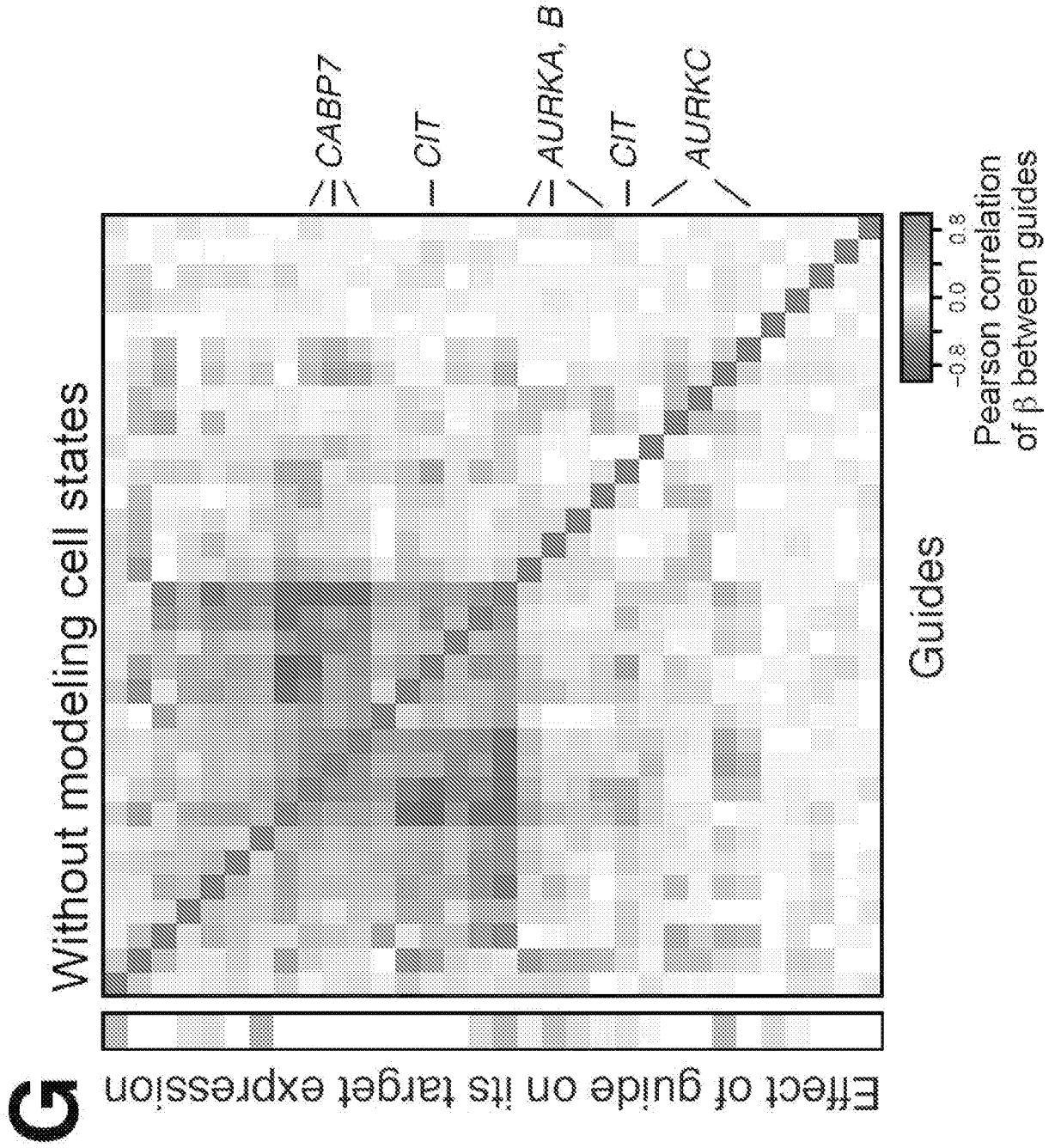
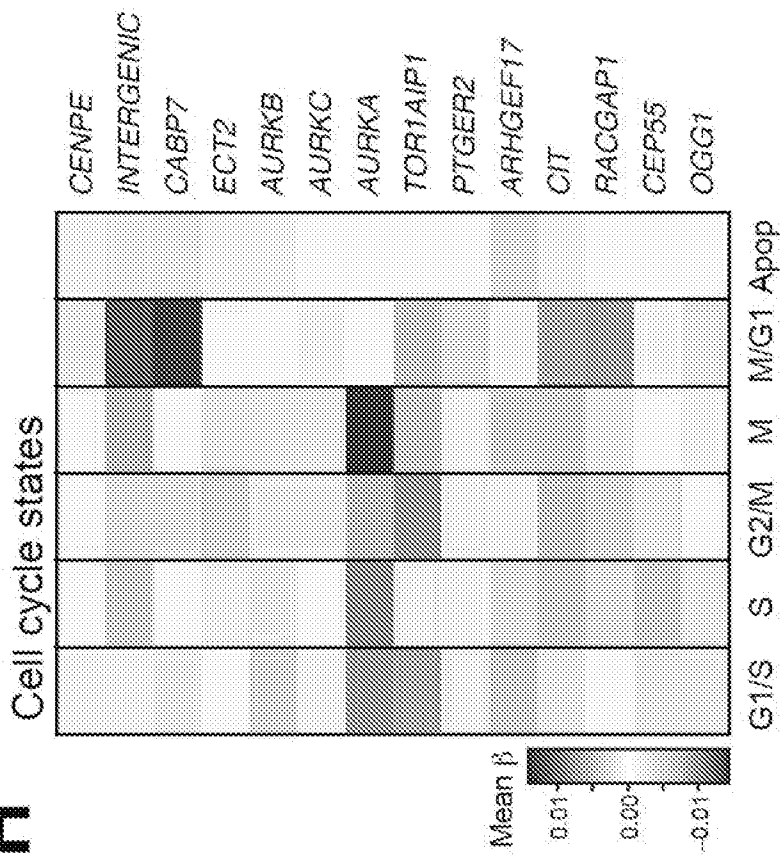


FIG. 53G

FIG. 53H-I

H



I

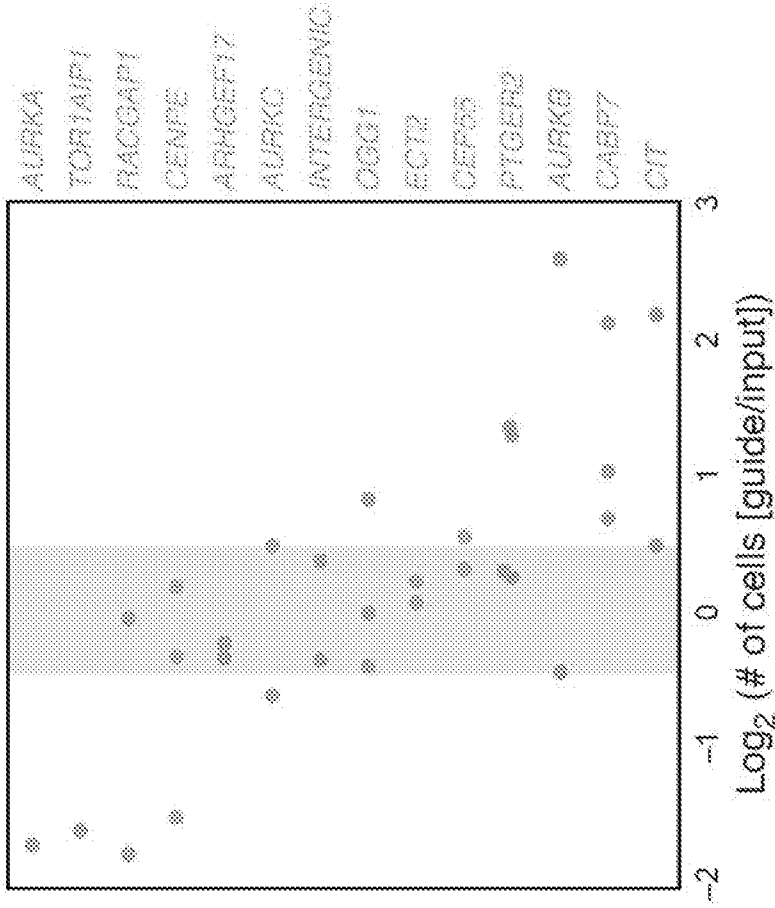


FIG. 54A-B

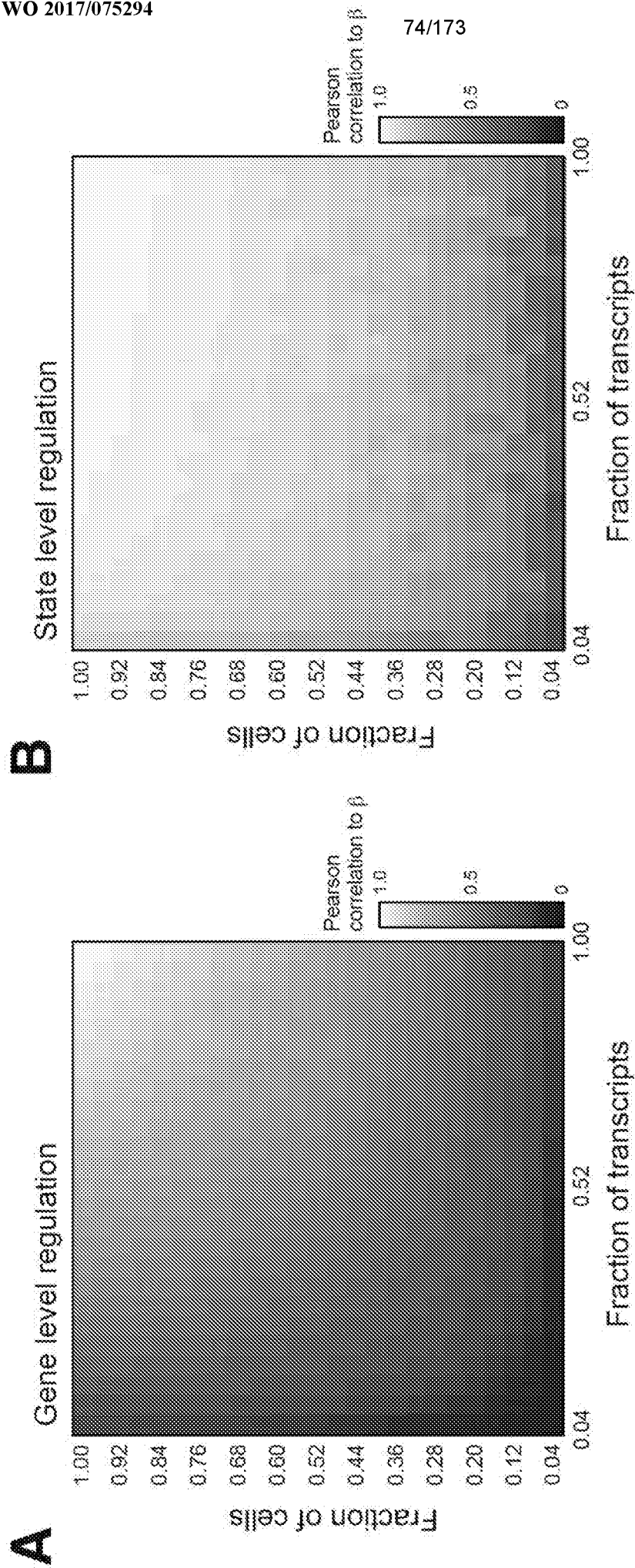


FIG. 54C

C

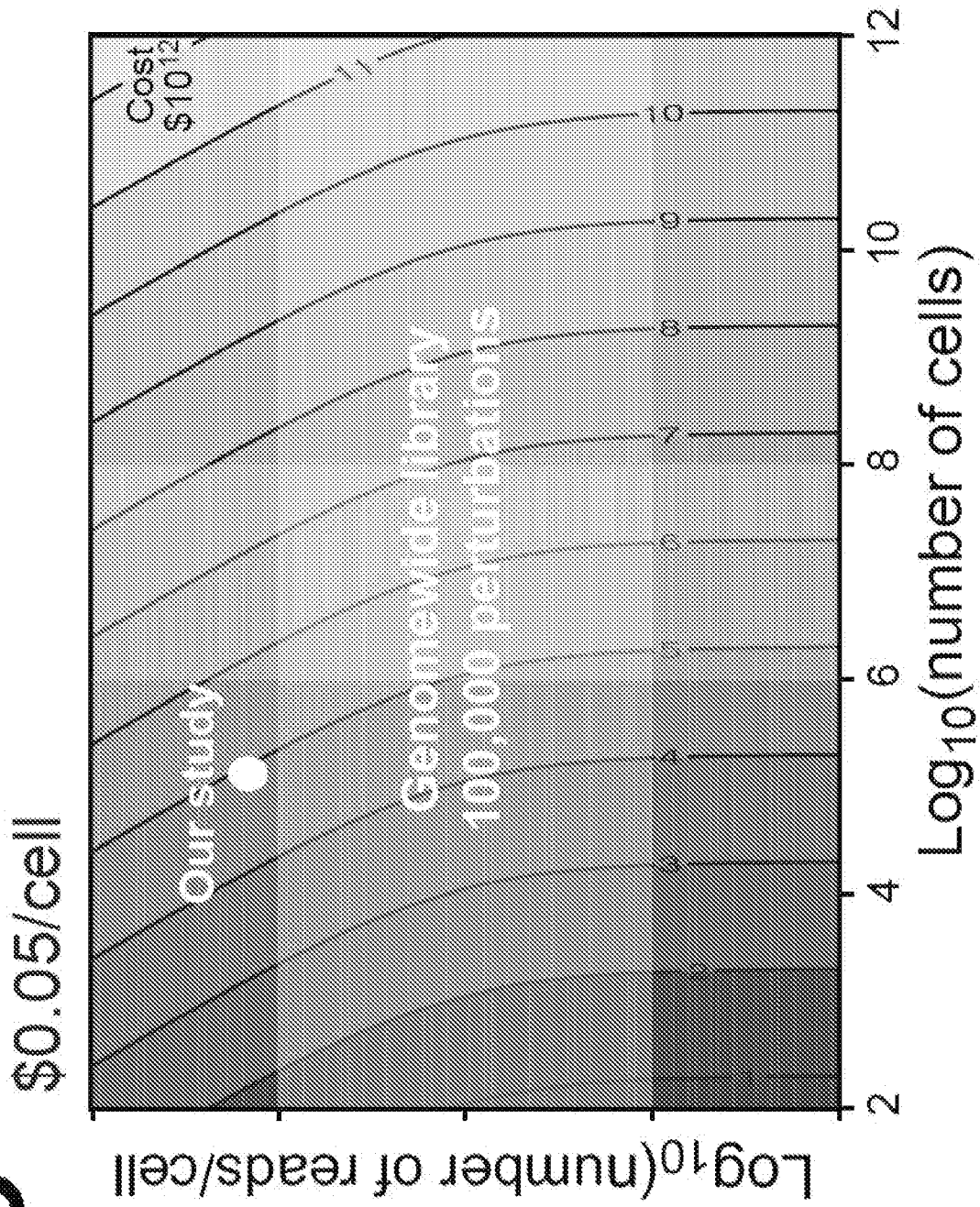


FIG. 54D

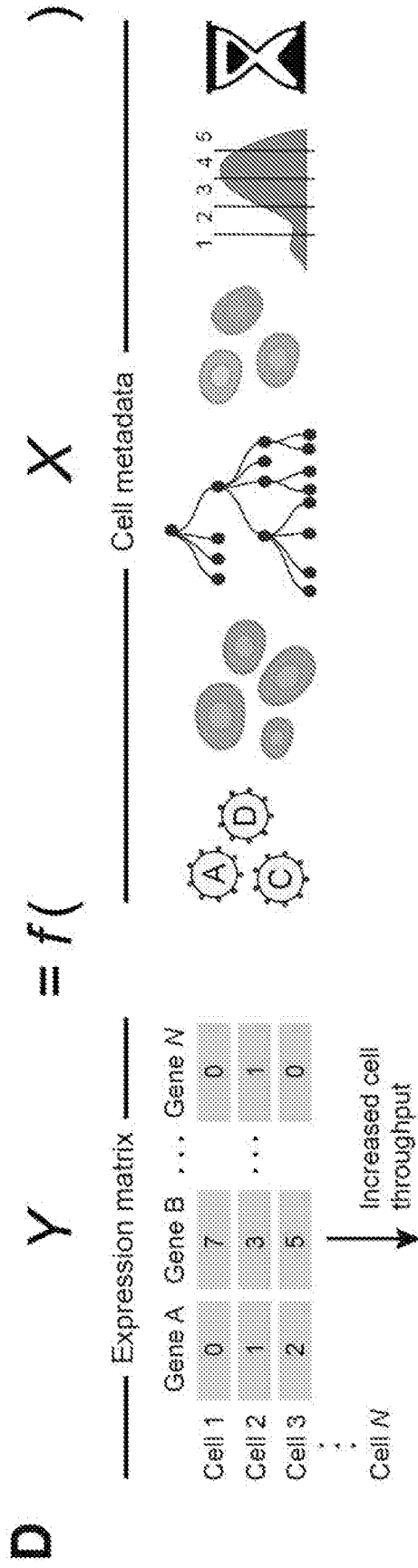


FIG. 55A-C

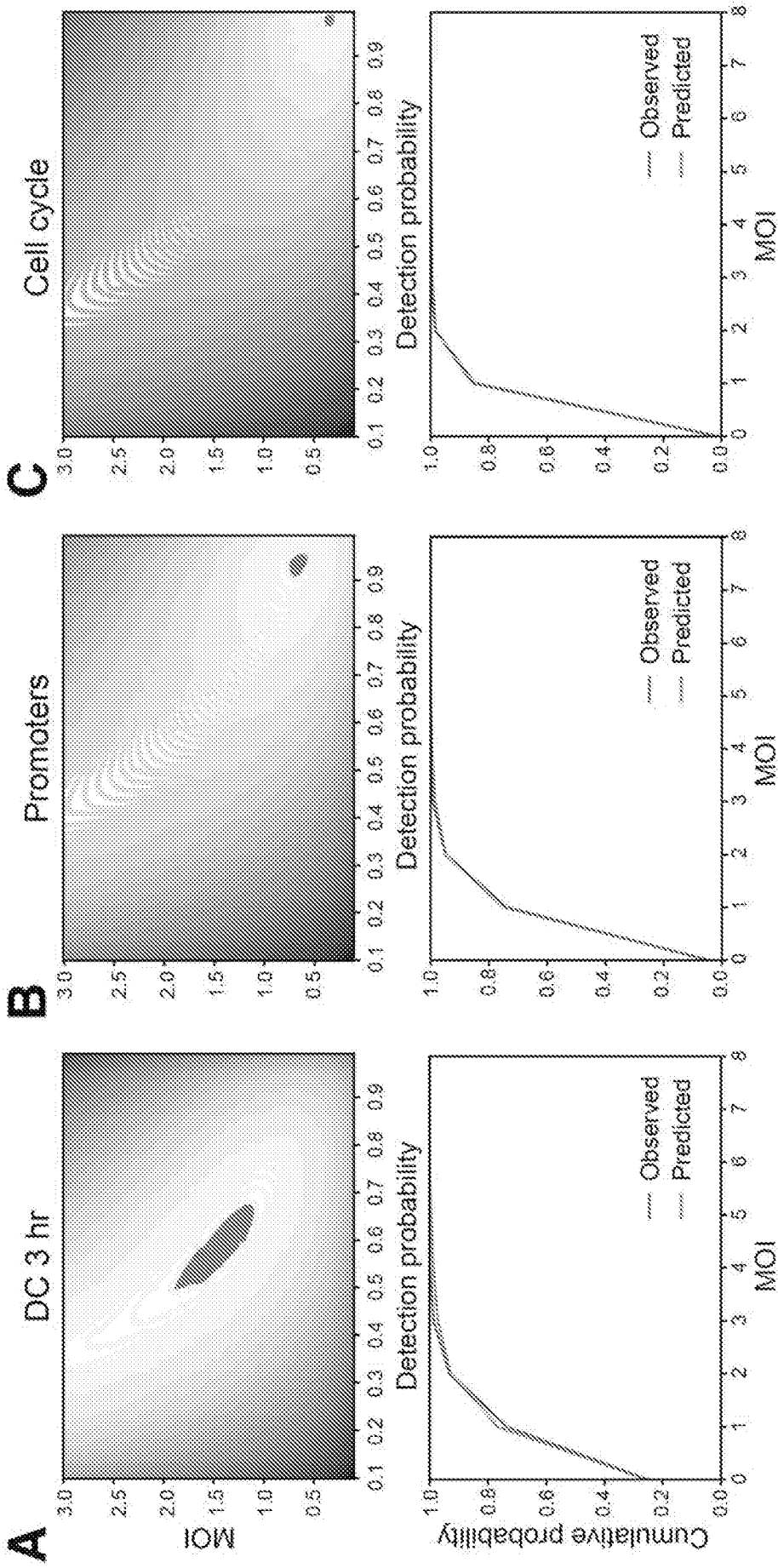
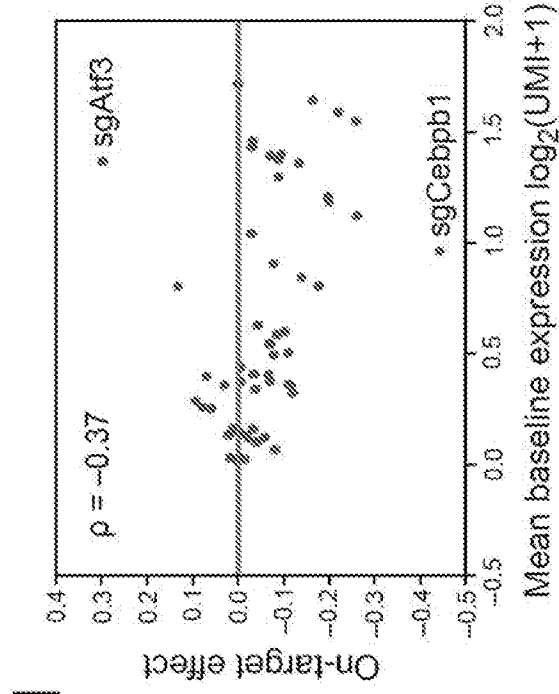
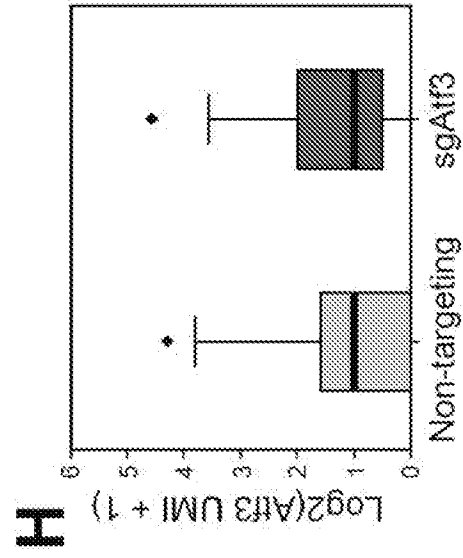
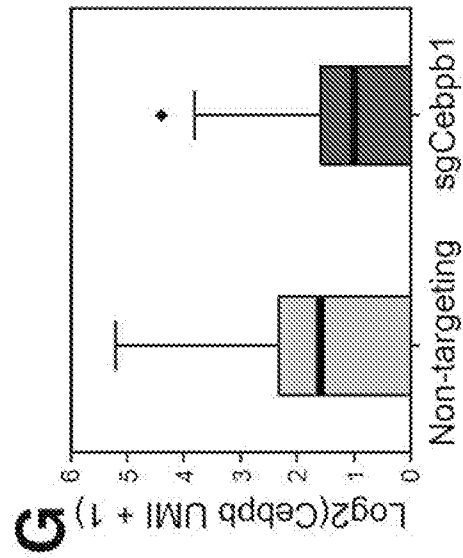
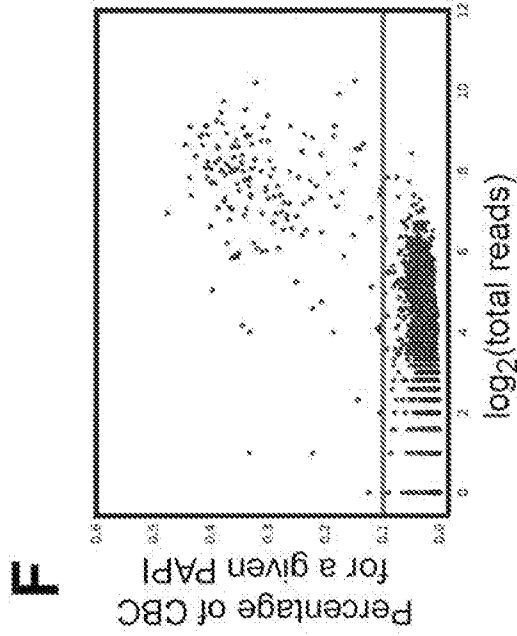
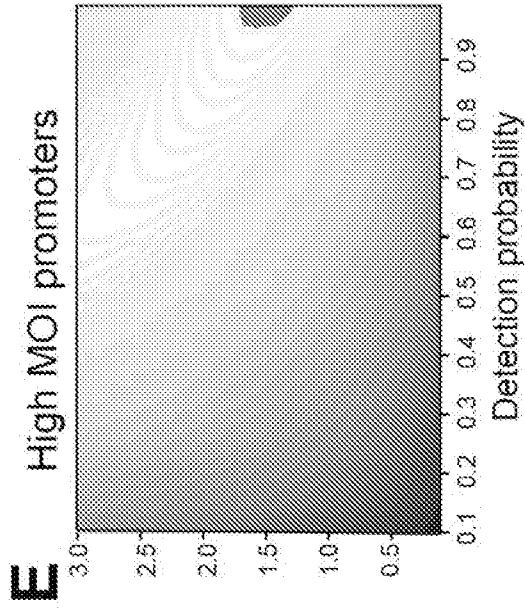
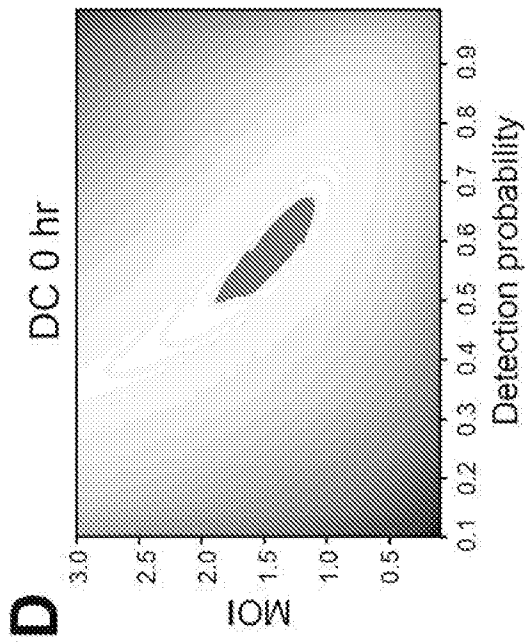
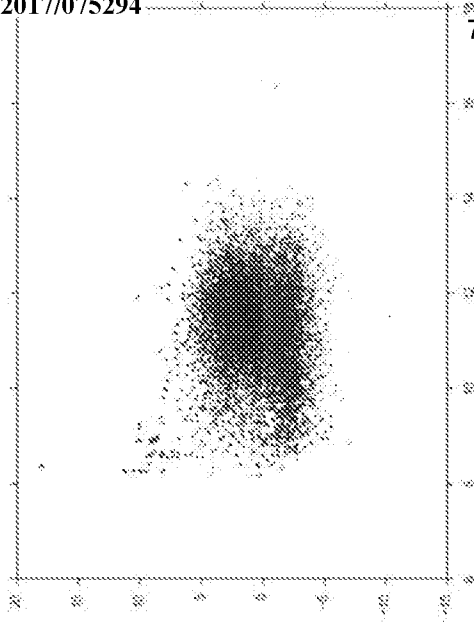


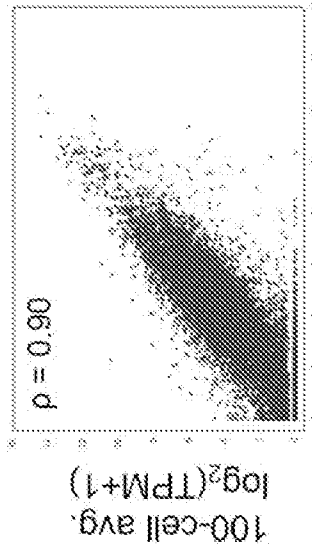
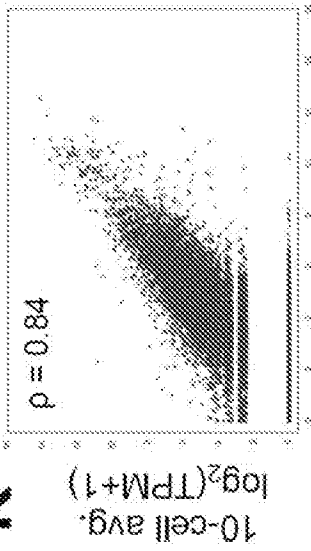
FIG. 55D-I



L Difference between population and single cell avg. expression



K



J

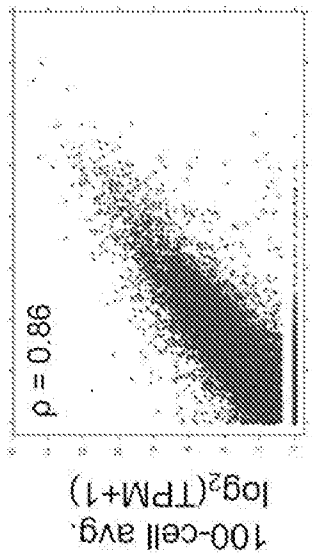
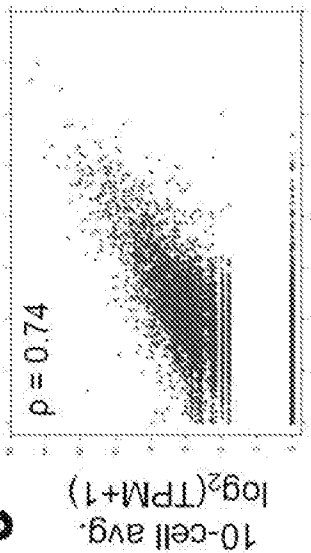


FIG. 55J-L

FIG. 56A-C

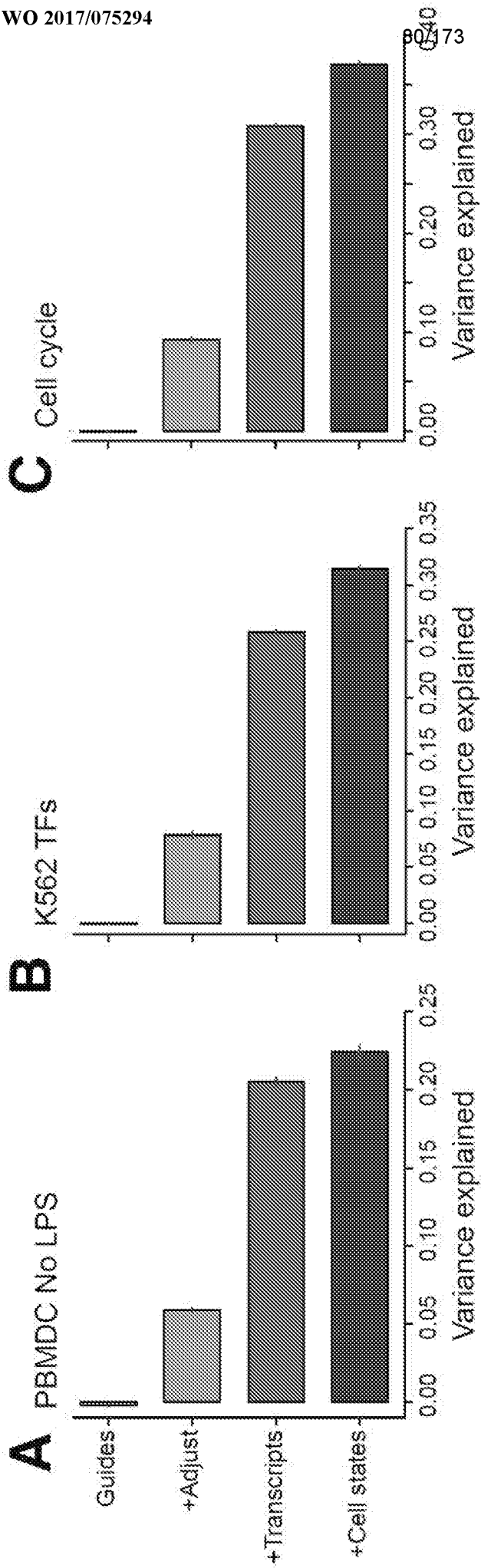


FIG. 56D-G

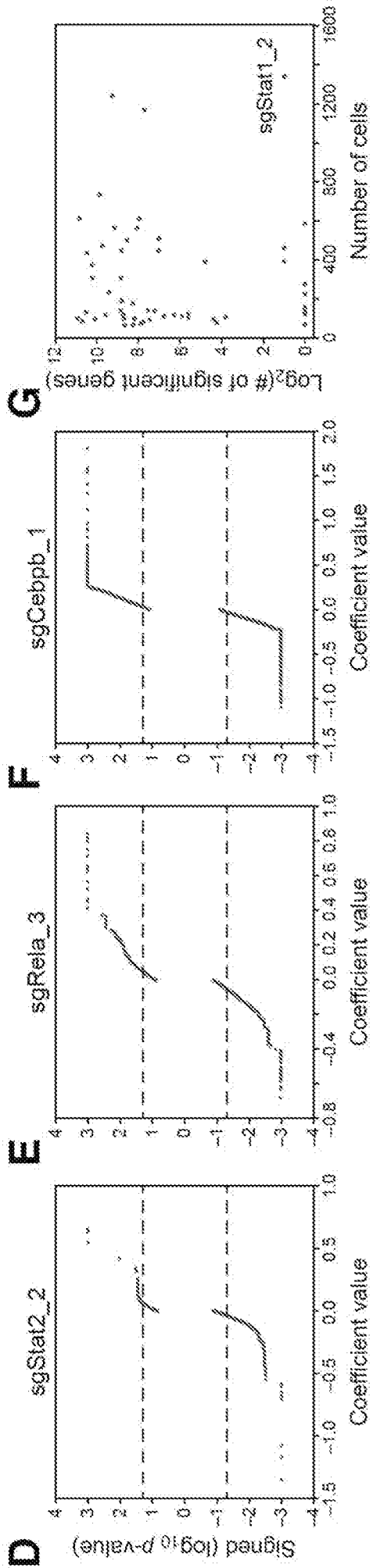


FIG. 56H

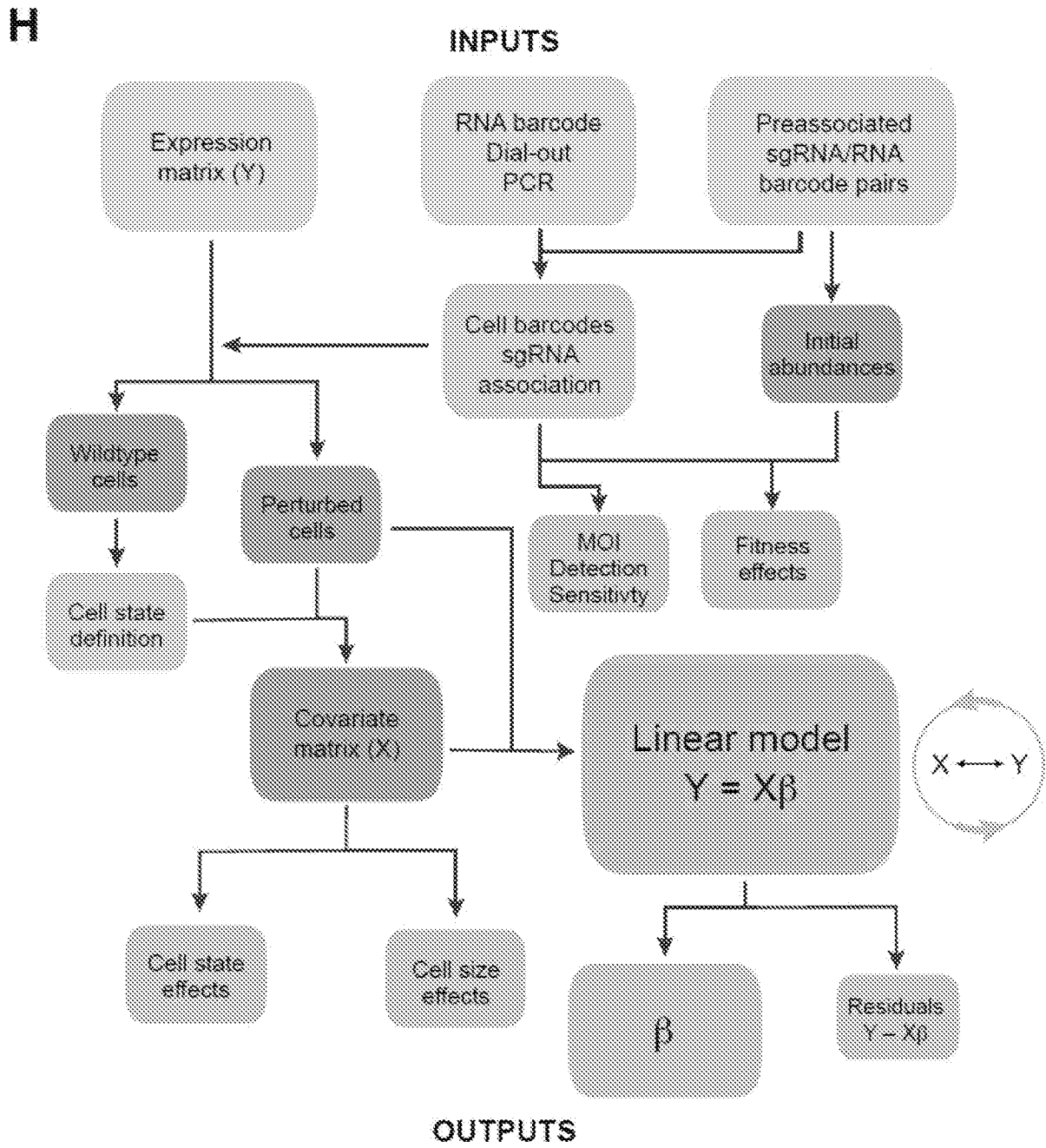


FIG. 56I

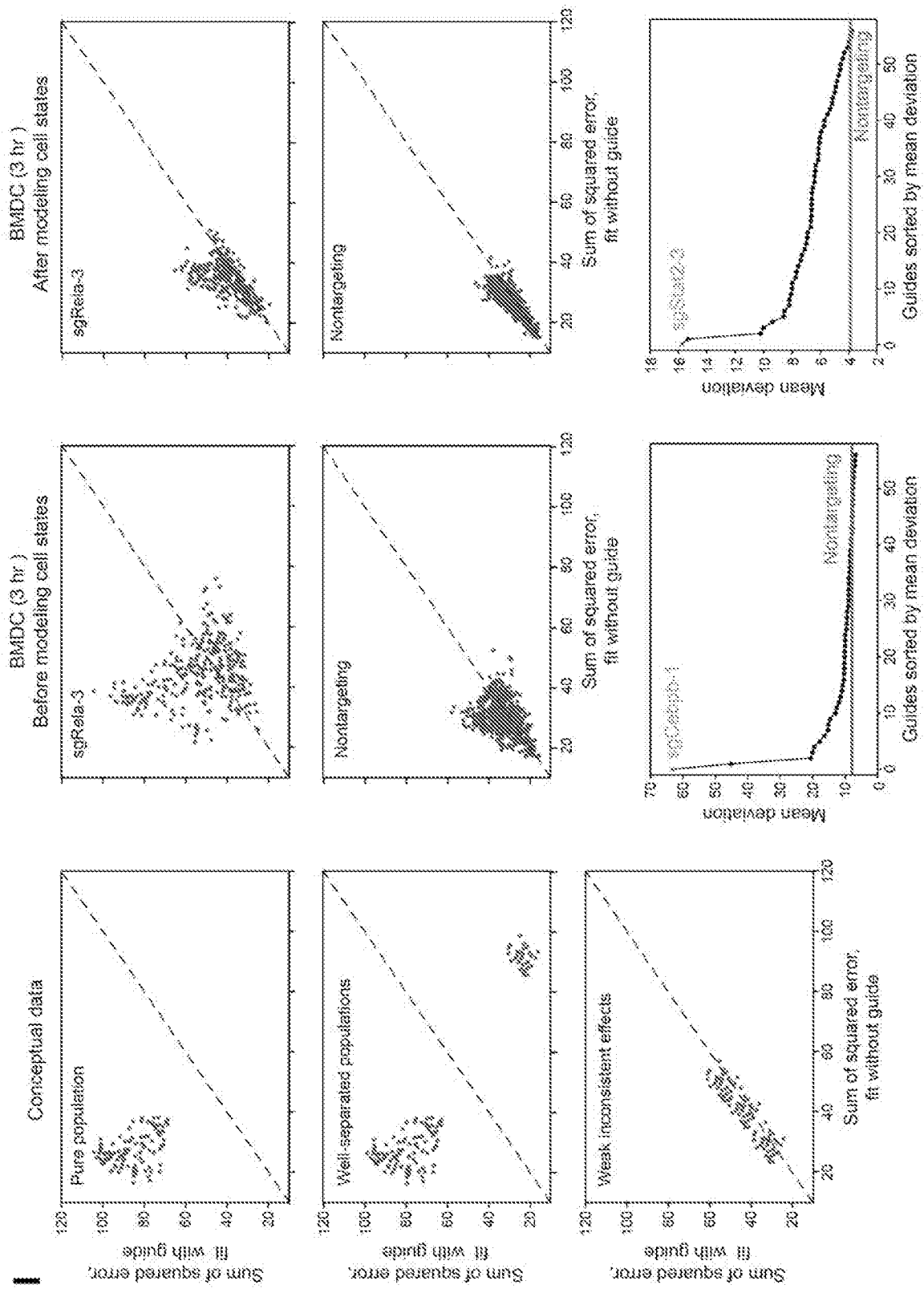


FIG. 57A-B

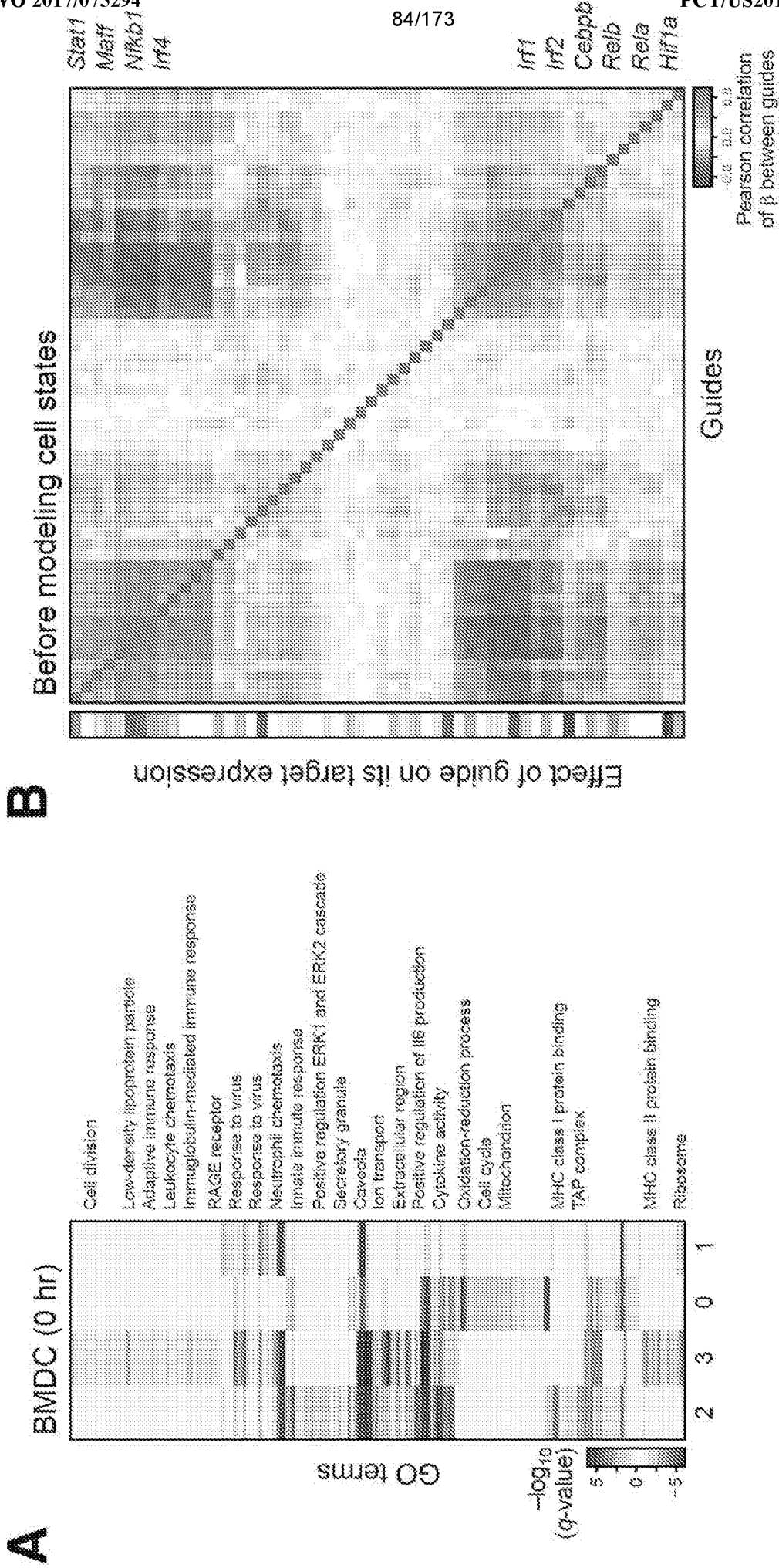


FIG. 57C

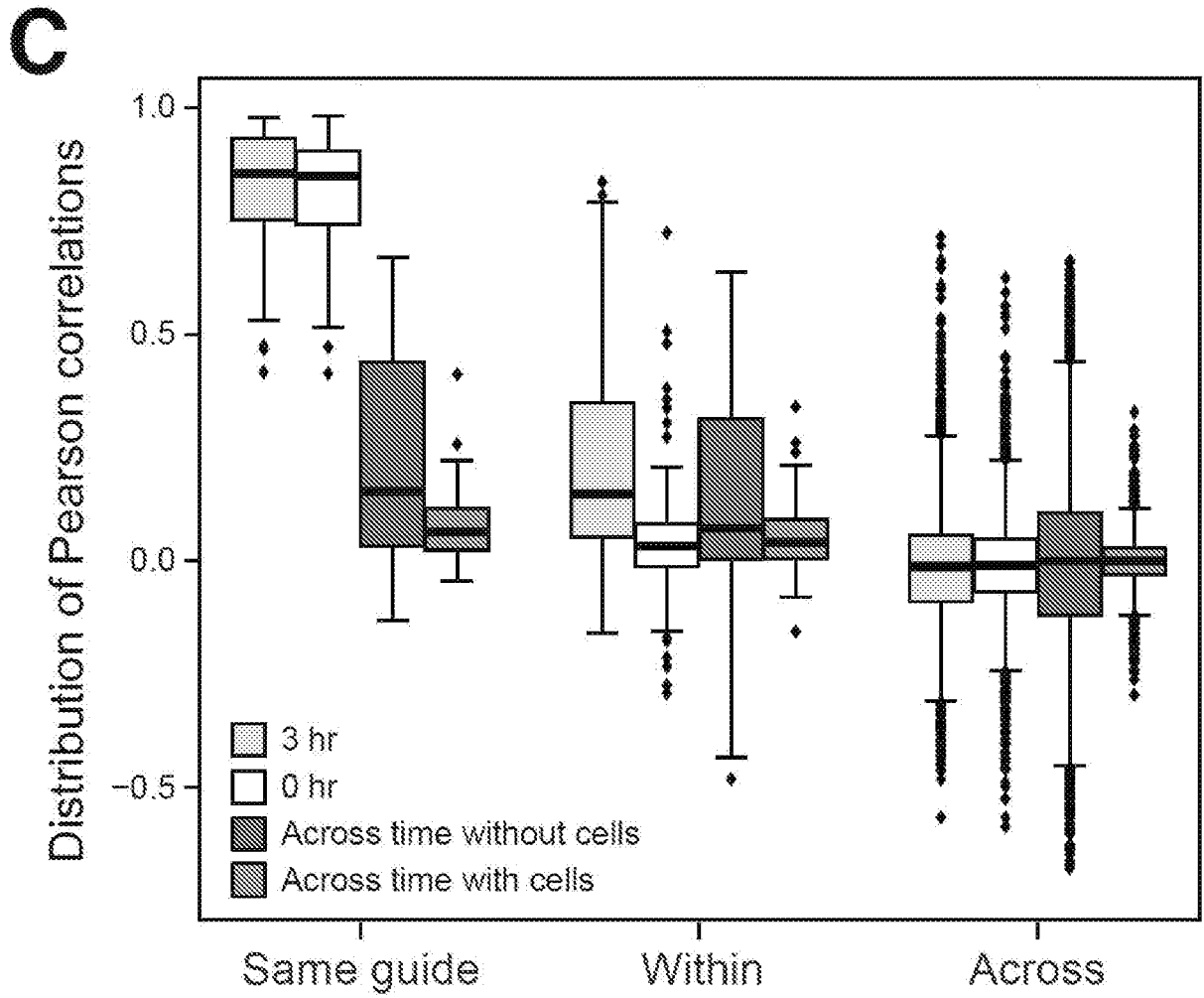


FIG. 57D

D

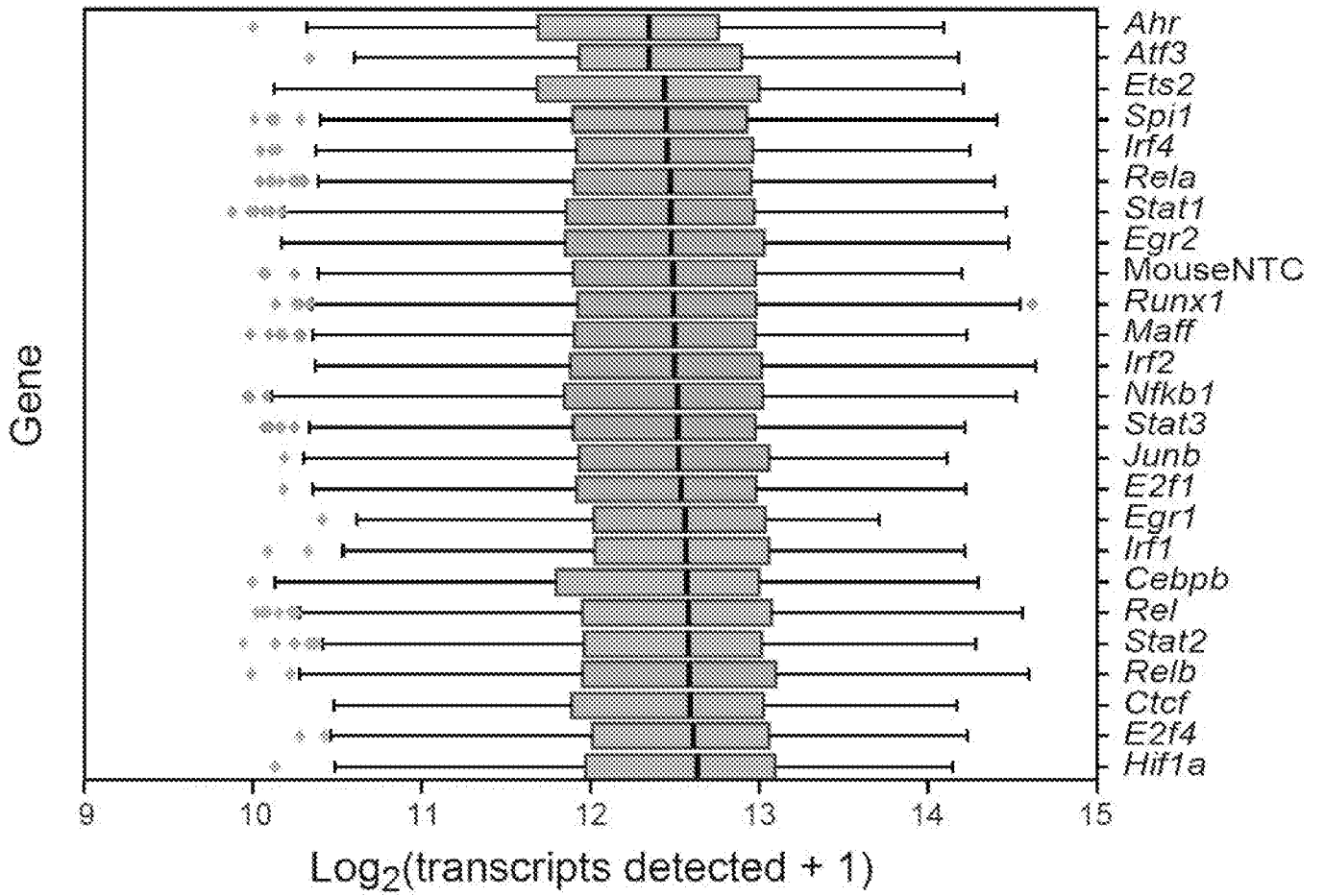


FIG. 57E

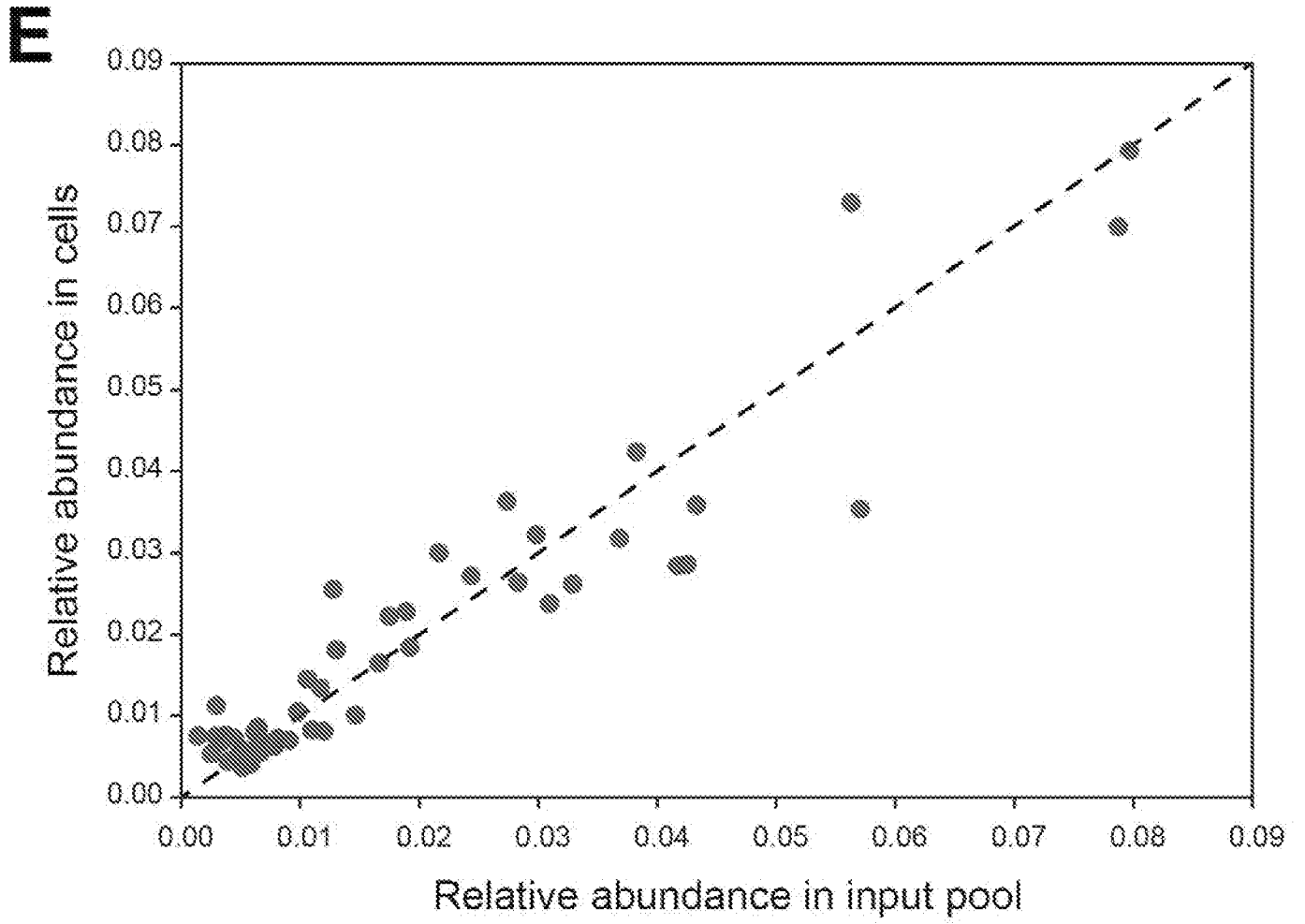


FIG. 57F

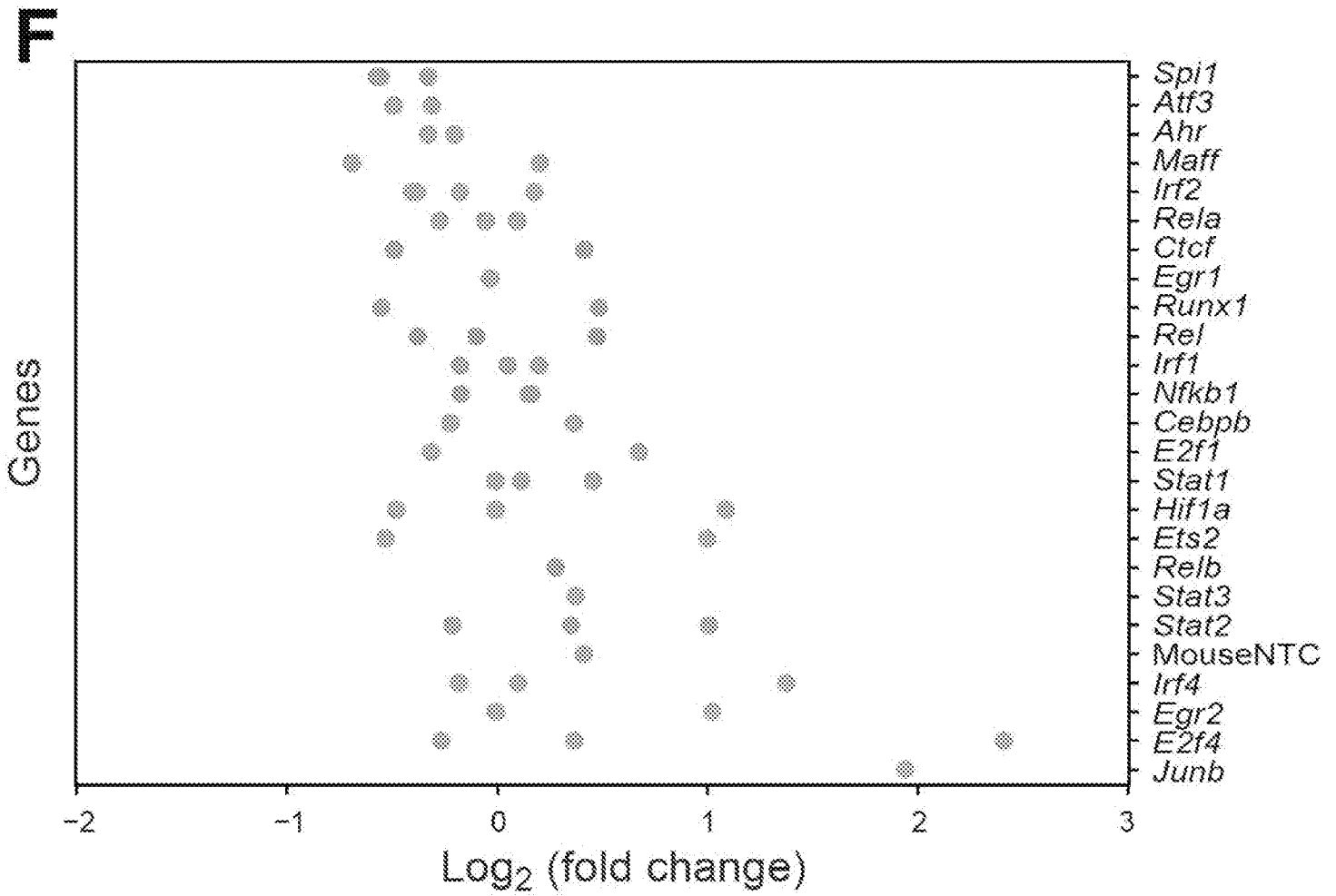


FIG. 58A

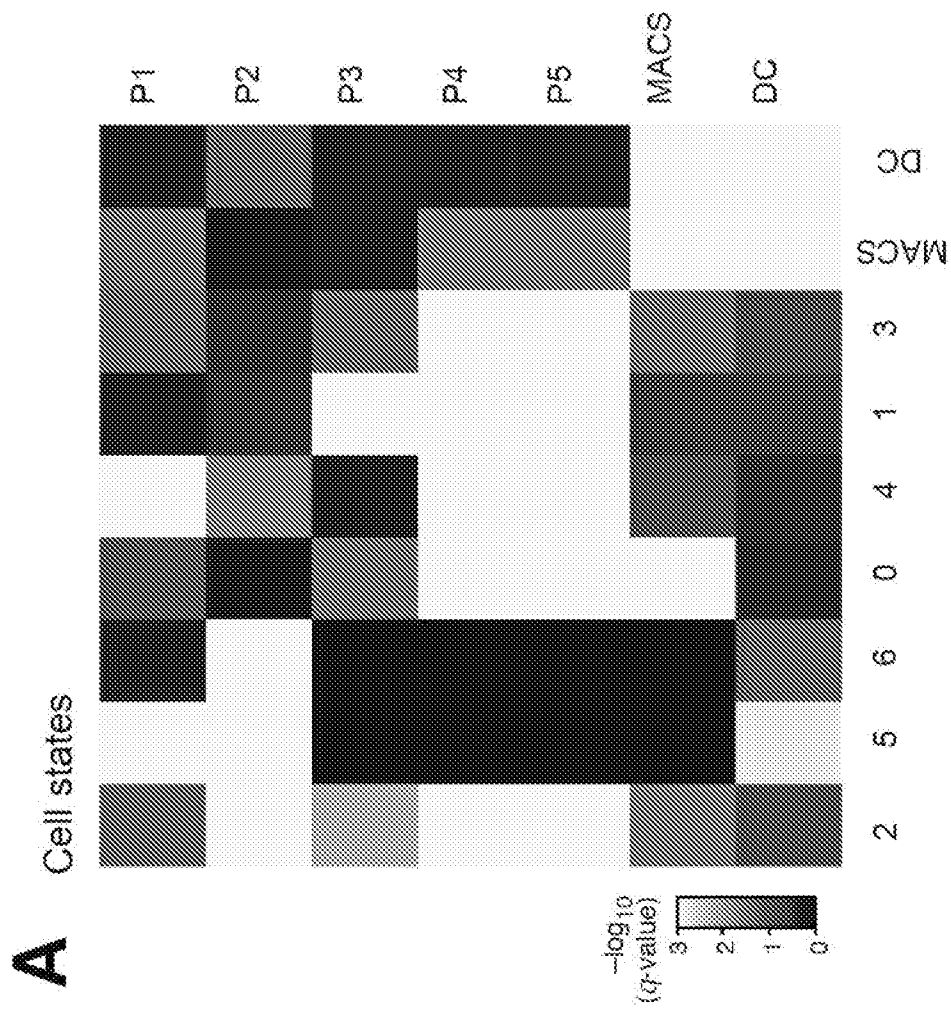


FIG. 58C

C

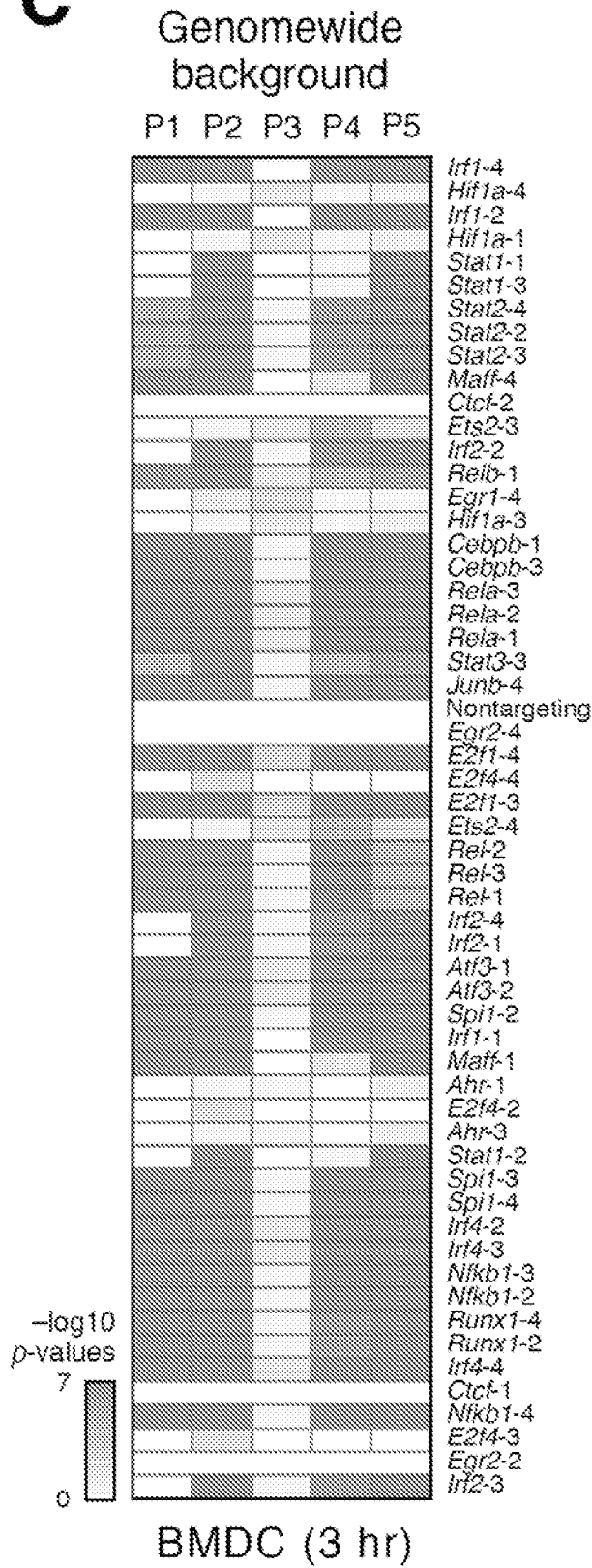


FIG. 58D

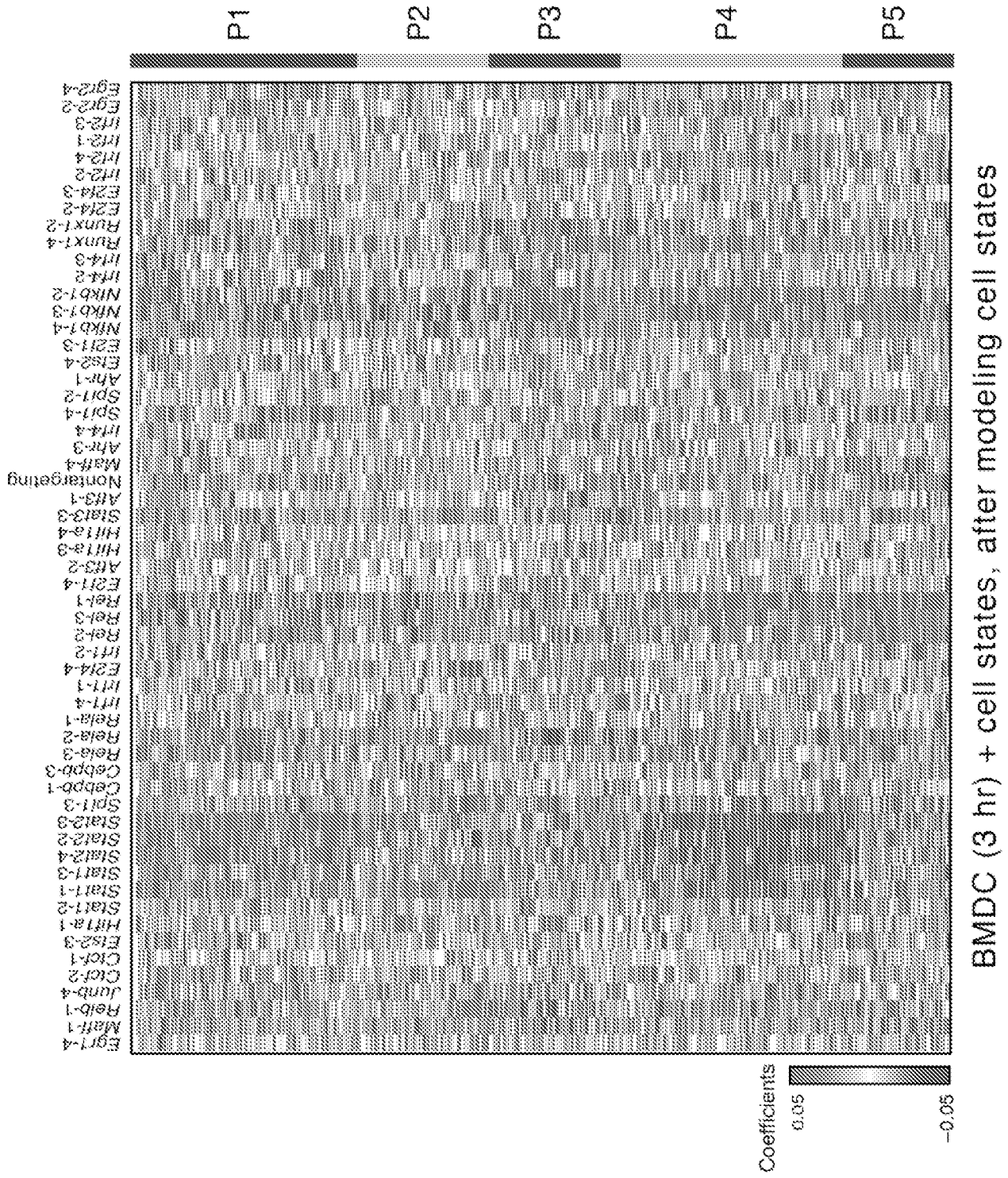
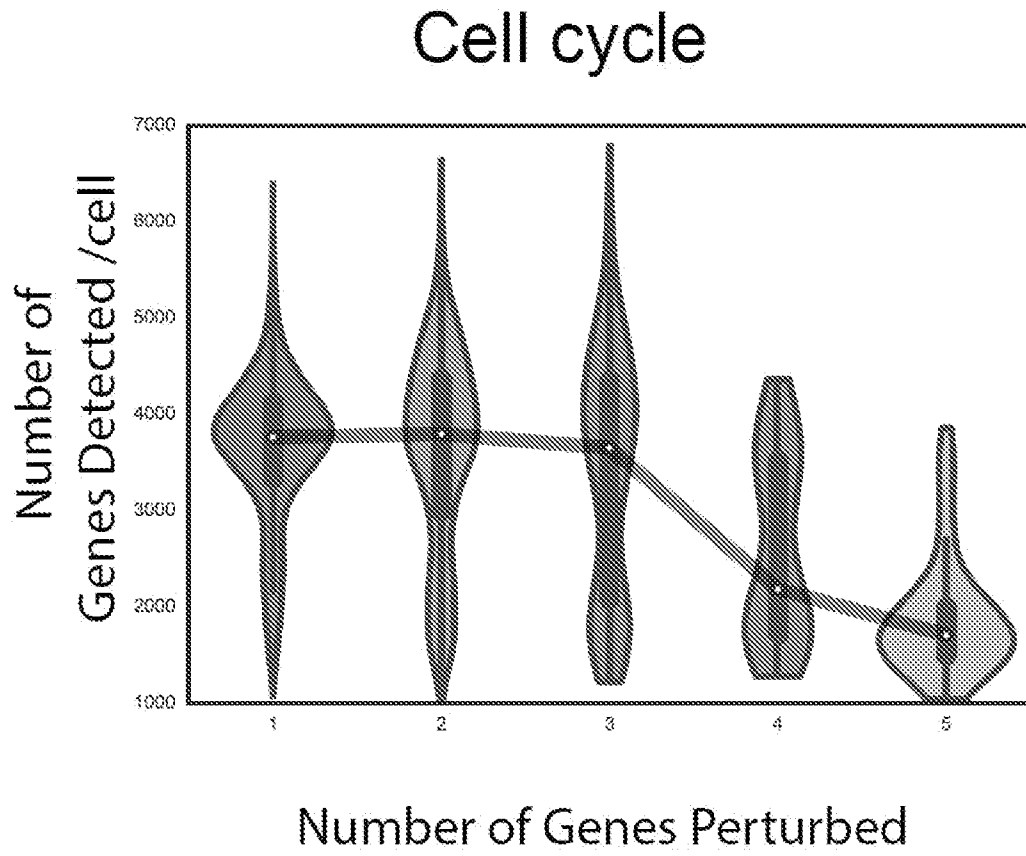


FIG. 59



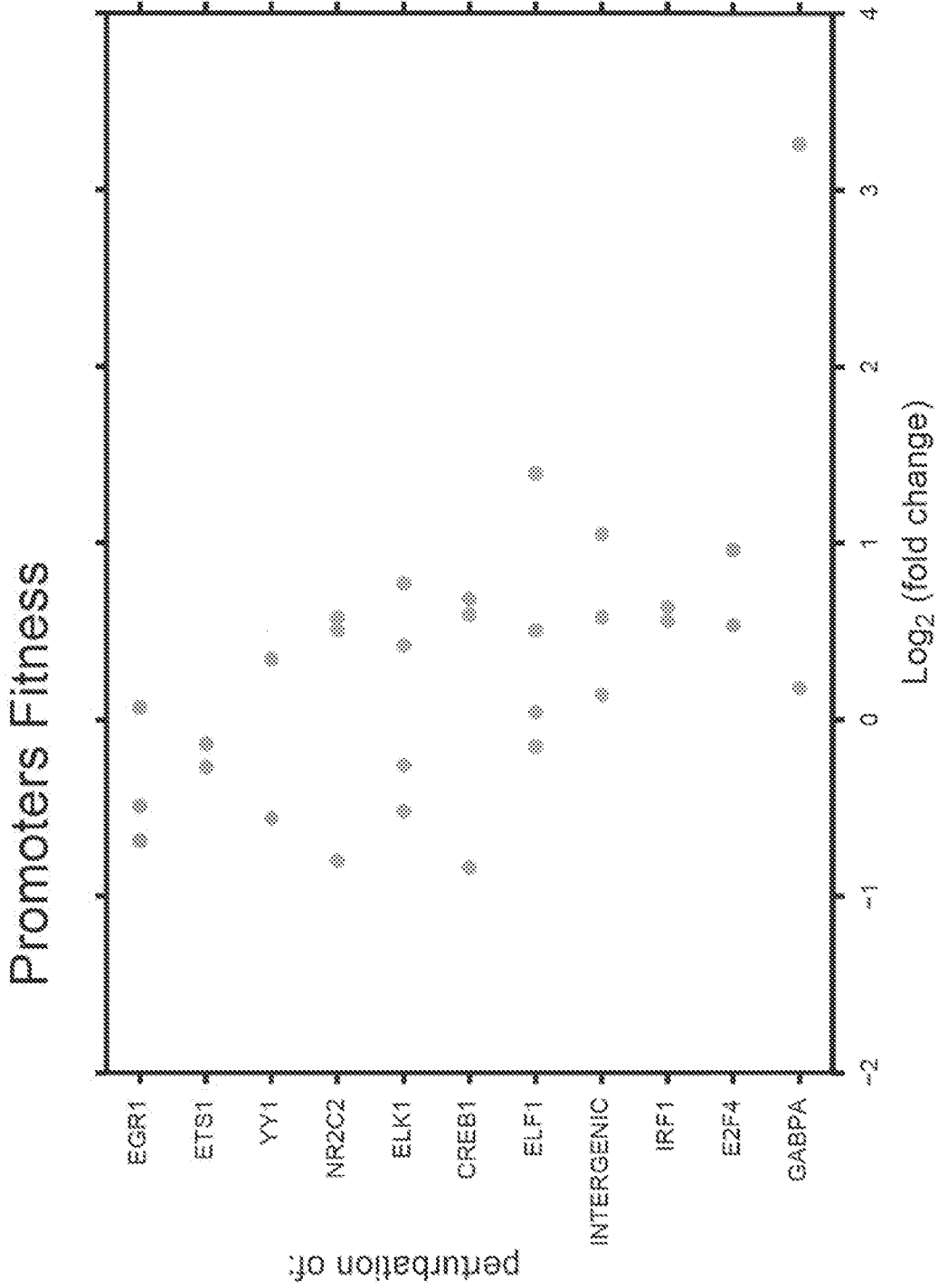
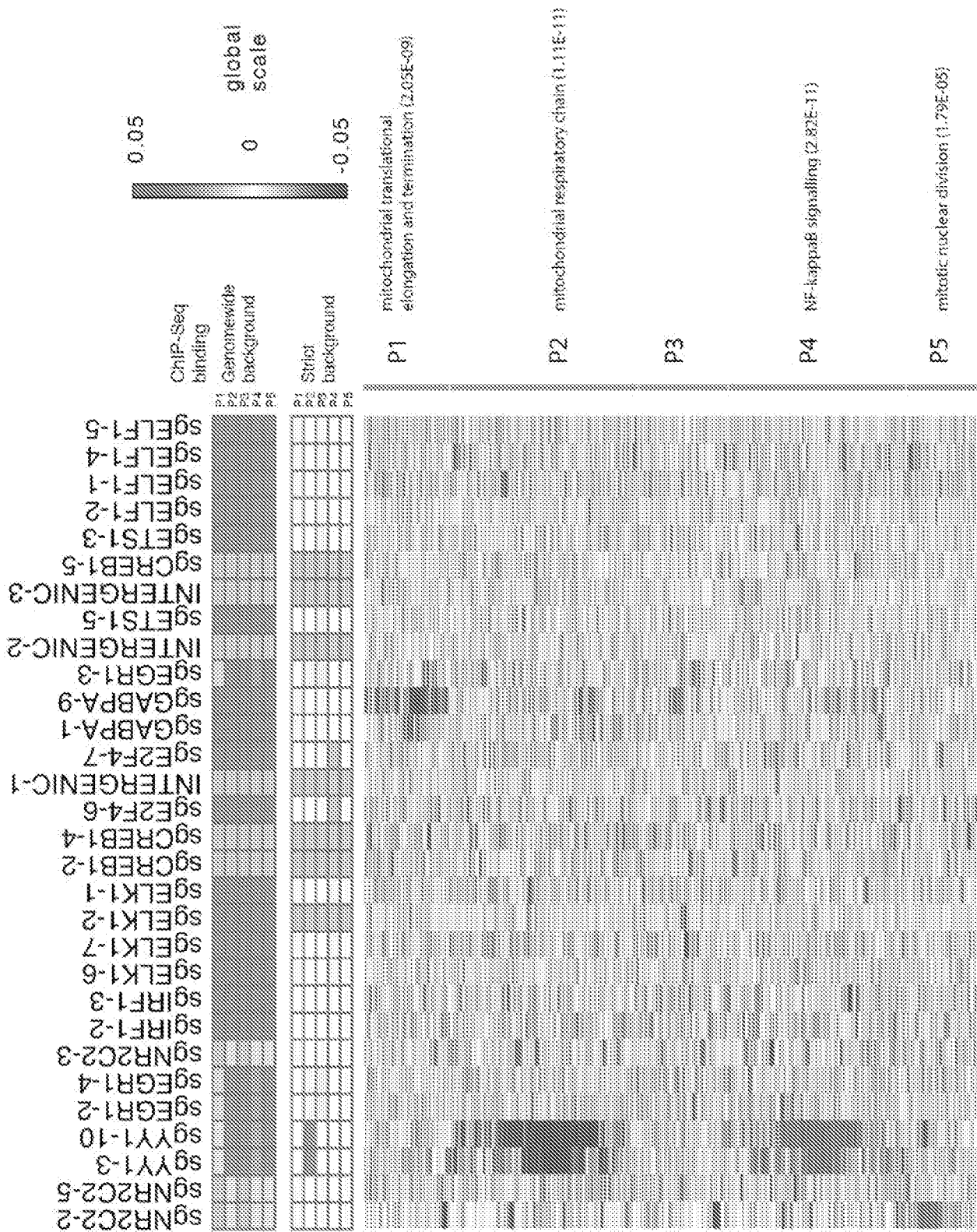


FIG. 60A

FIG. 60B



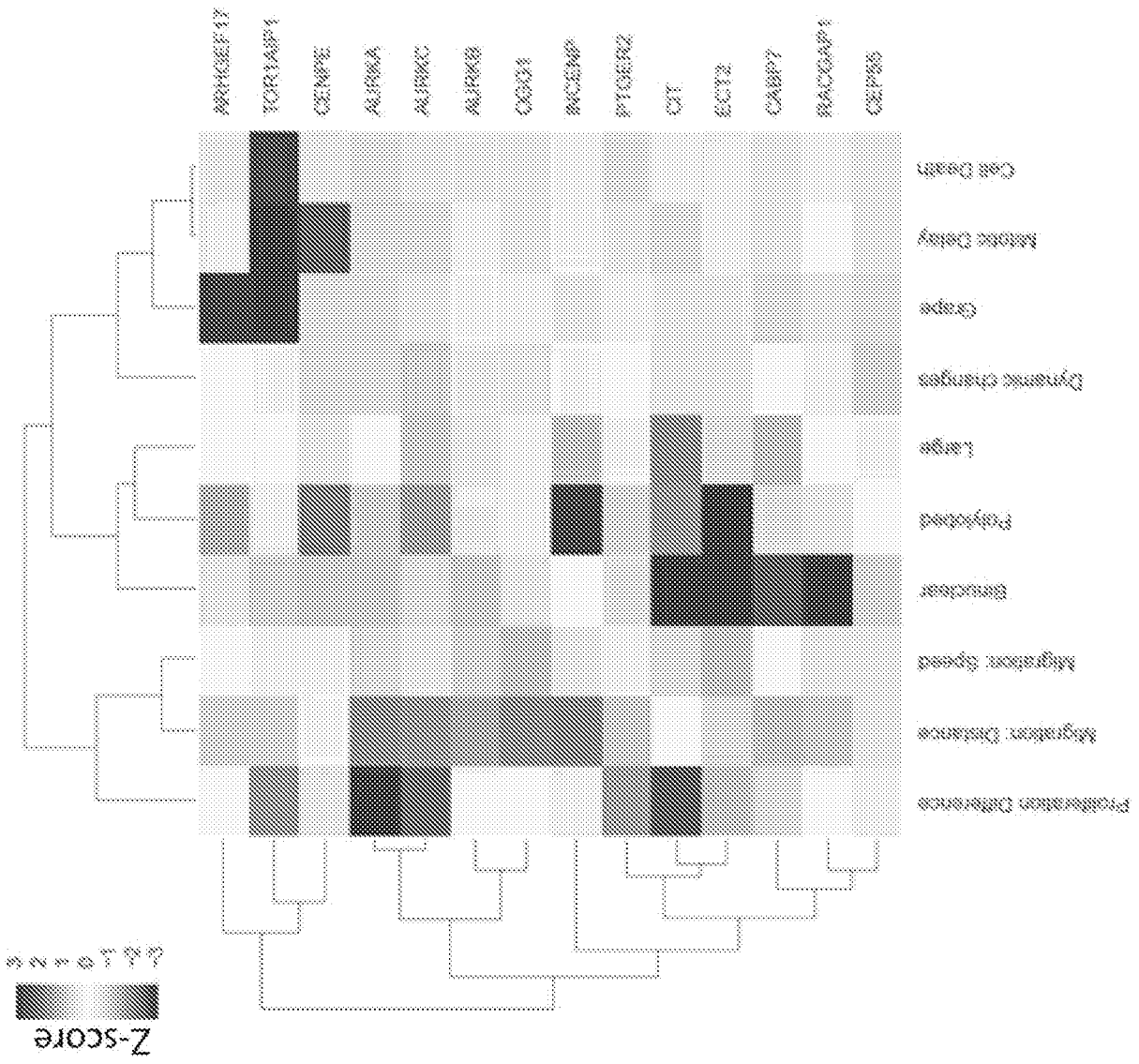


FIG. 60C

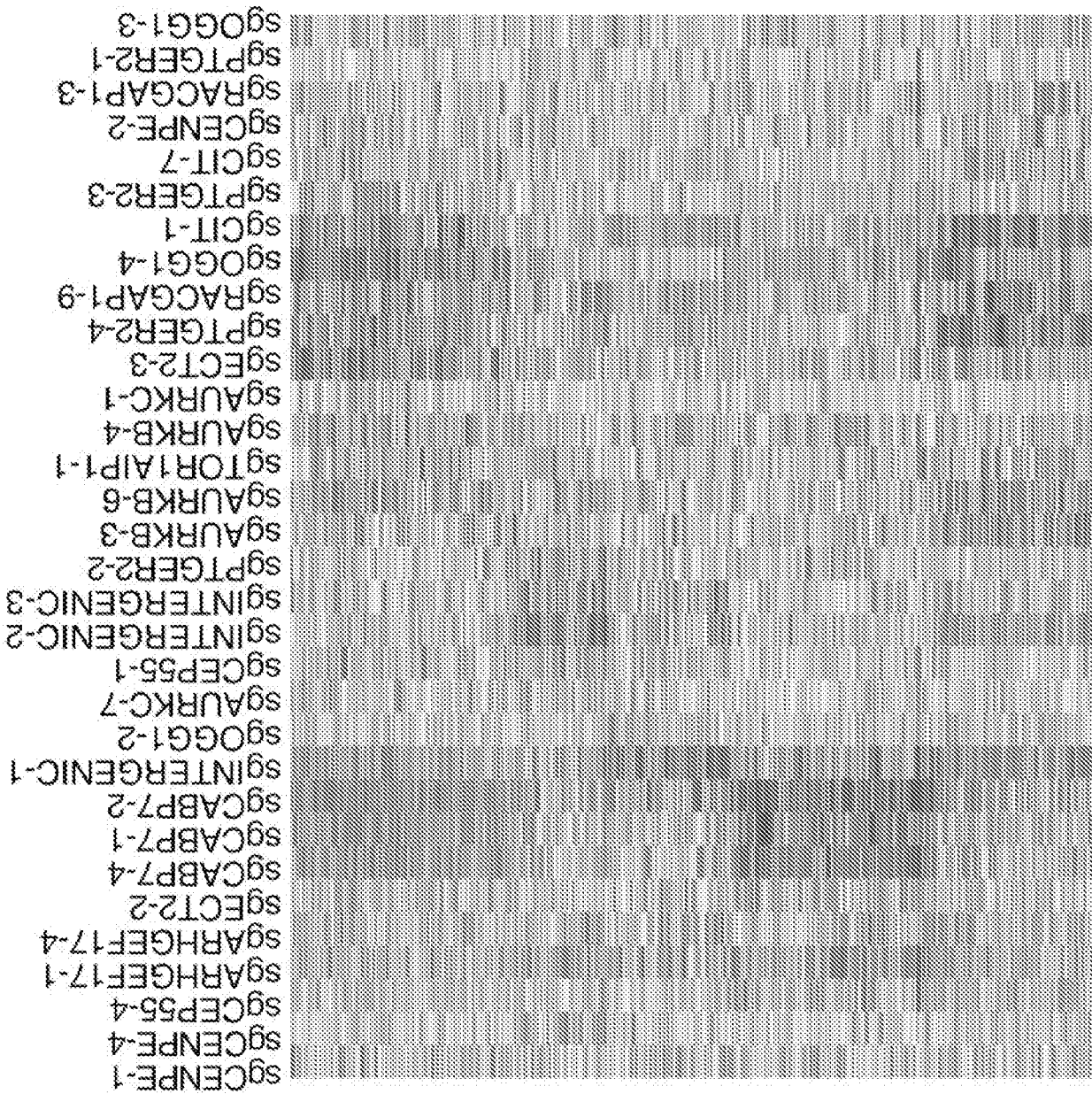
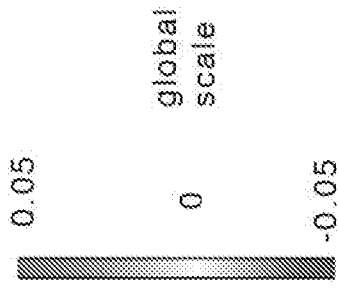


FIG. 60D

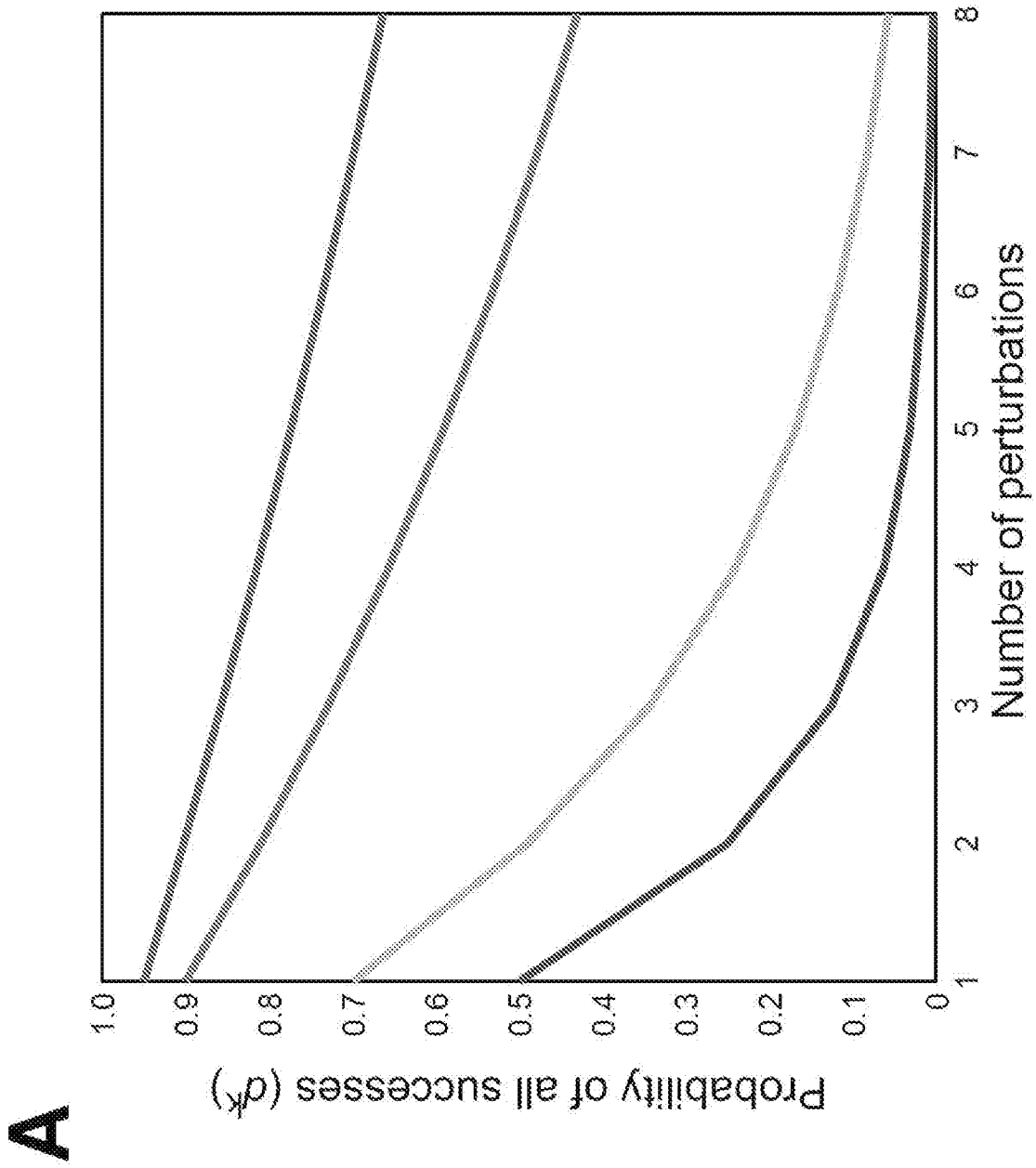
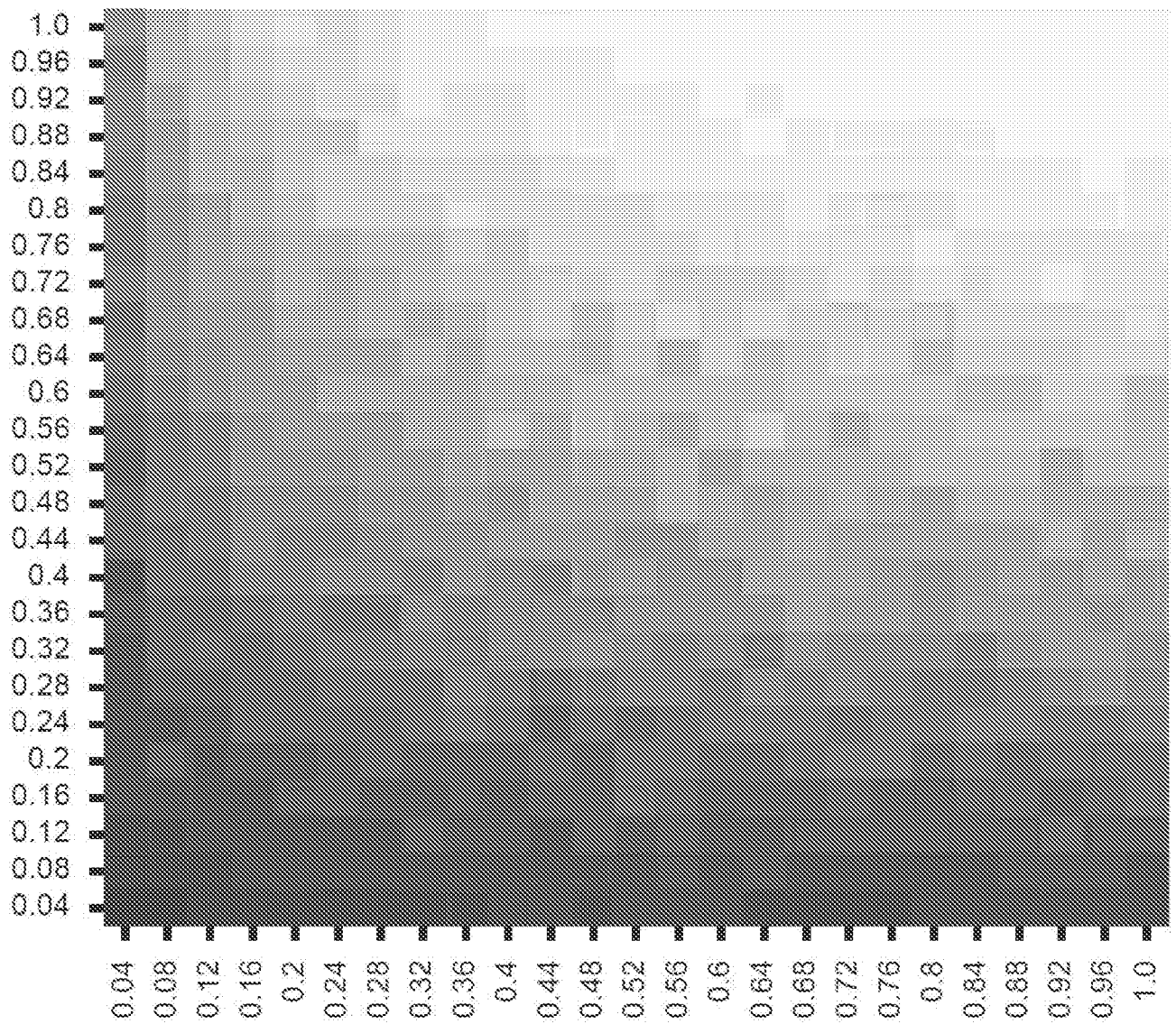


FIG. 61A

FIG. 61B

PCA Scores



Fraction of transcripts

Sensitivity
FDR<0.1

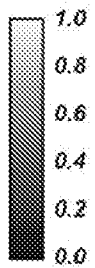


FIG. 61C

PCA Scores

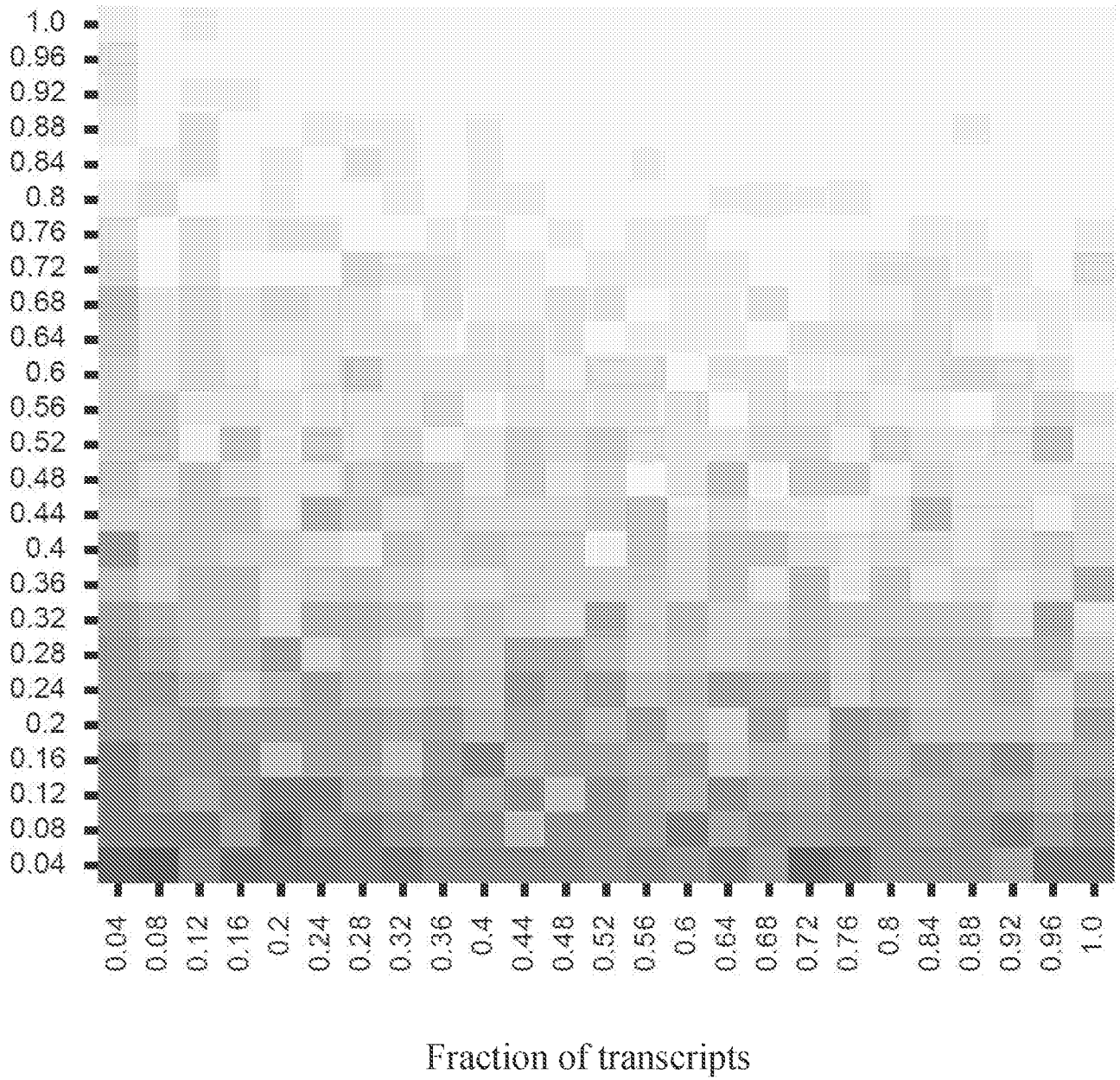


FIG. 61D

PCA Scores

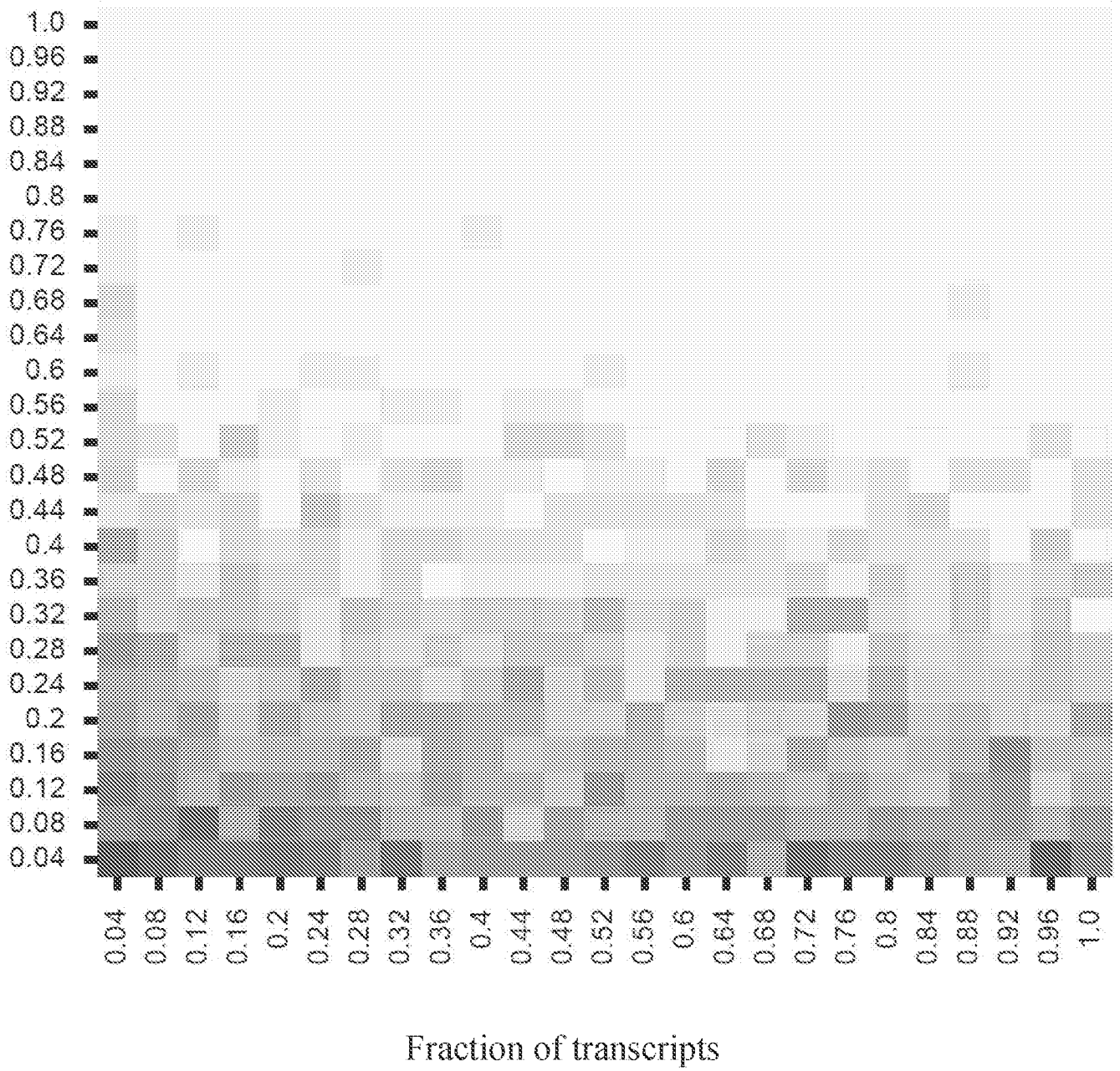
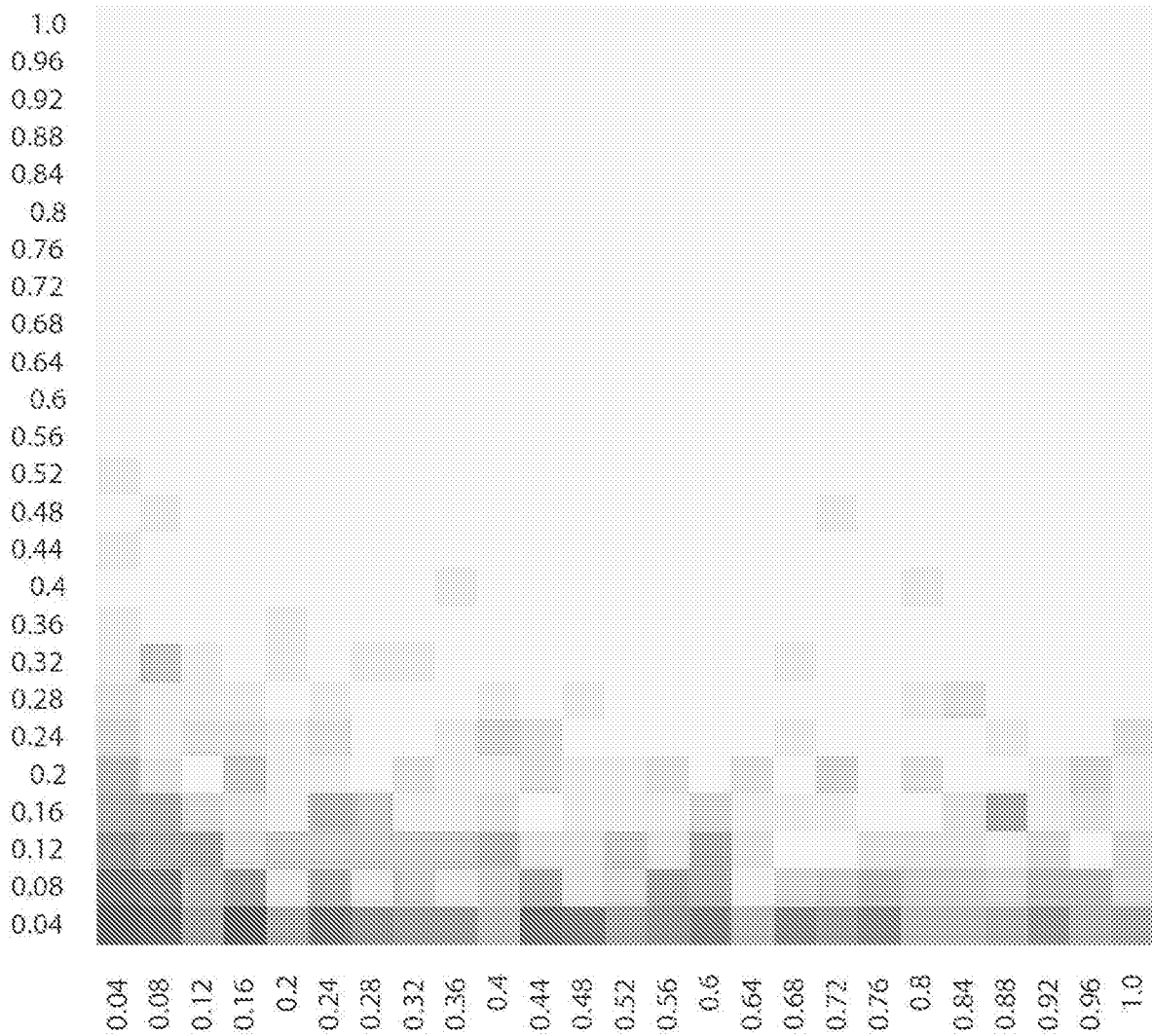


FIG. 61E

Gene Signatures

<-0.1

E



Fraction of transcripts

FIG. 61F

Gene Signatures

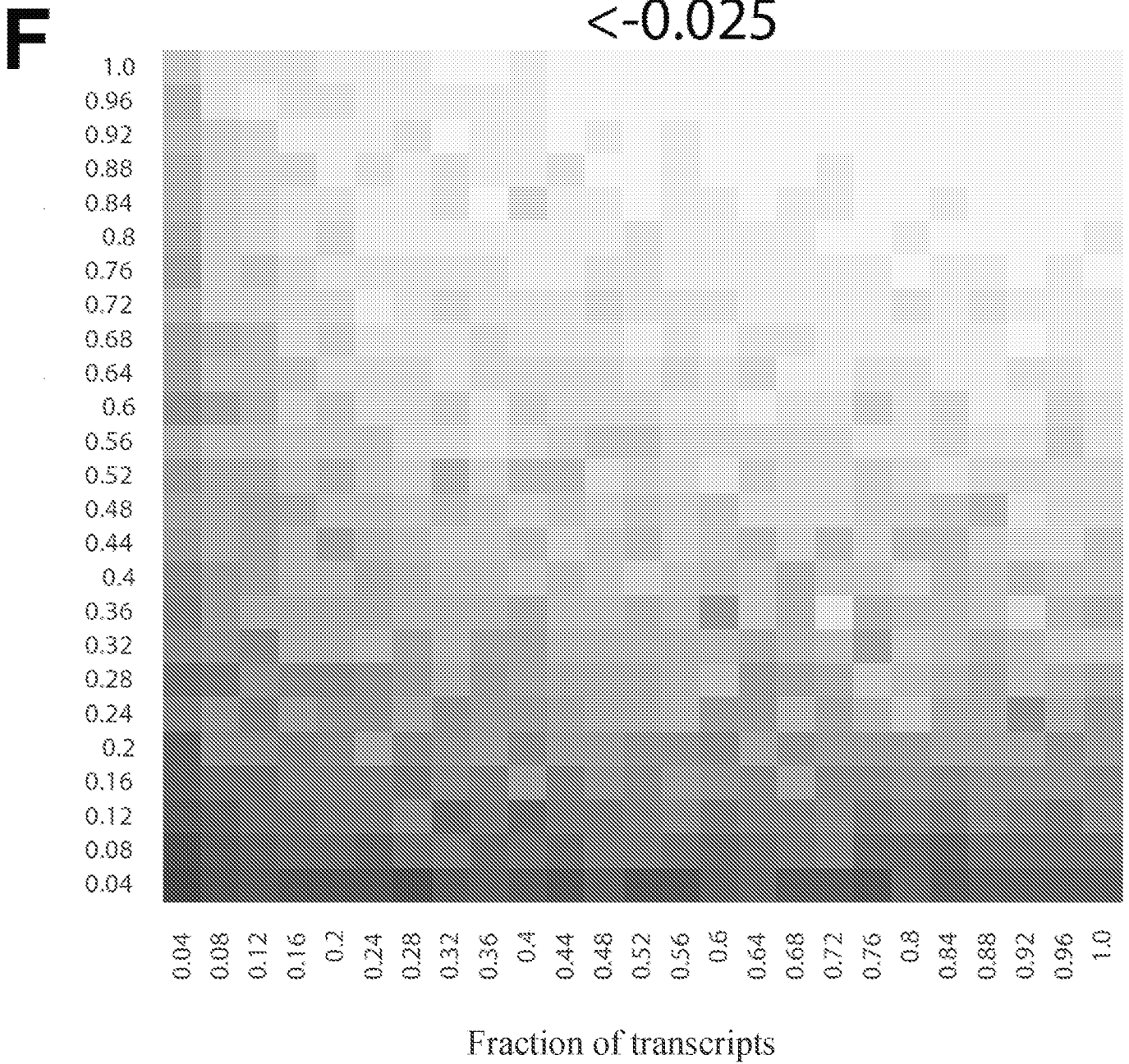


FIG. 61G

Gene Signatures

>0.025

G

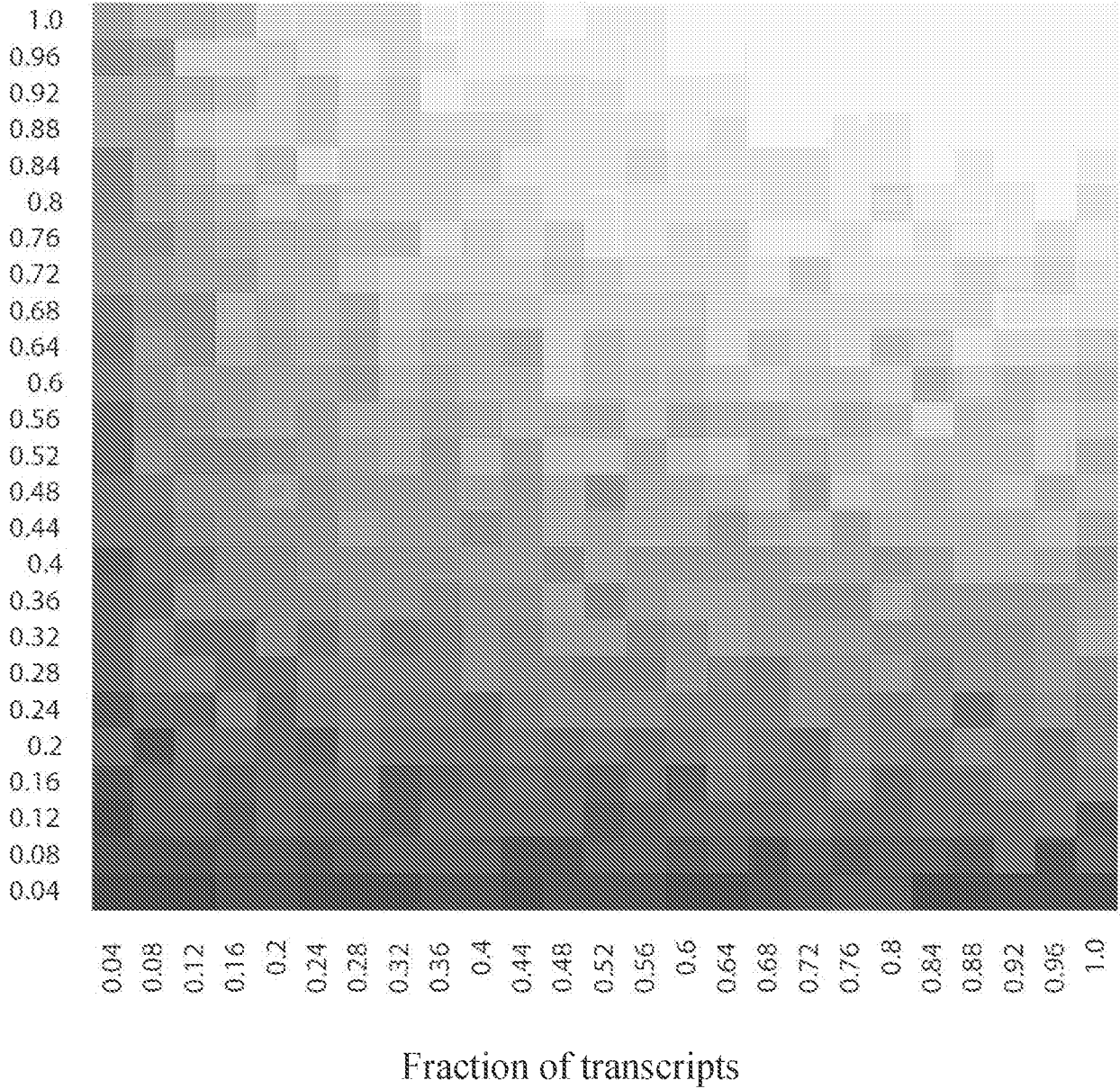


FIG. 61H

Gene Signatures

H

>0.1

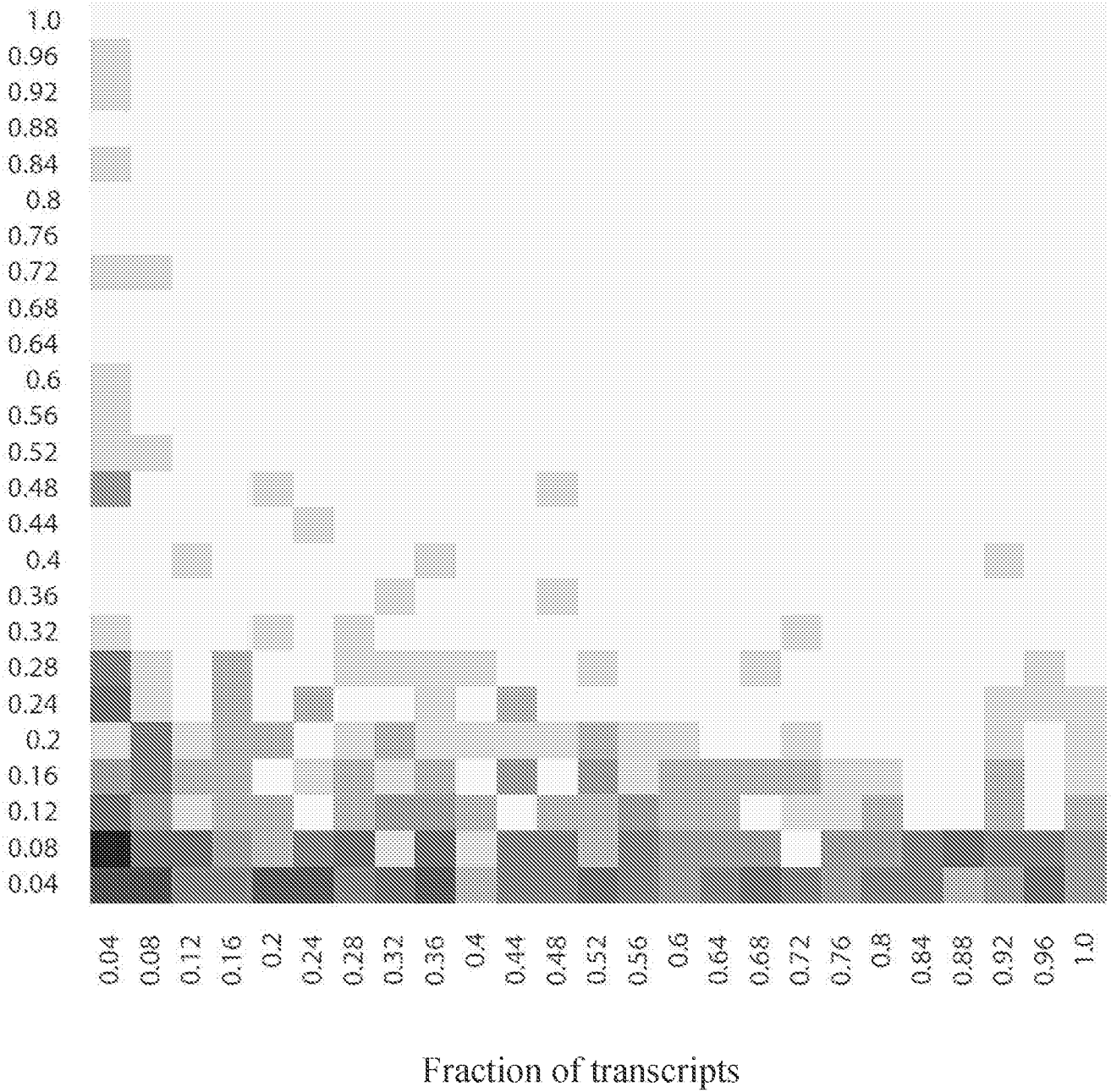
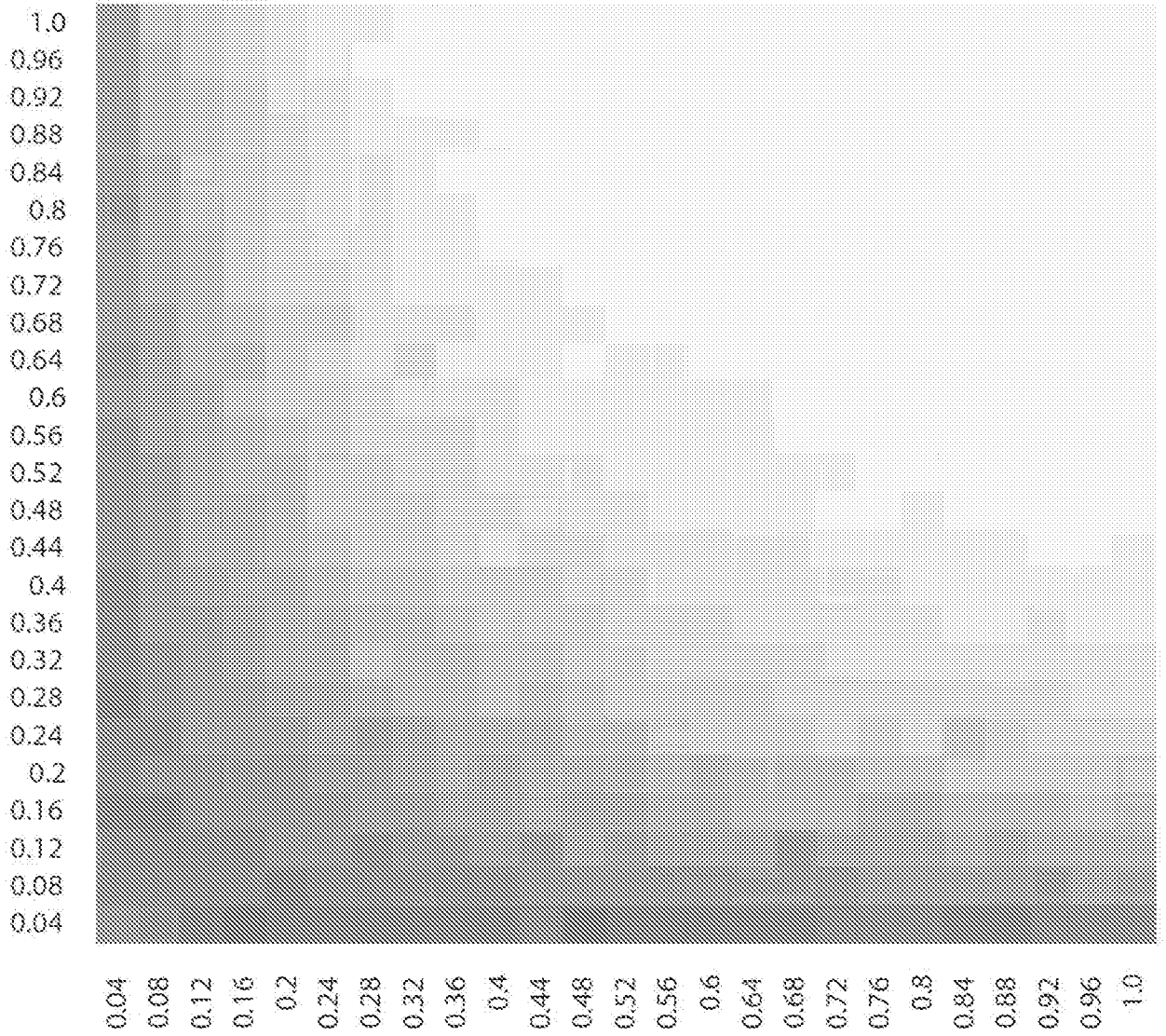


FIG. 6II

Individual Gene

<-0.25



Fraction of transcripts

FIG. 61J

Individual Gene

< -0.05

J

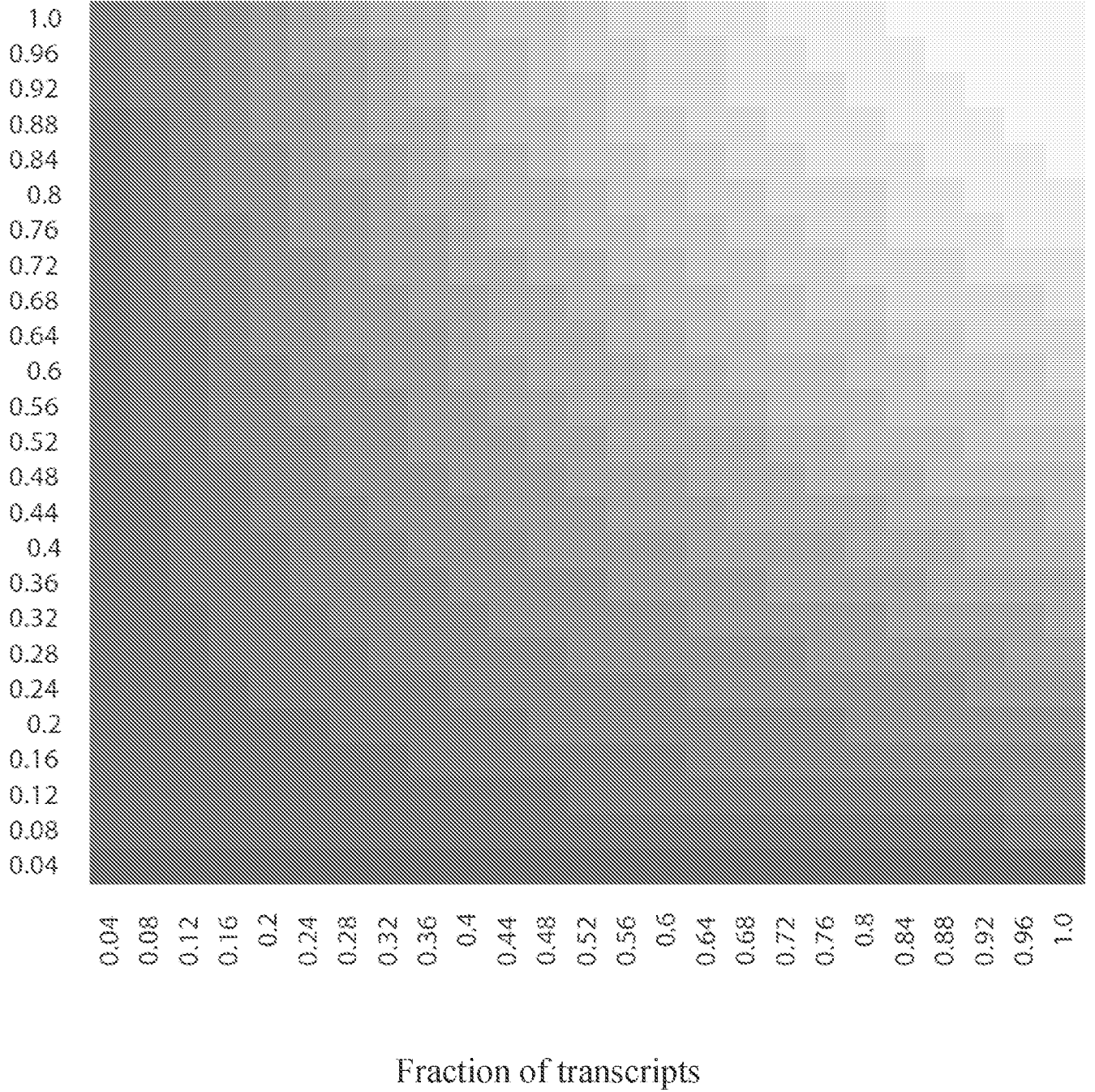


FIG. 61K

Individual Gene

K

>0.05

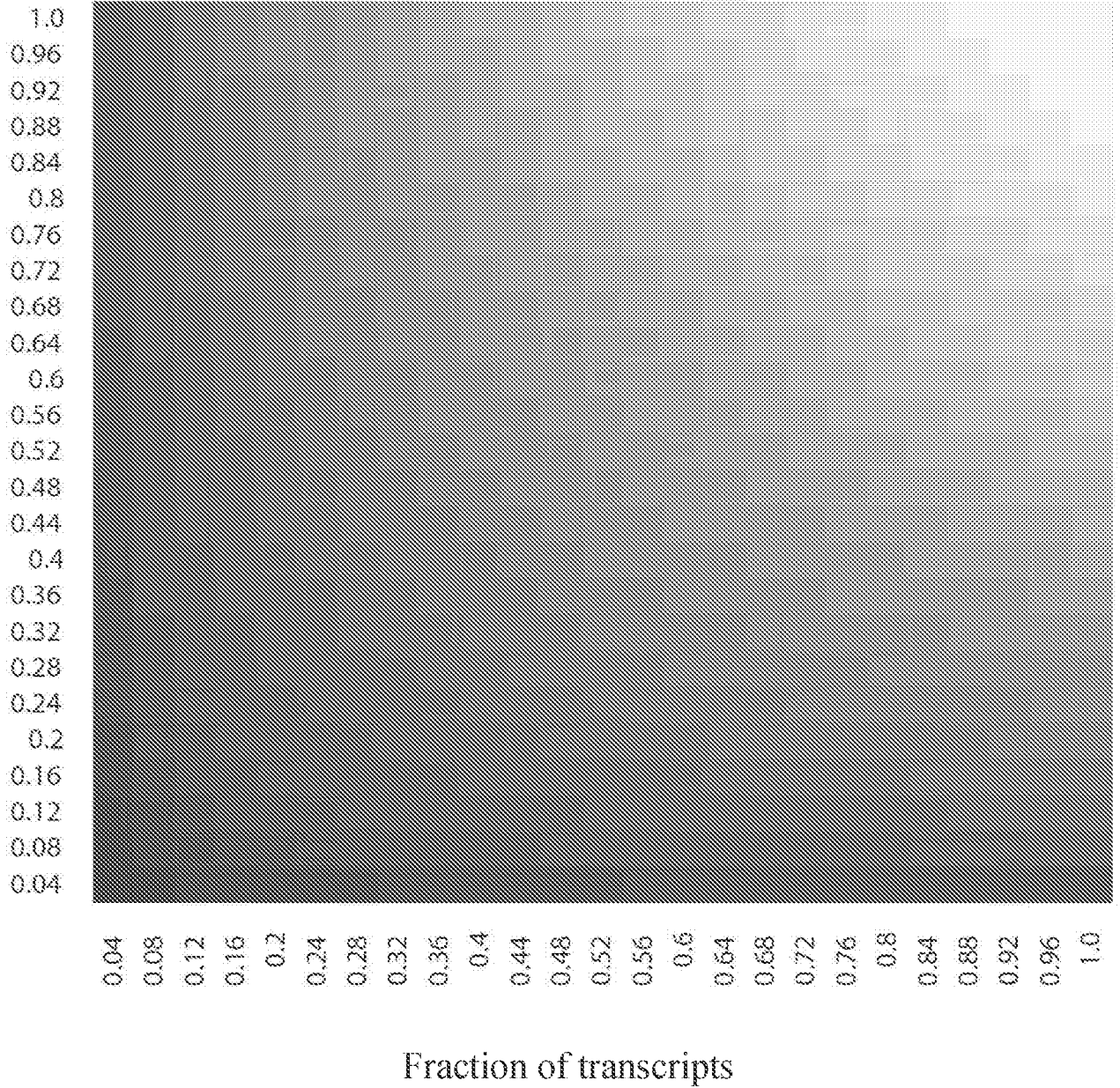
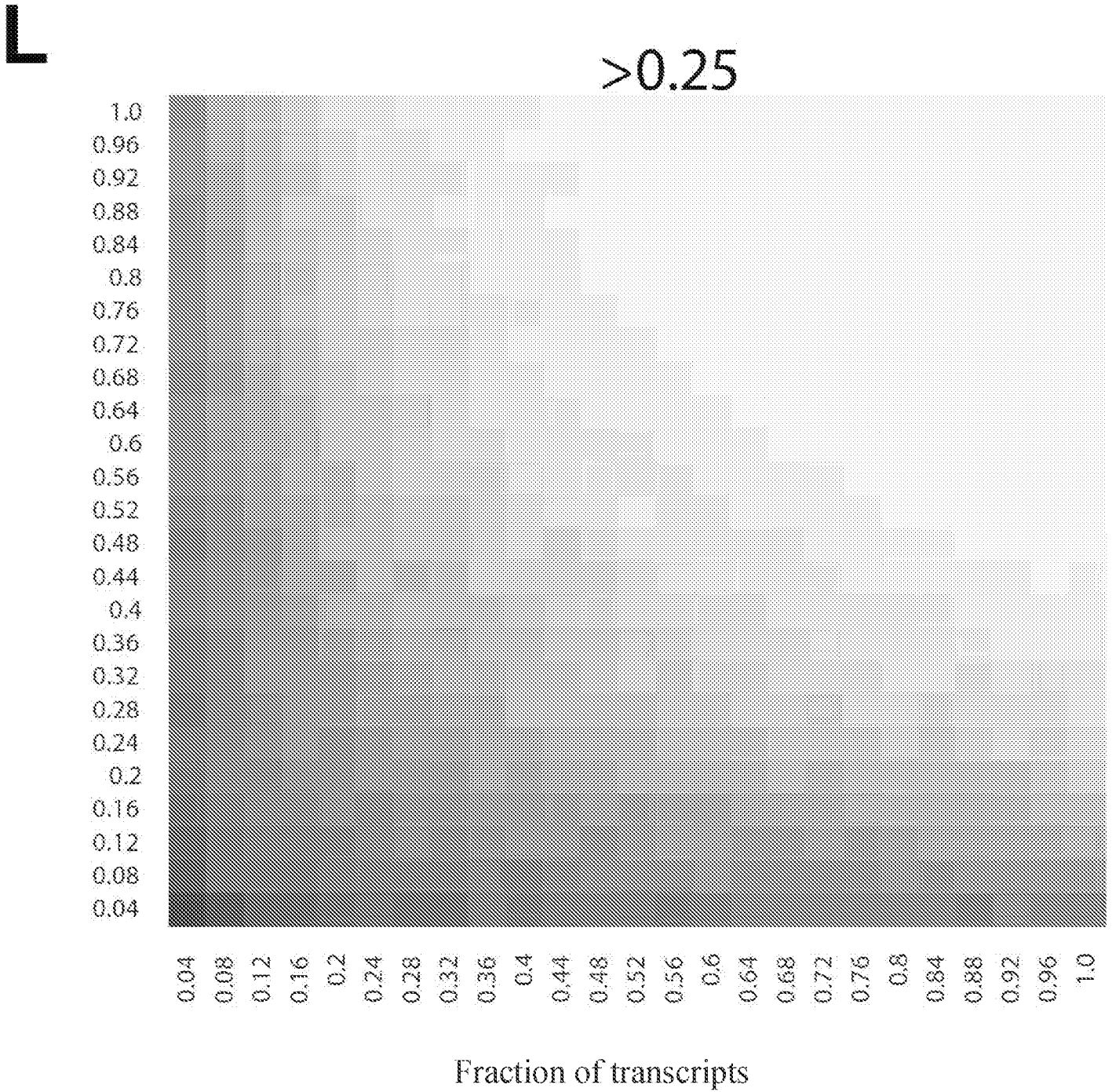


FIG. 61L

Individual Gene



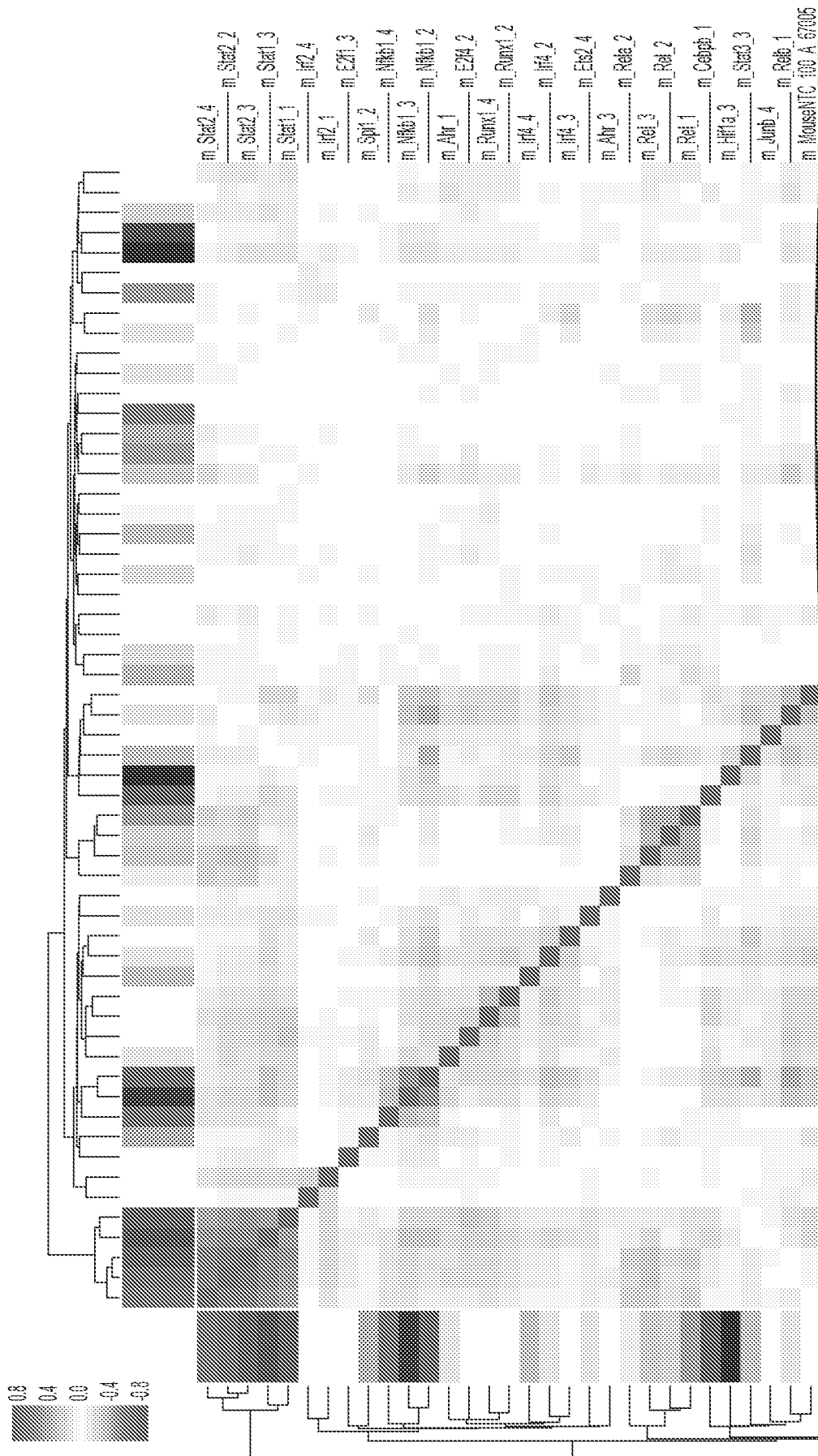


FIG. 62

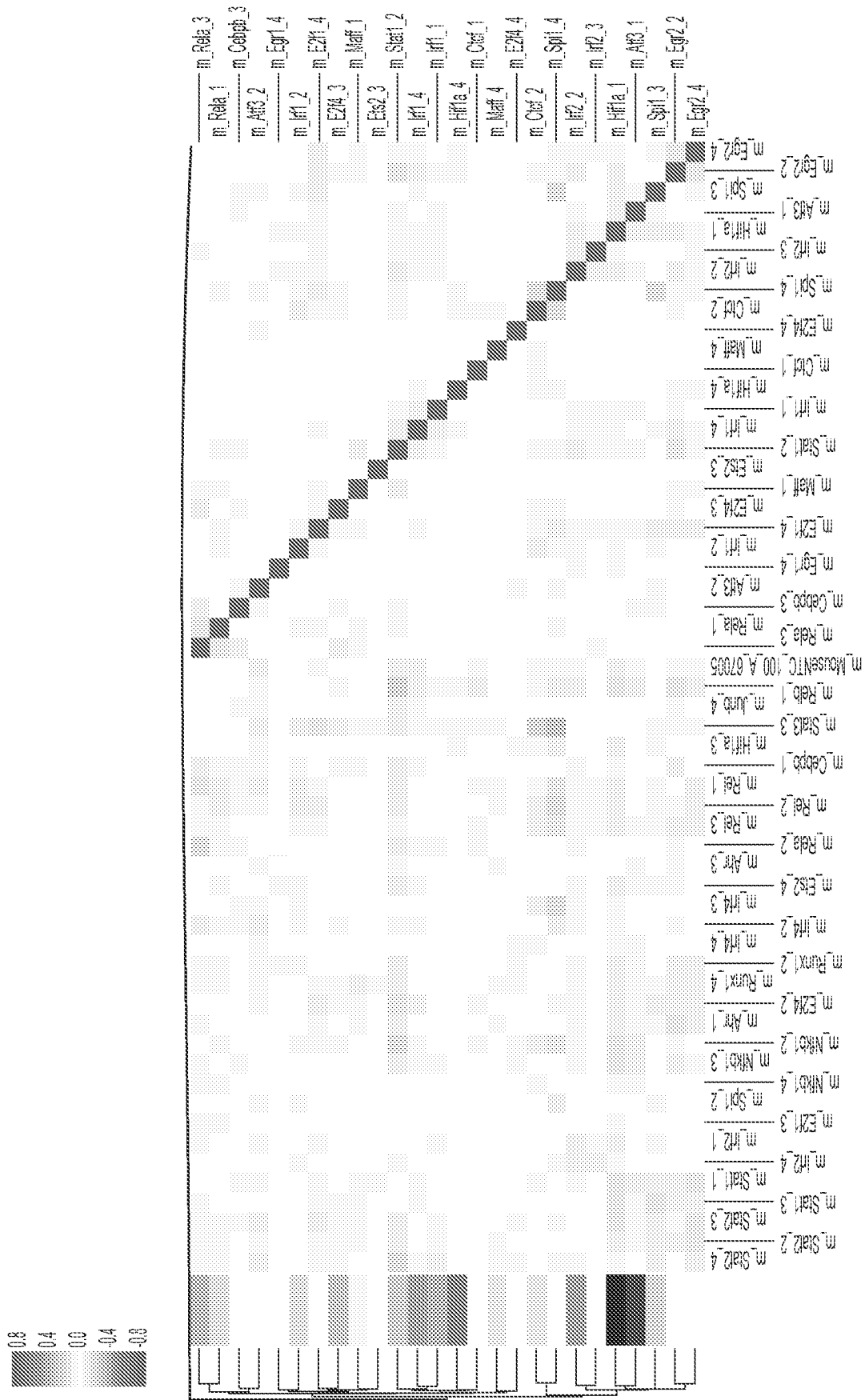
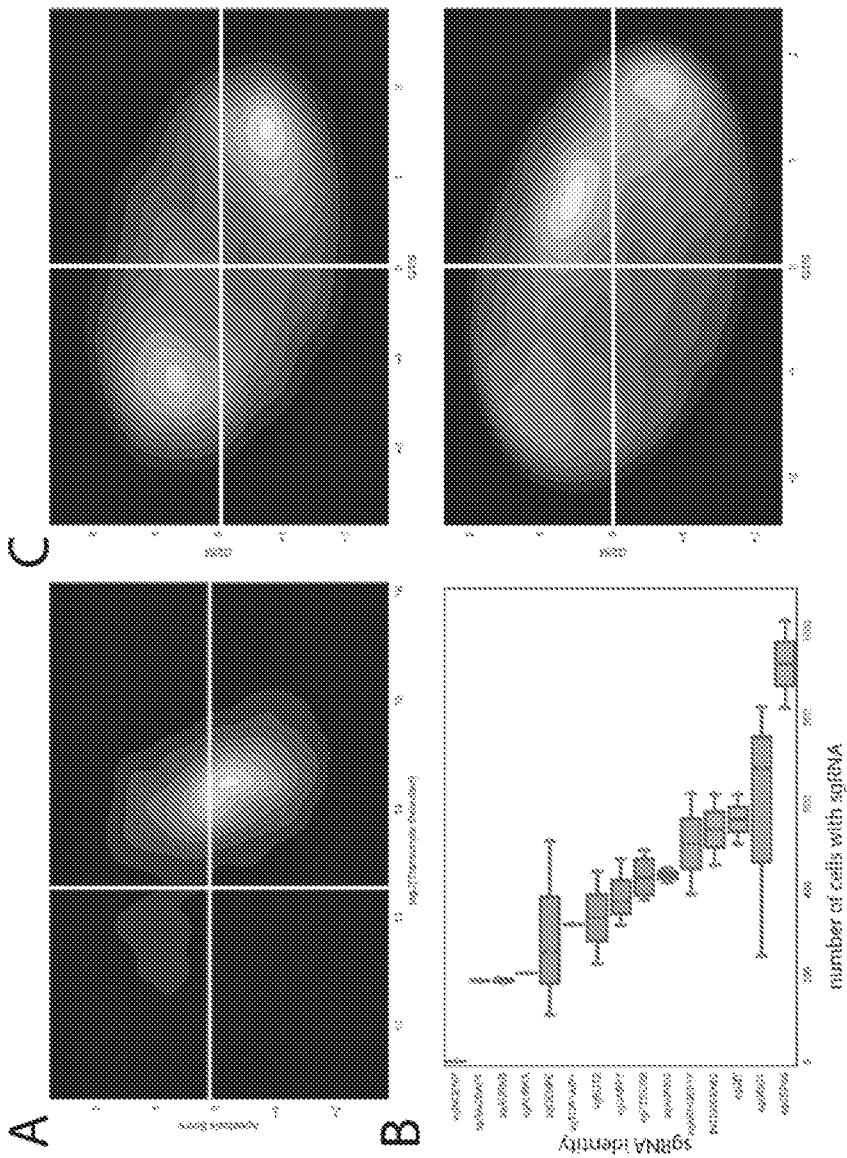


FIG. 62 continued

Figure 63A-C.



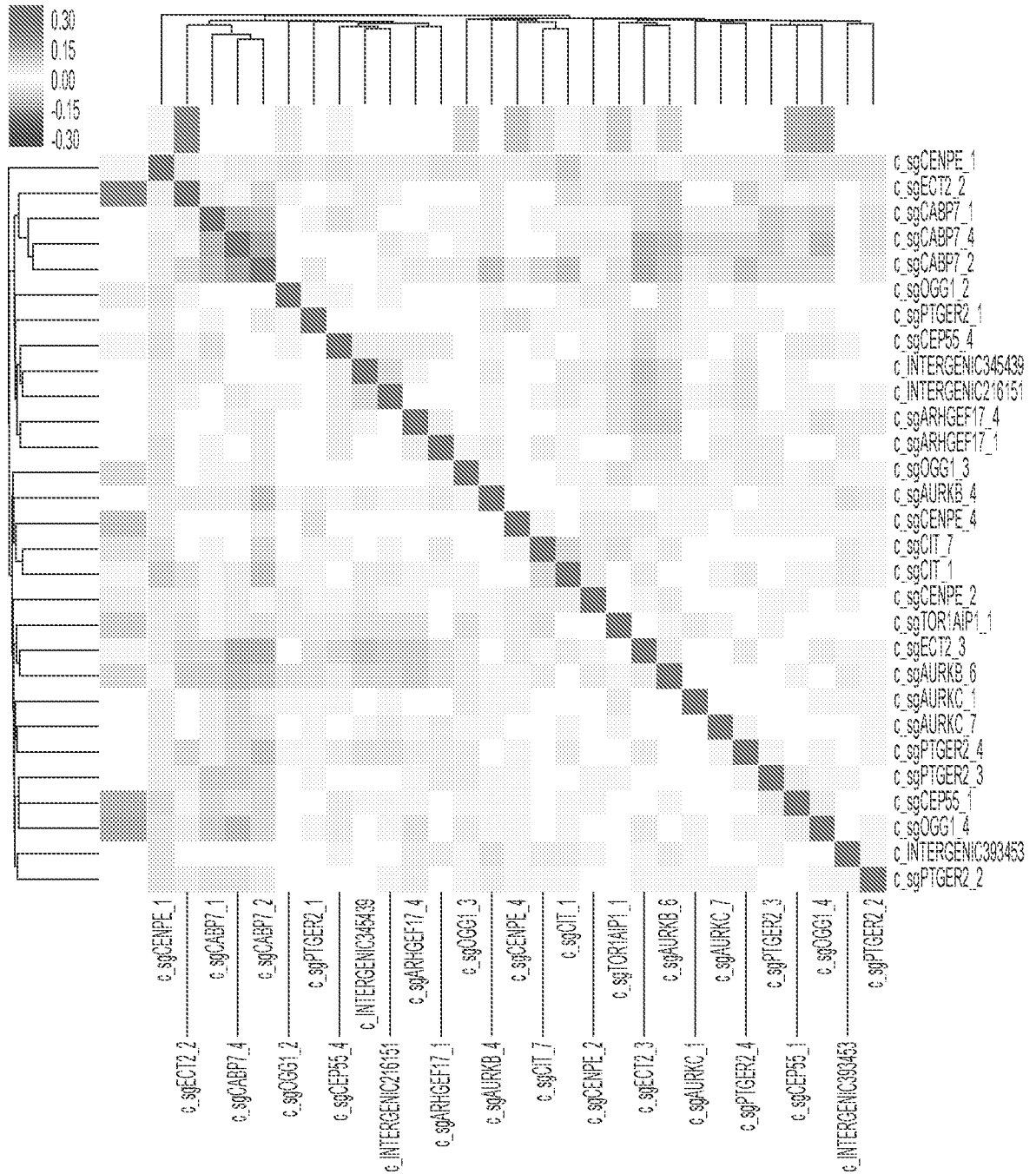


FIG. 63D

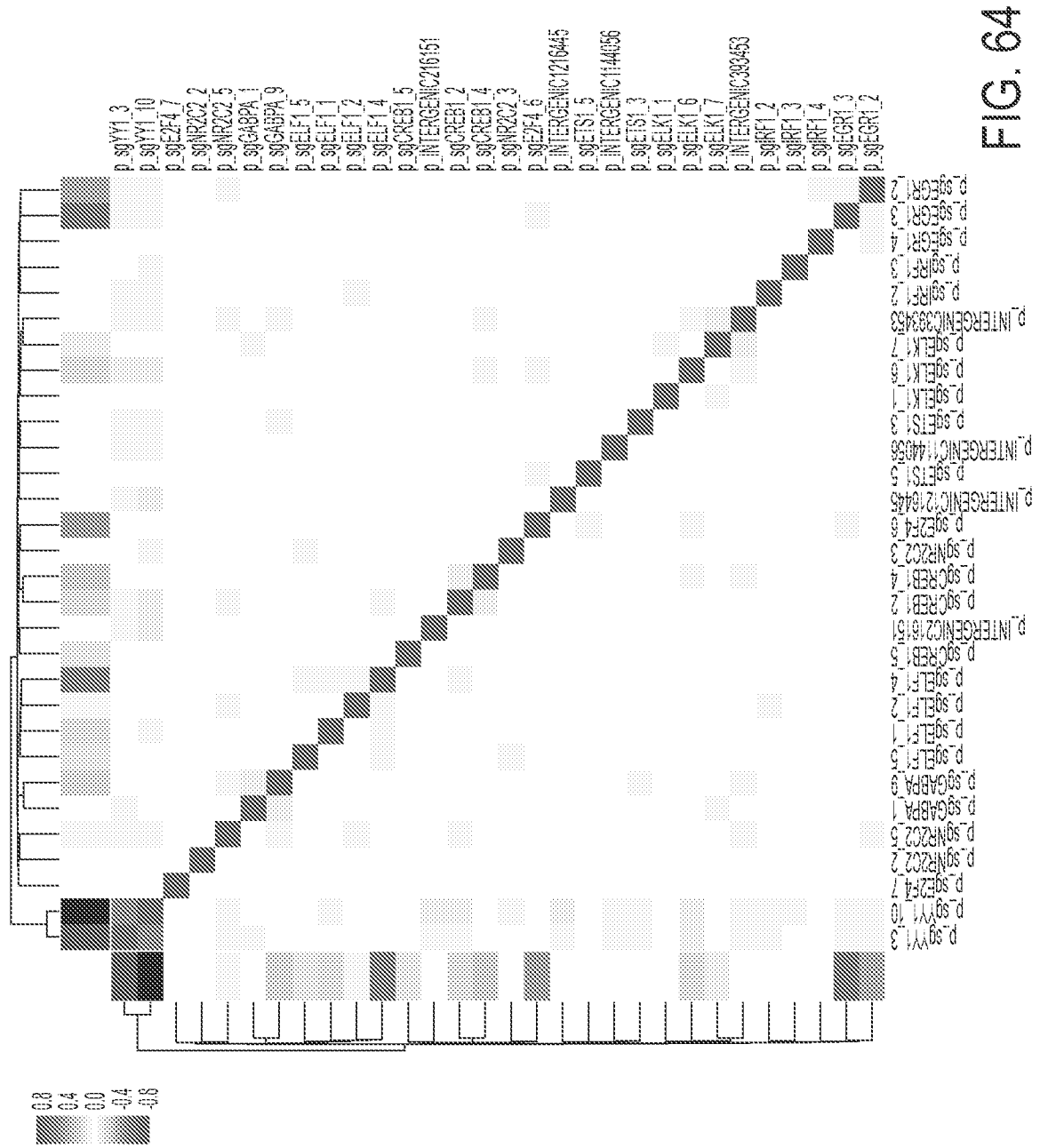
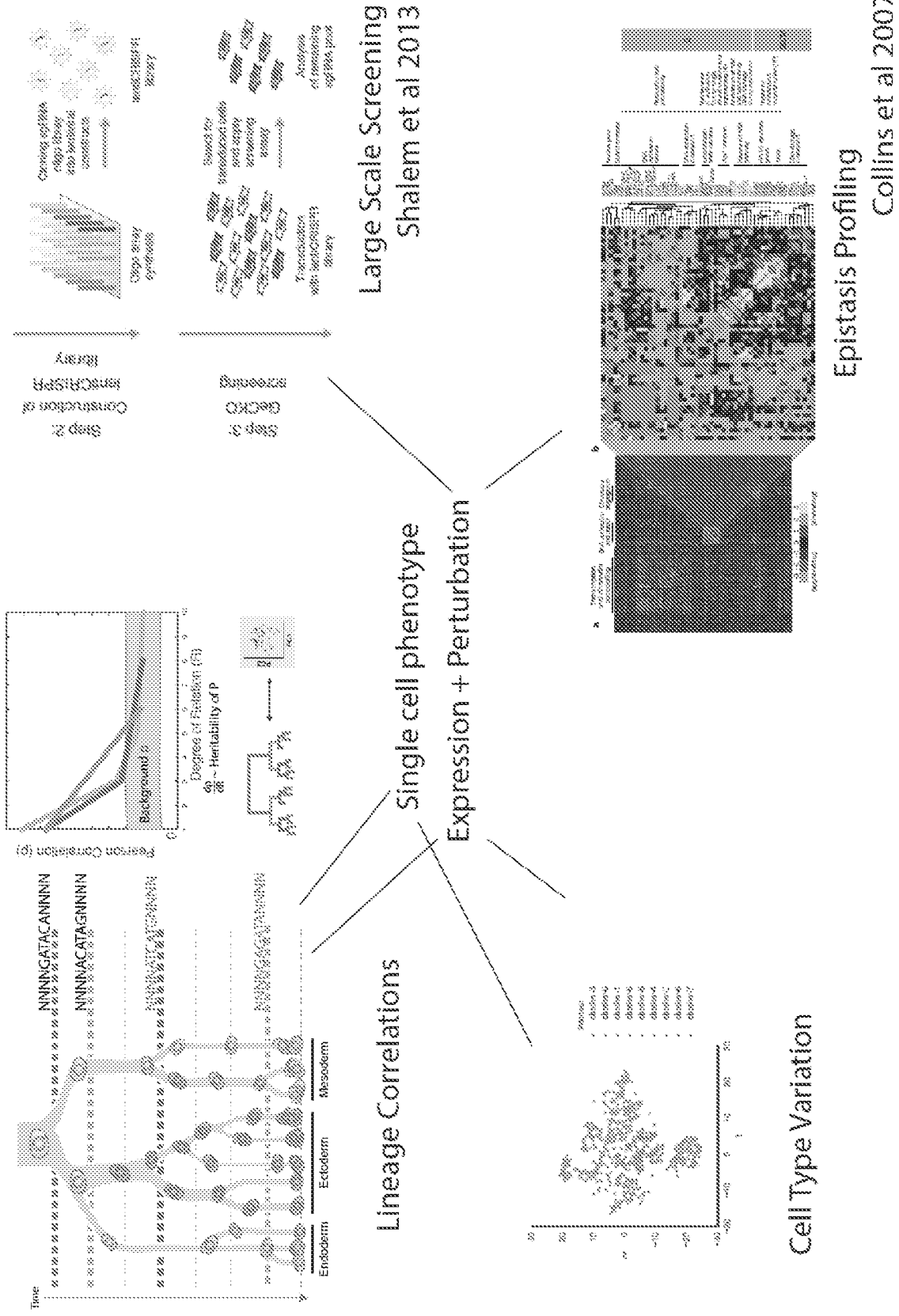


FIG. 64

Figure 65



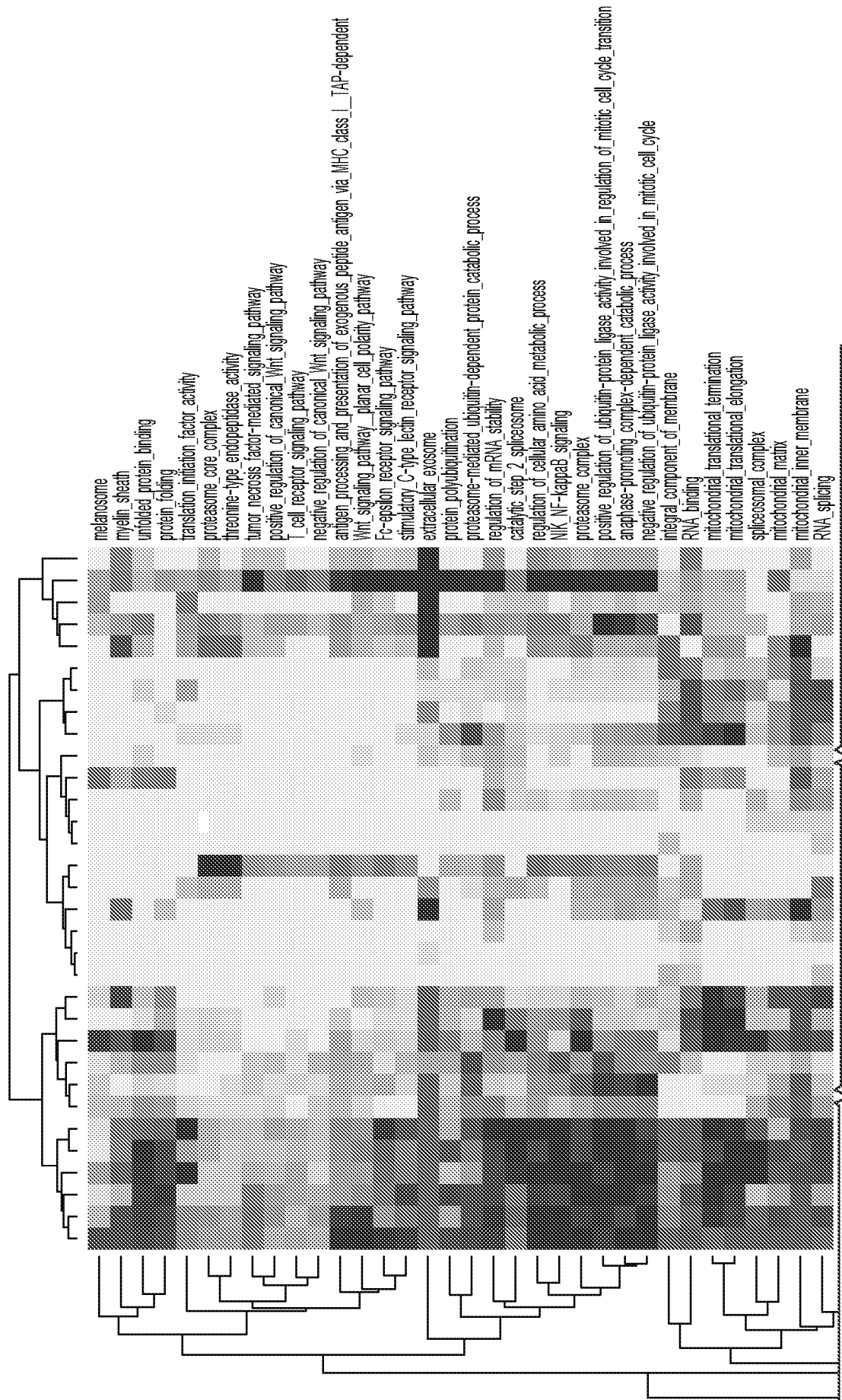


FIG. 66

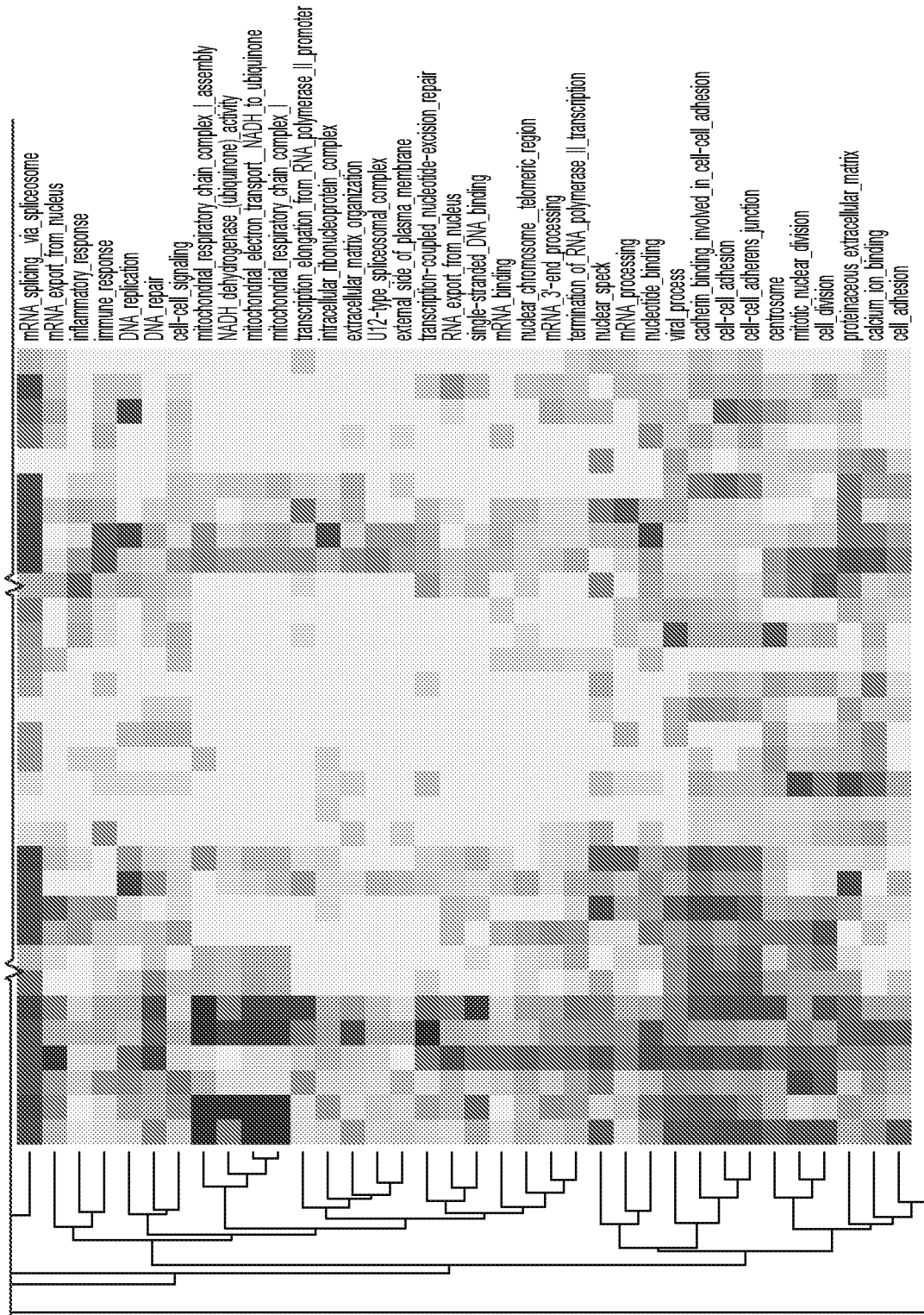


FIG. 66 continued

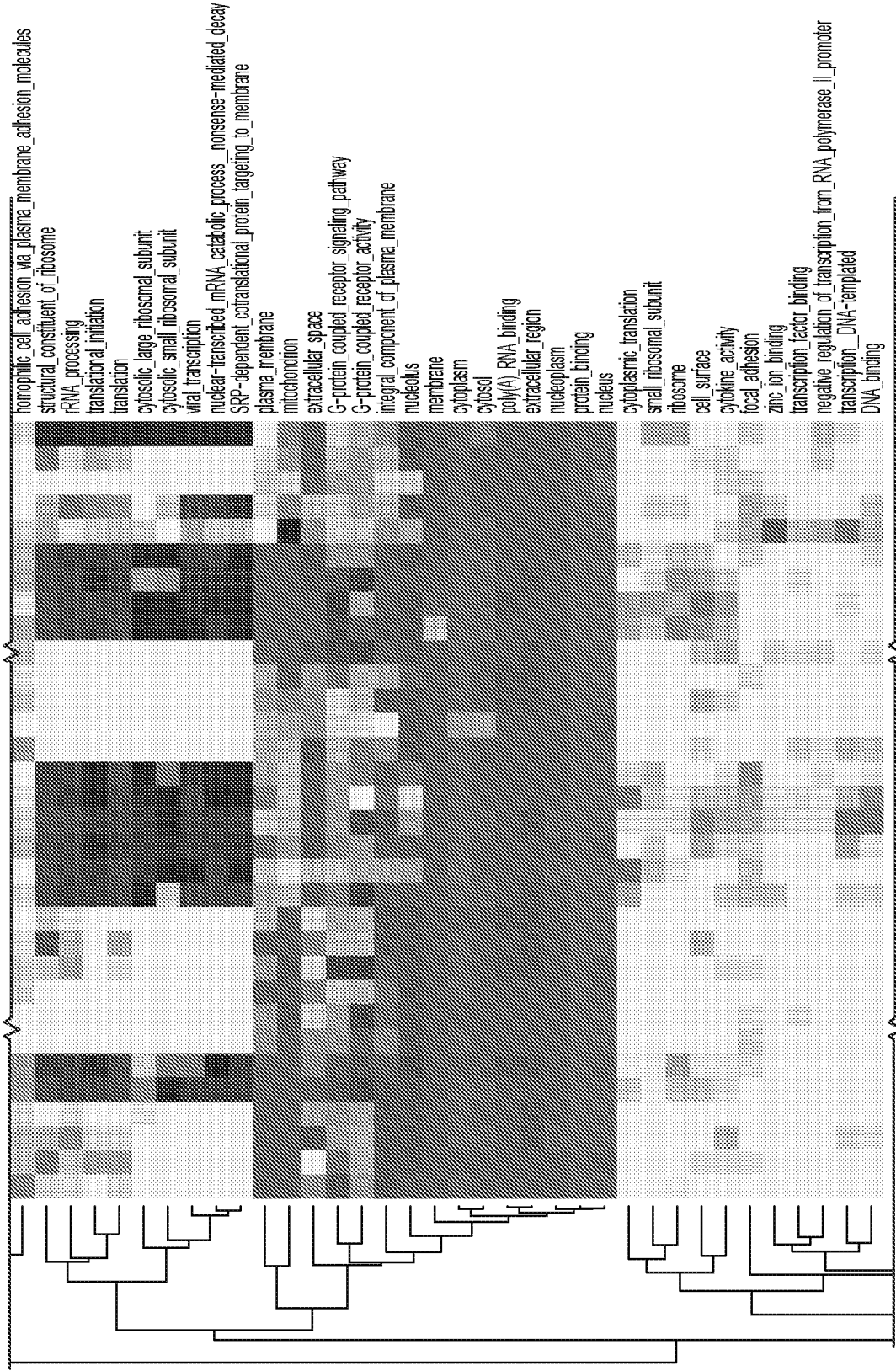


FIG. 66 continued

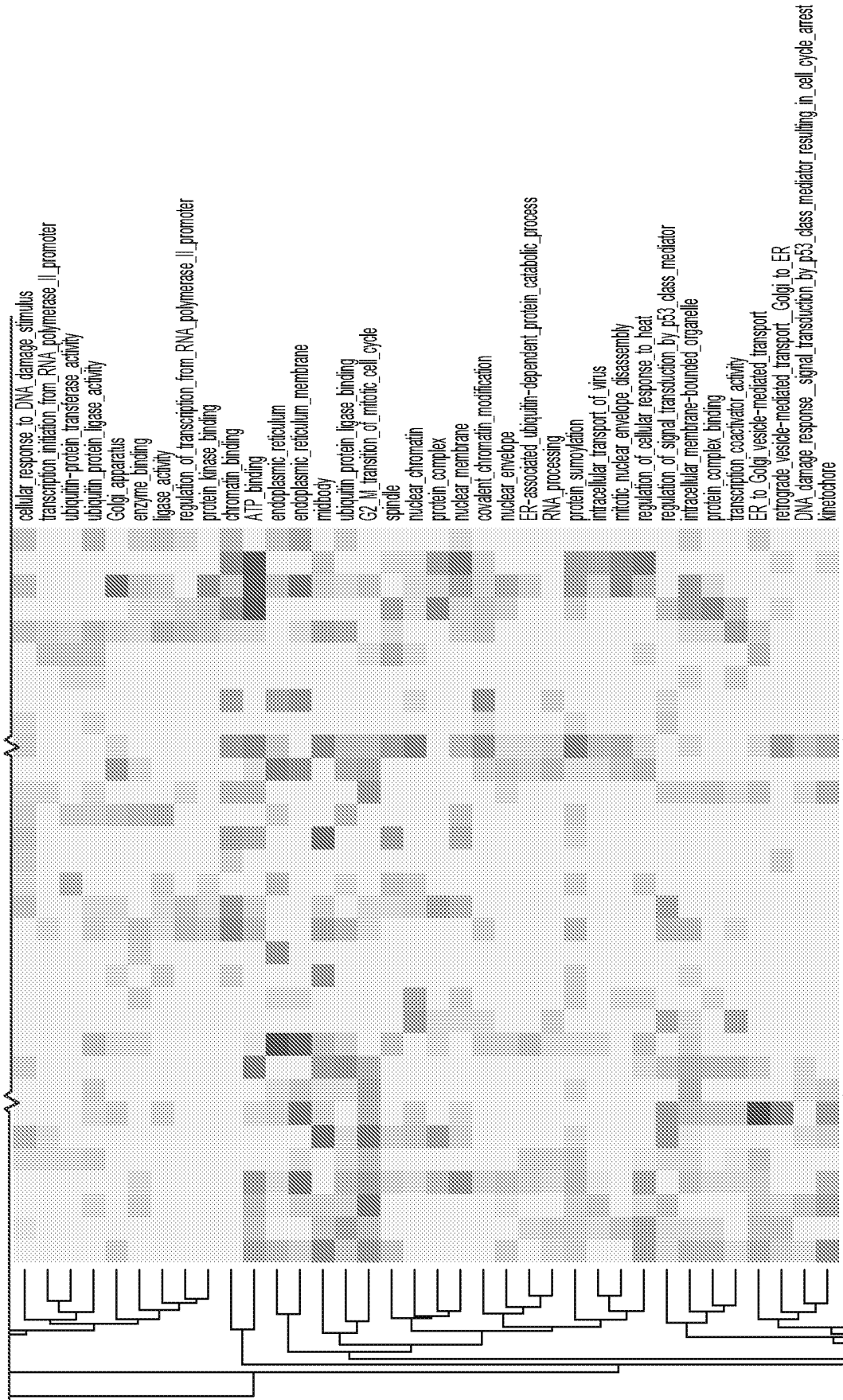


FIG. 66 continued

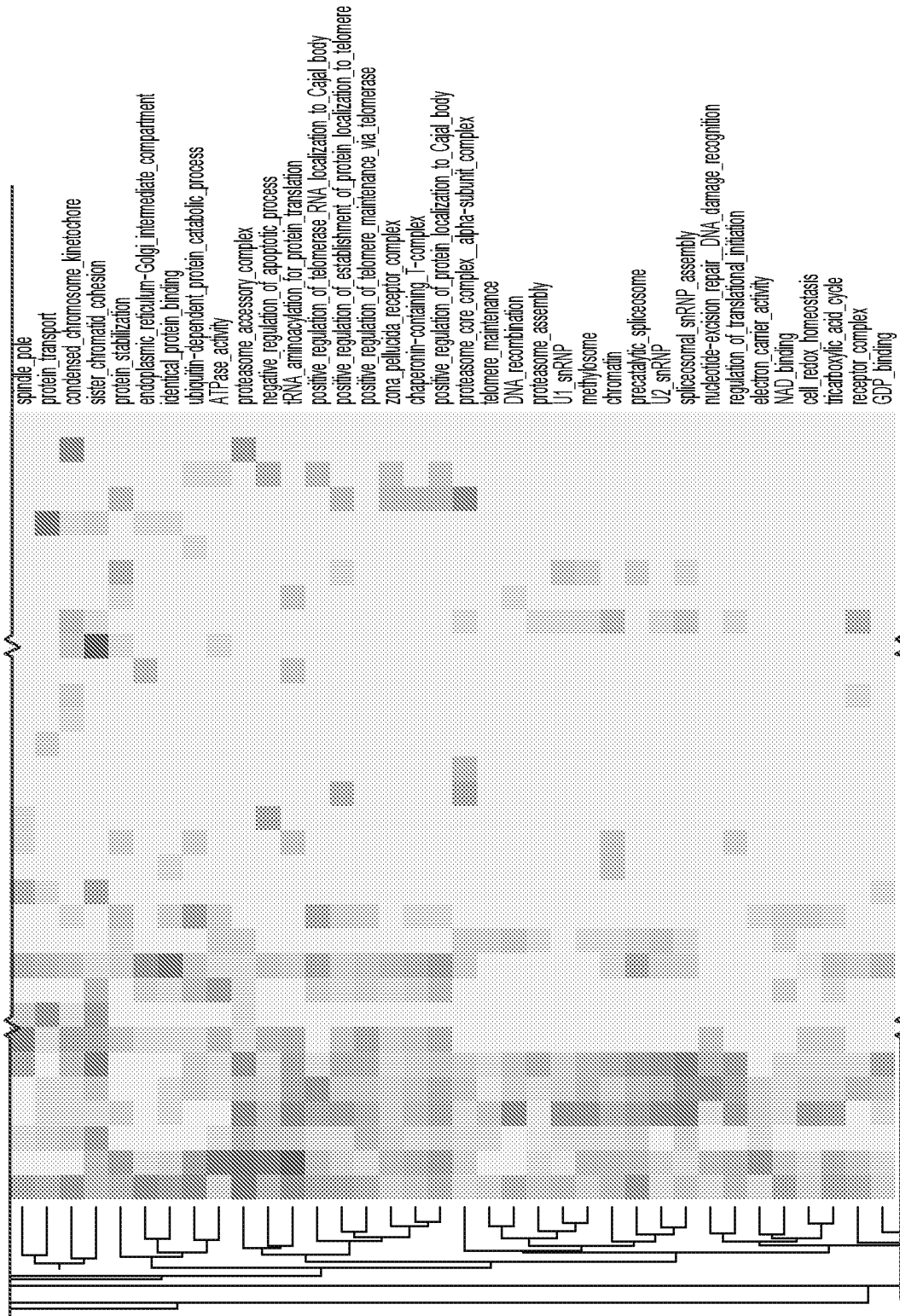


FIG. 66 continued

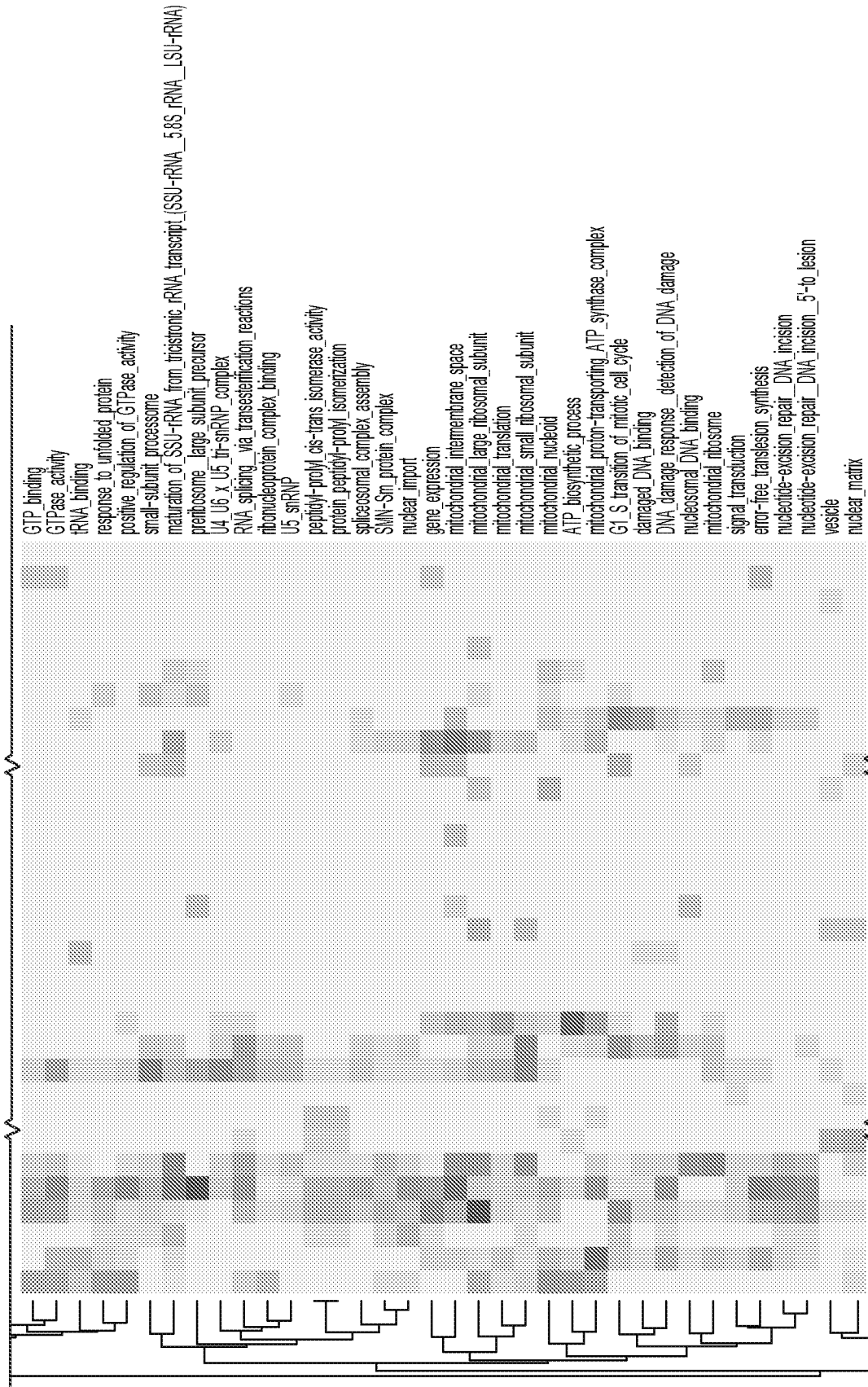


FIG. 66 continued



FIG. 66 continued

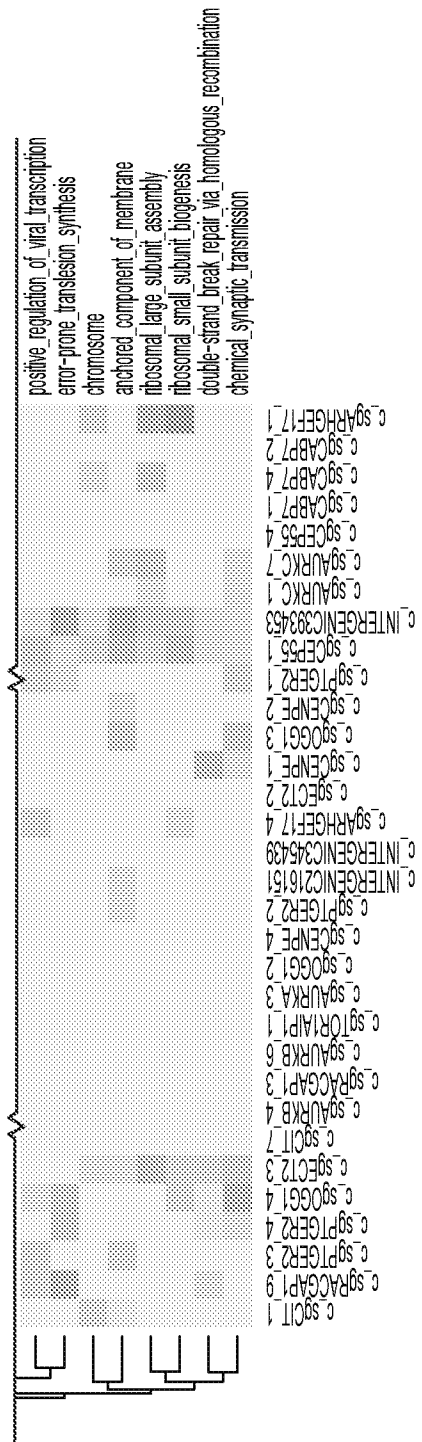


FIG. 66 continued

Figure 67

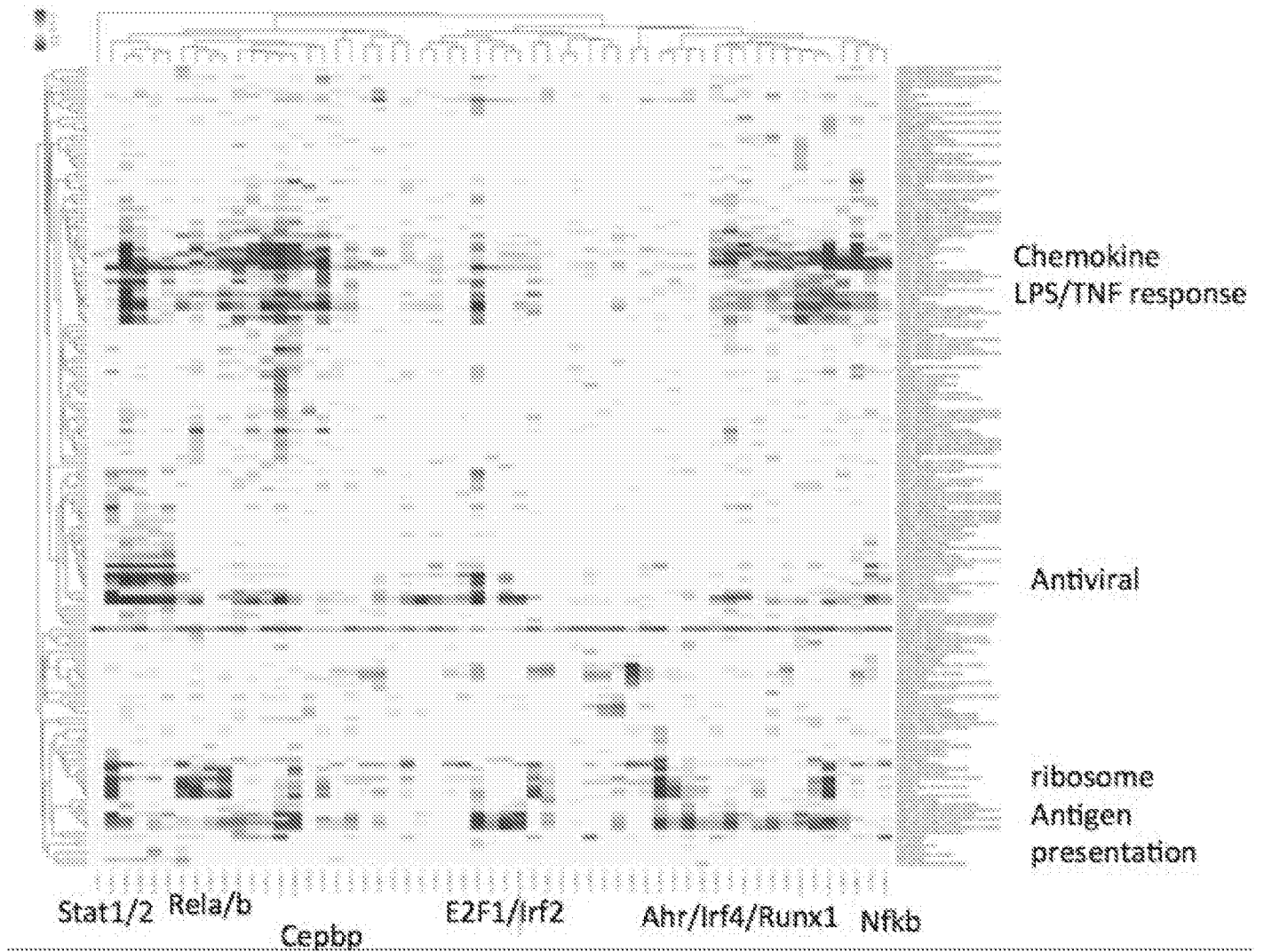


Figure 68

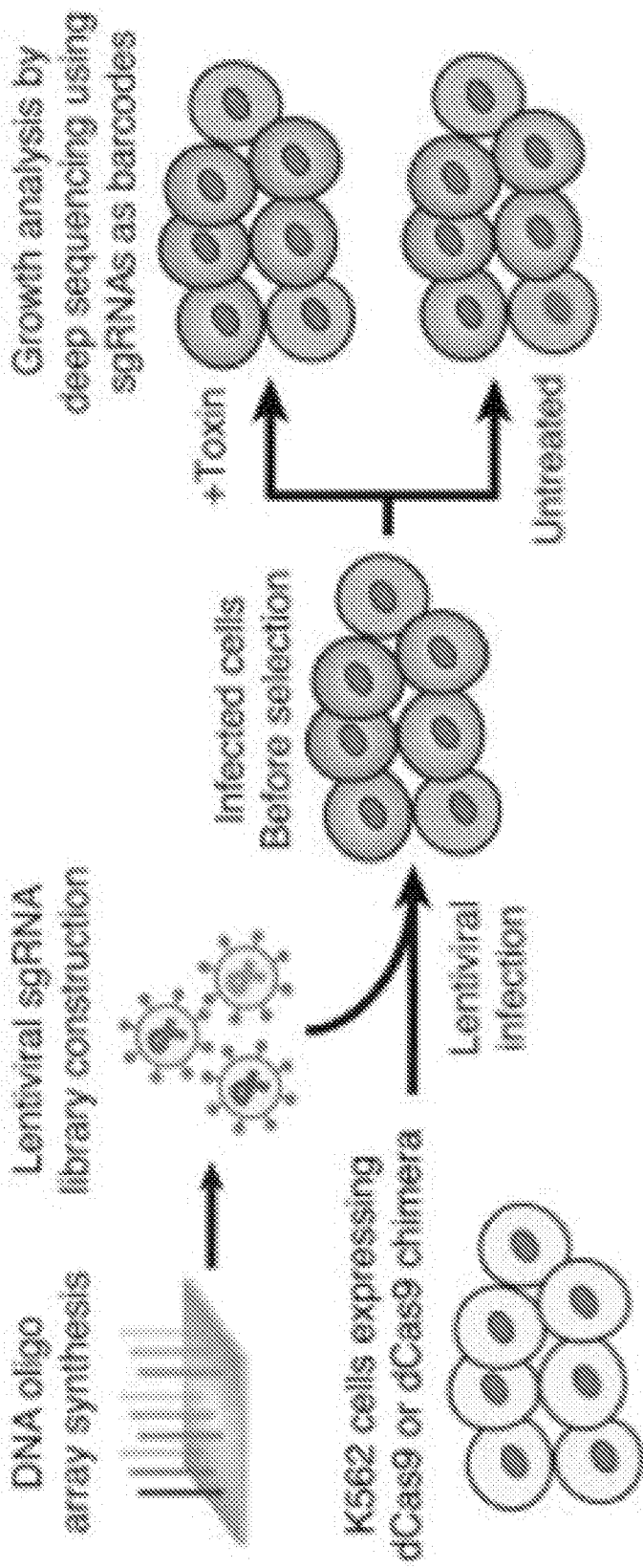


Figure 69

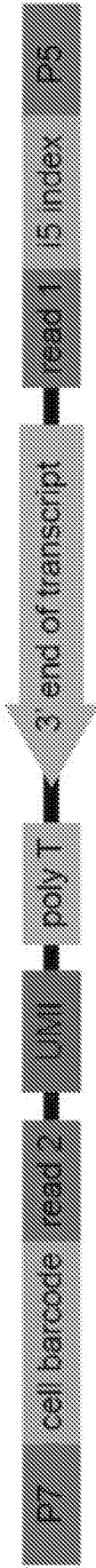


Figure 70

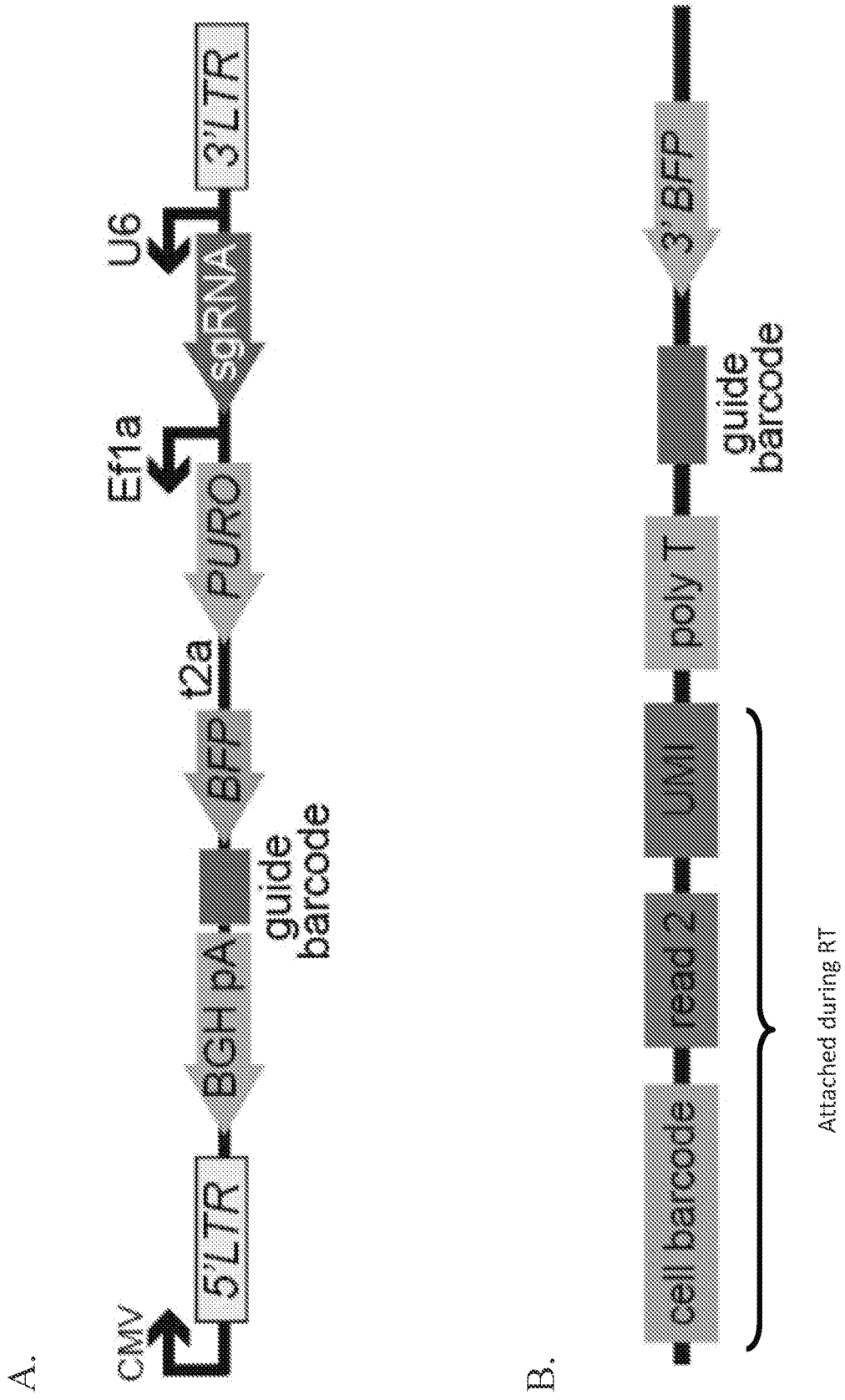


Figure 71

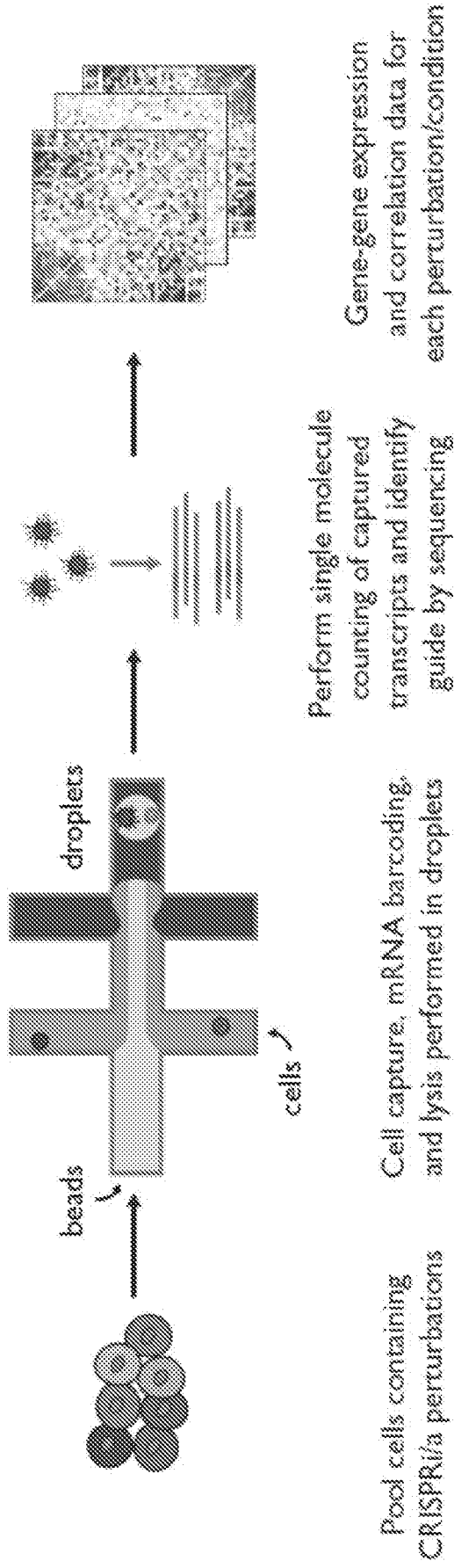


Figure 73

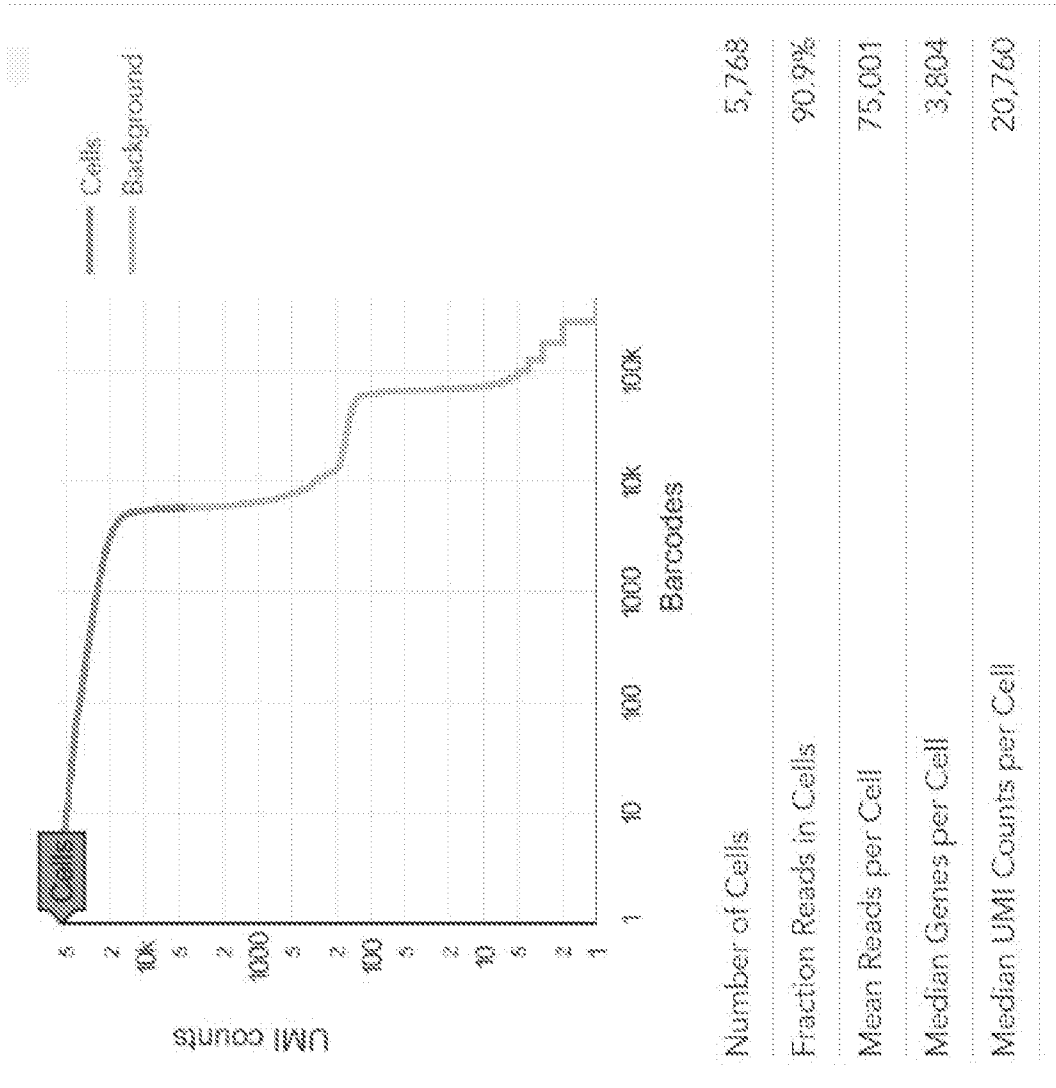


Figure 74

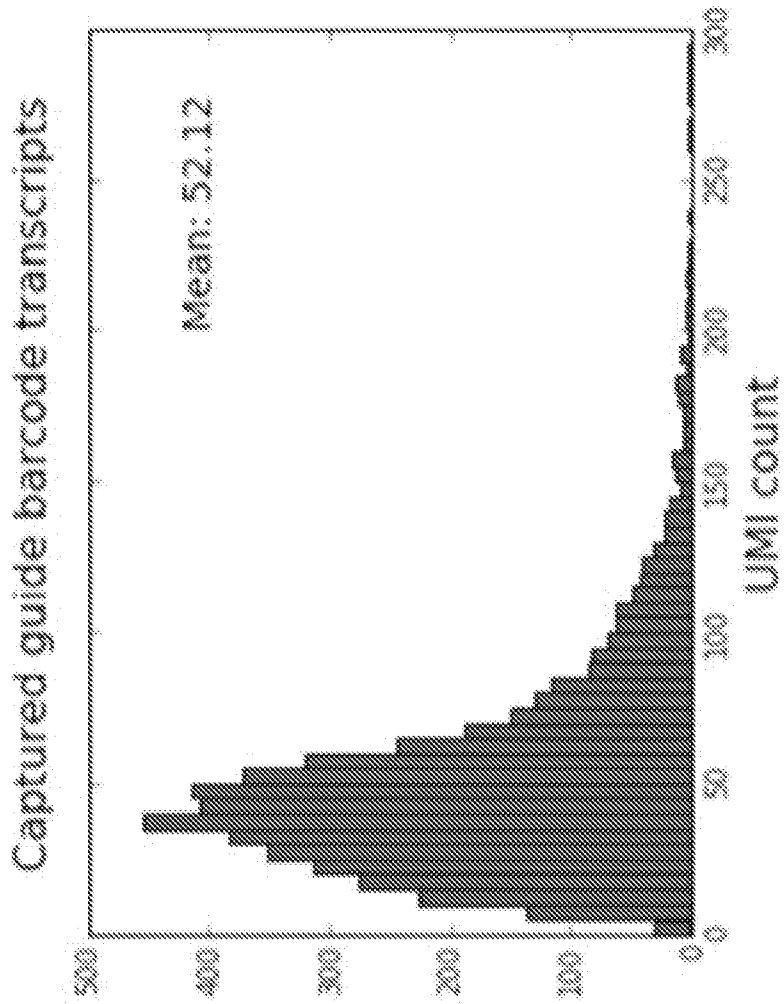
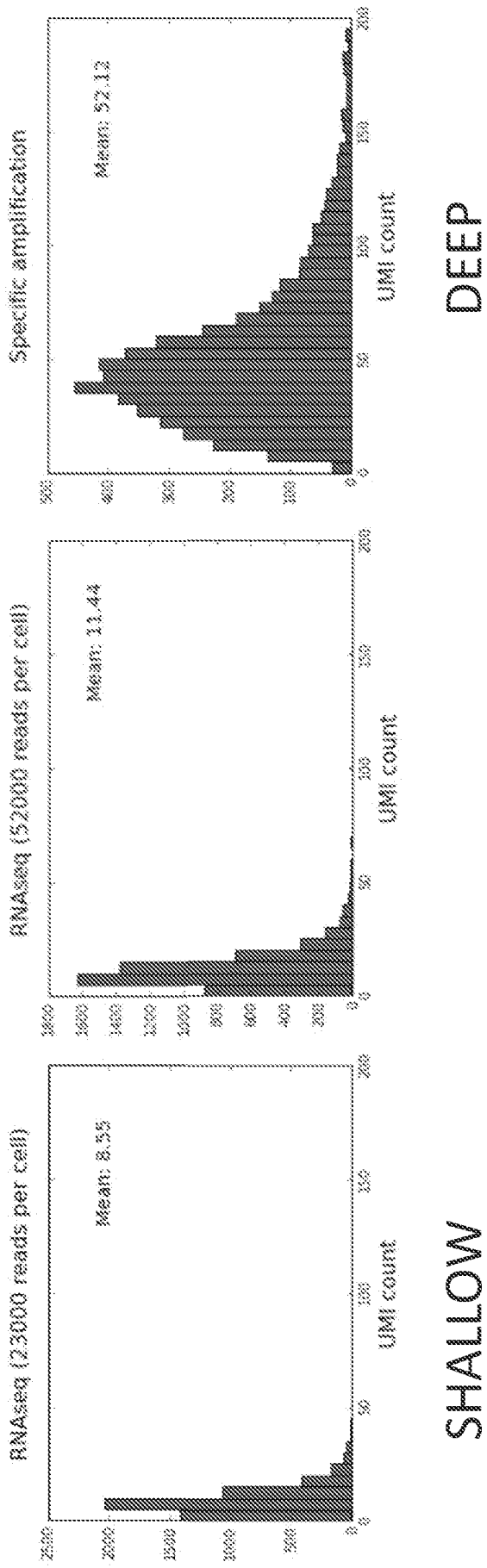


Figure 75



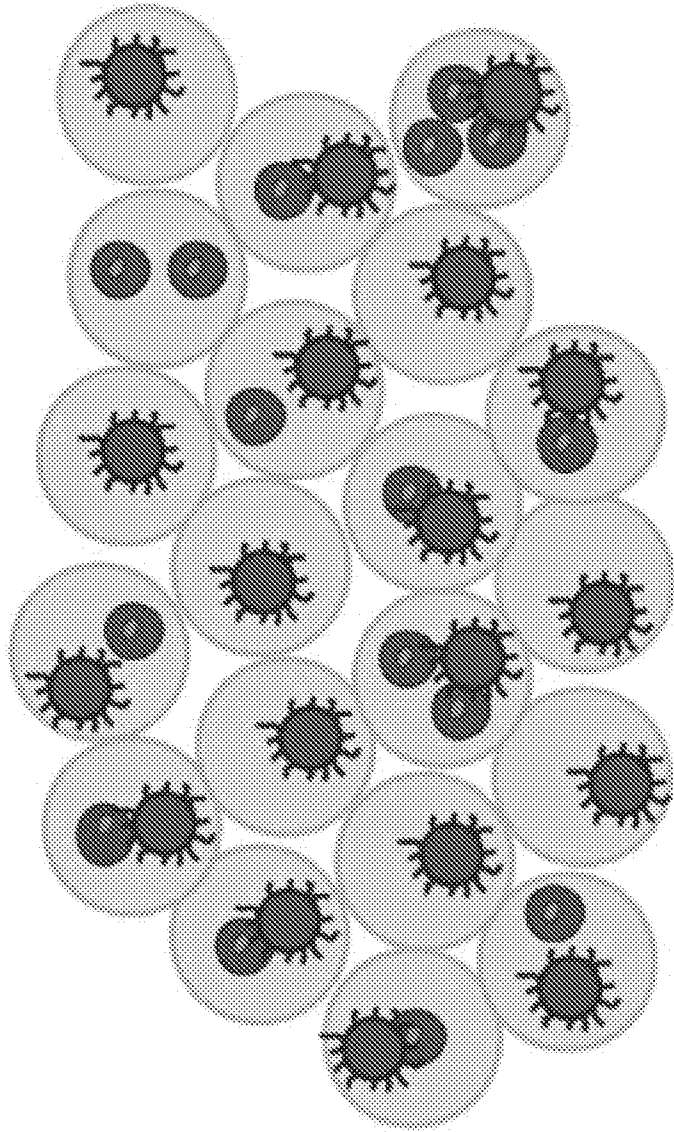
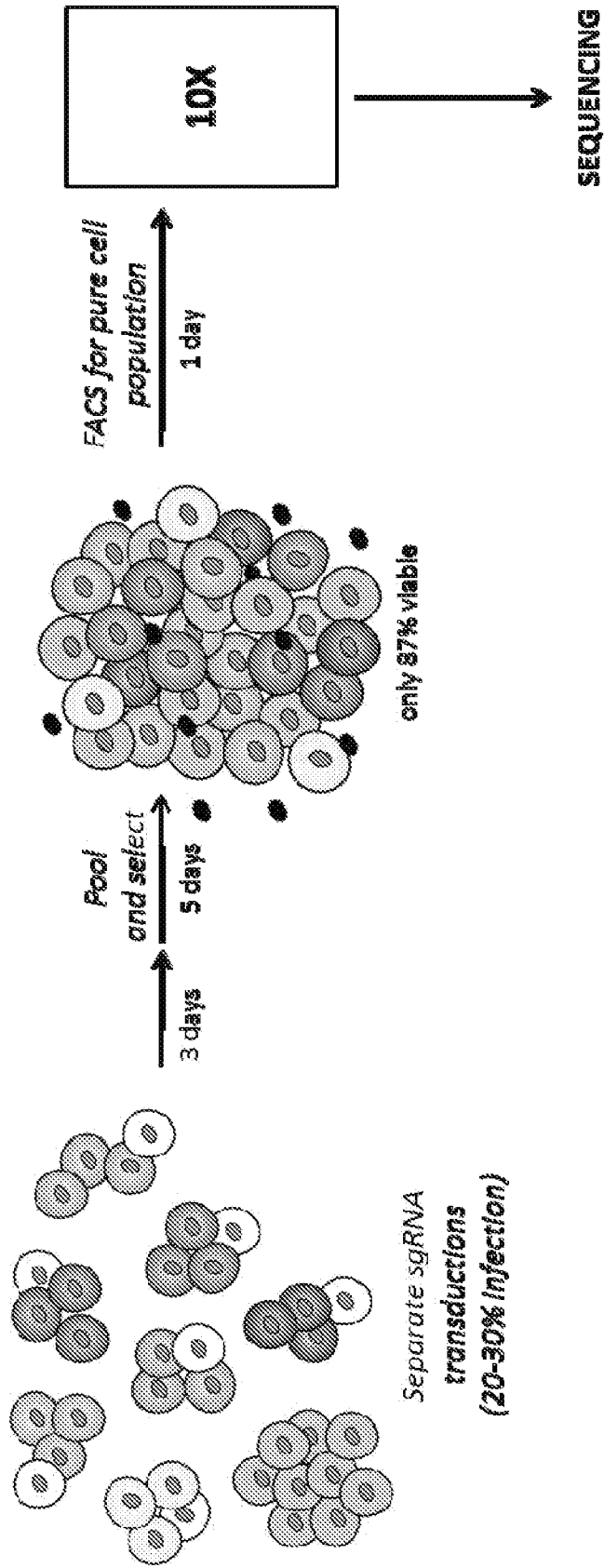


Figure 76

Figure 77



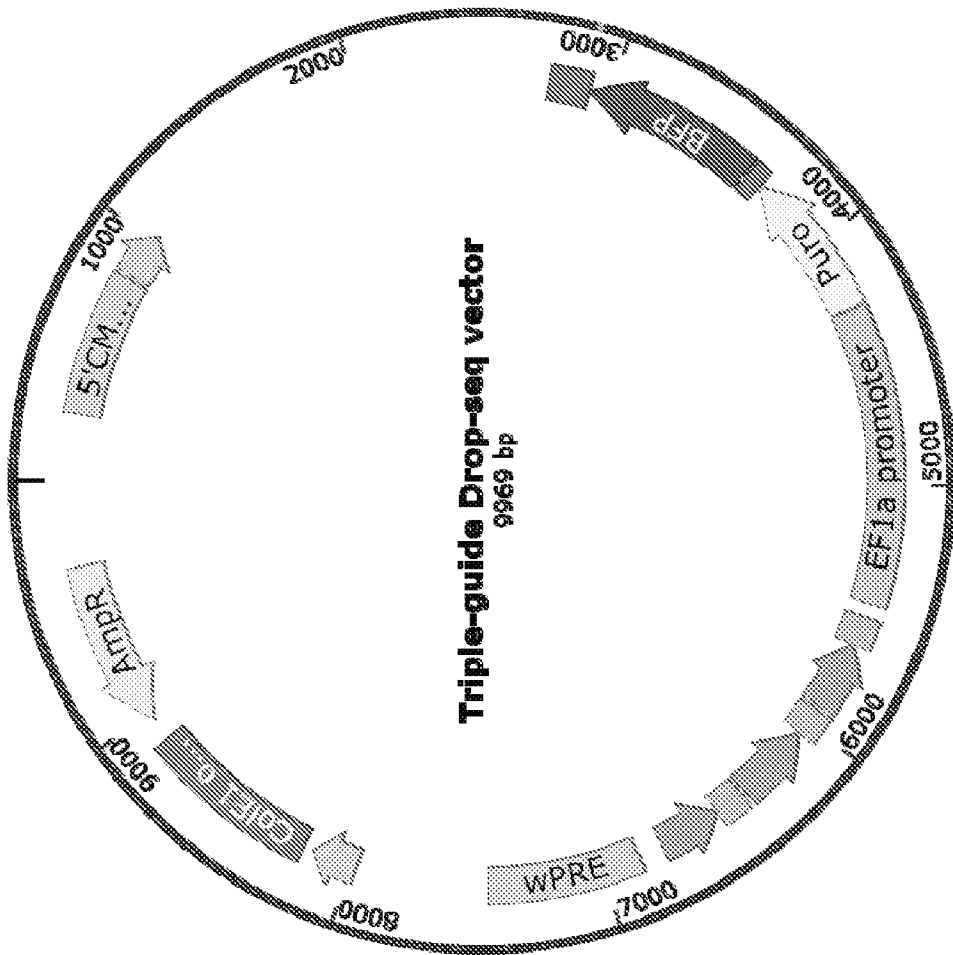


Figure 78

Figure 79

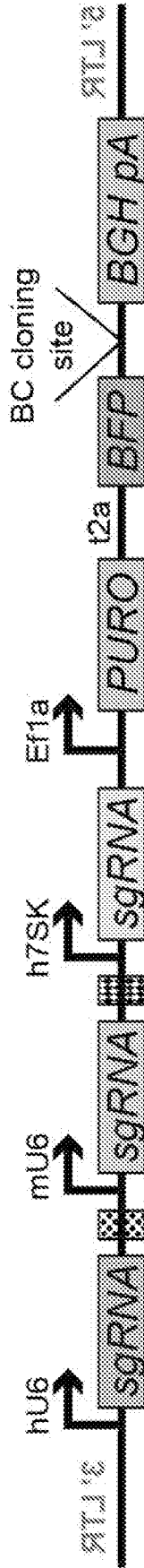


Figure 80

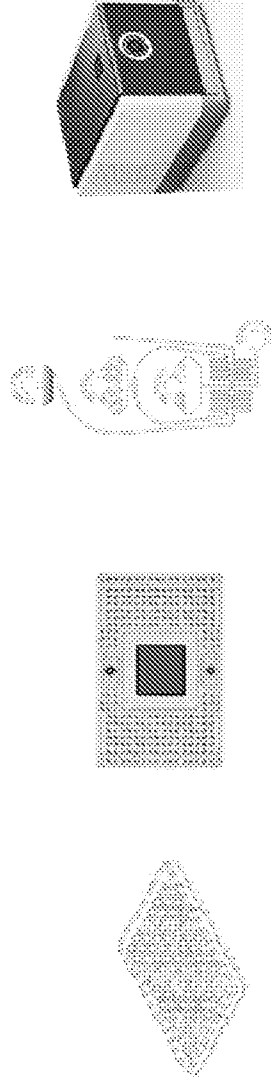


	Plate-based	Fluidigm CI Autoprep	Drop-seq/InDrop	10X GemCode
Capture method	FACS into plate	Microfluidic wells	Random droplet encapsulation	Random droplet encapsulation
Throughput	96 per plate	96 (average: 72) or 800 per chip	~5000 – 10000 per hour	1000 – 6000 per 10 minutes
Volume "scale"	4 µL to 25 µL	4.5 nL to 5 µL	~1 nL to 50 nL	
mRNA capture efficiency	~20-40%	~48%	12.8% (DropSeq) 7% (InDrop)	Comparable to dropseq
Cost per cell	~\$10	\$3.50	\$0.07 (DropSeq)	

FIG. 81A

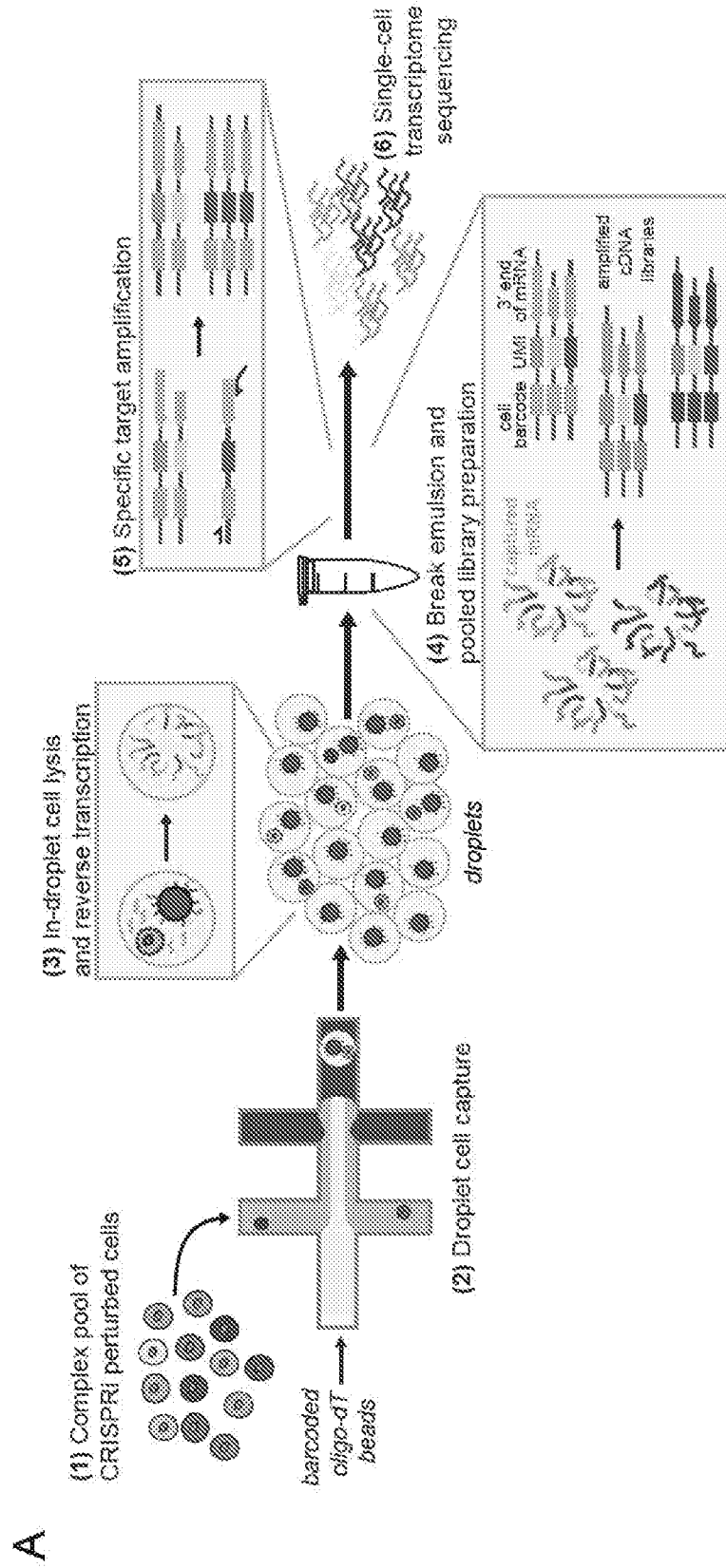


FIG. 81B

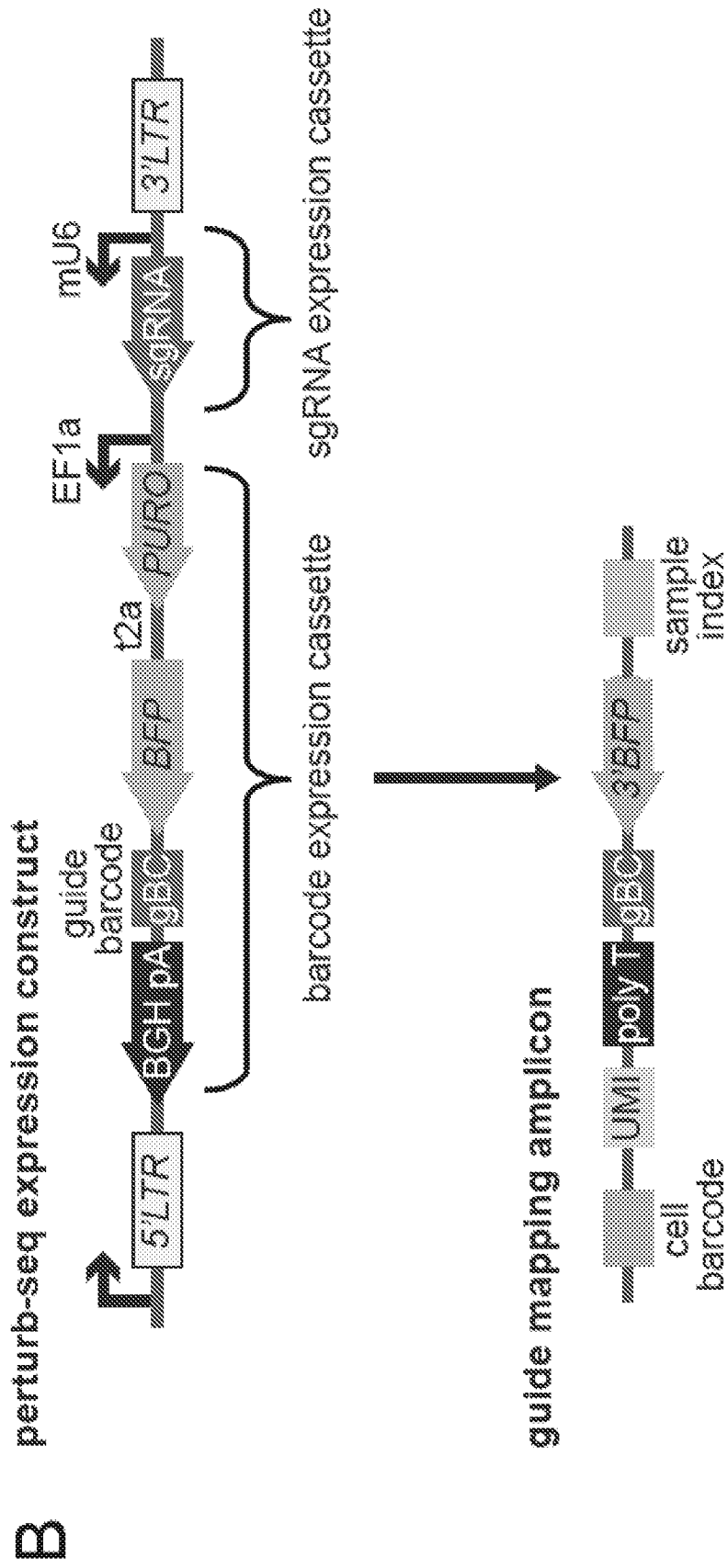


FIG. 81C-F

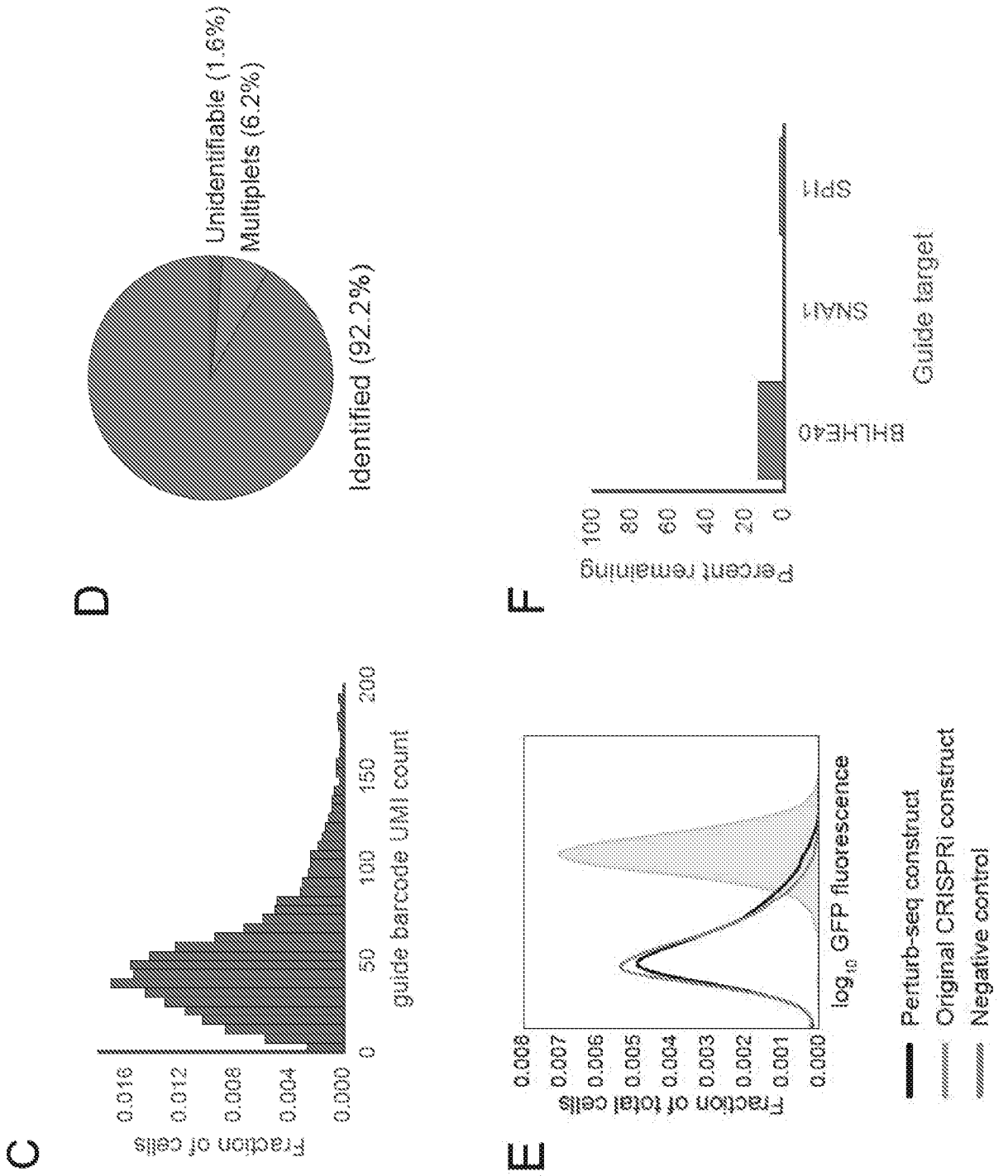


FIG. 82A-C

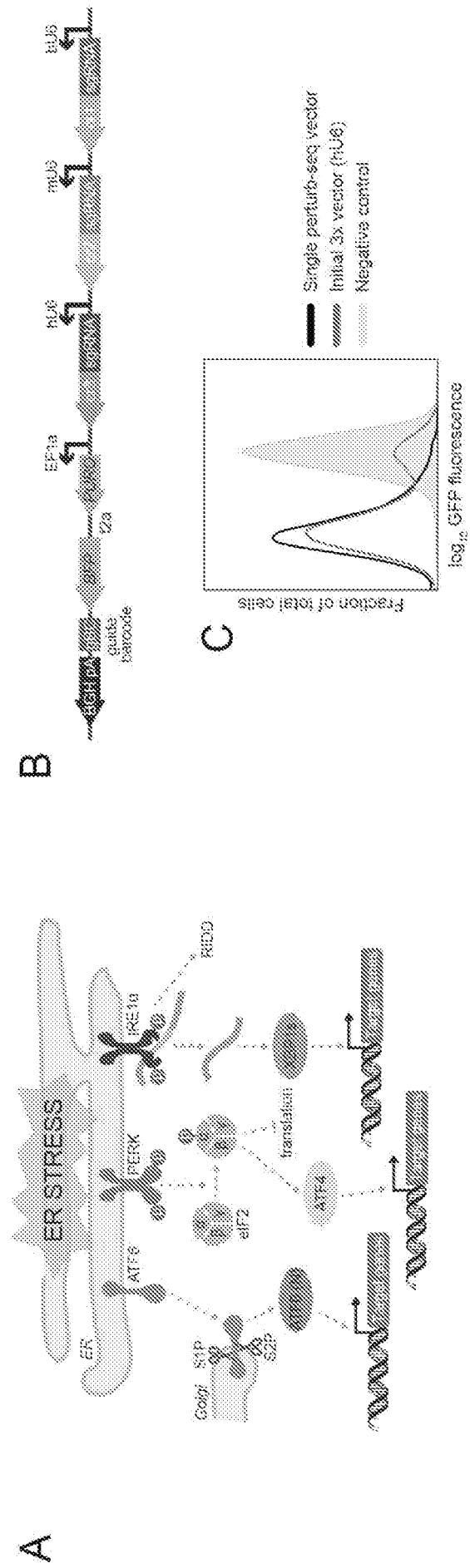


FIG. 82D

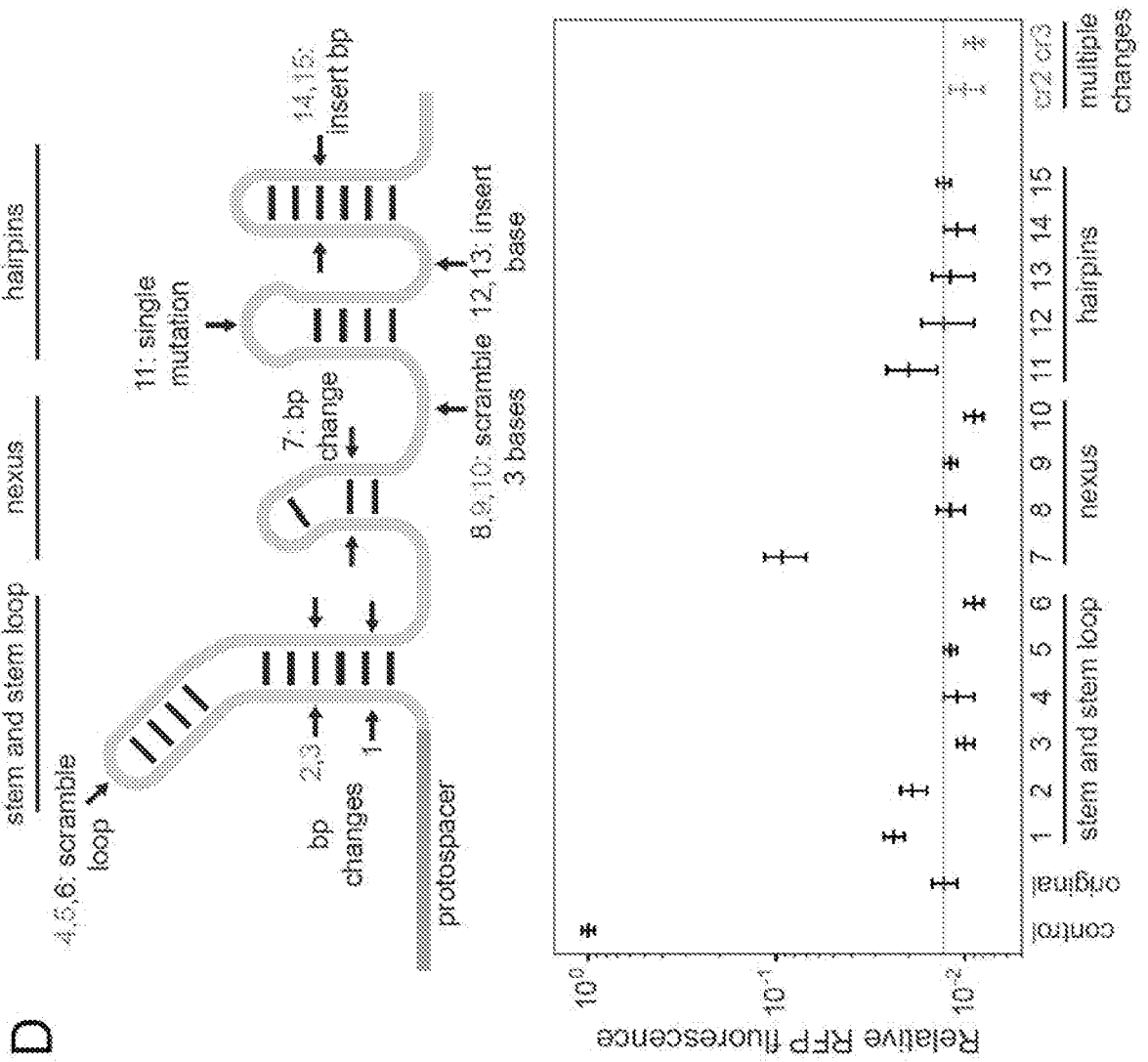


FIG. 82E-G

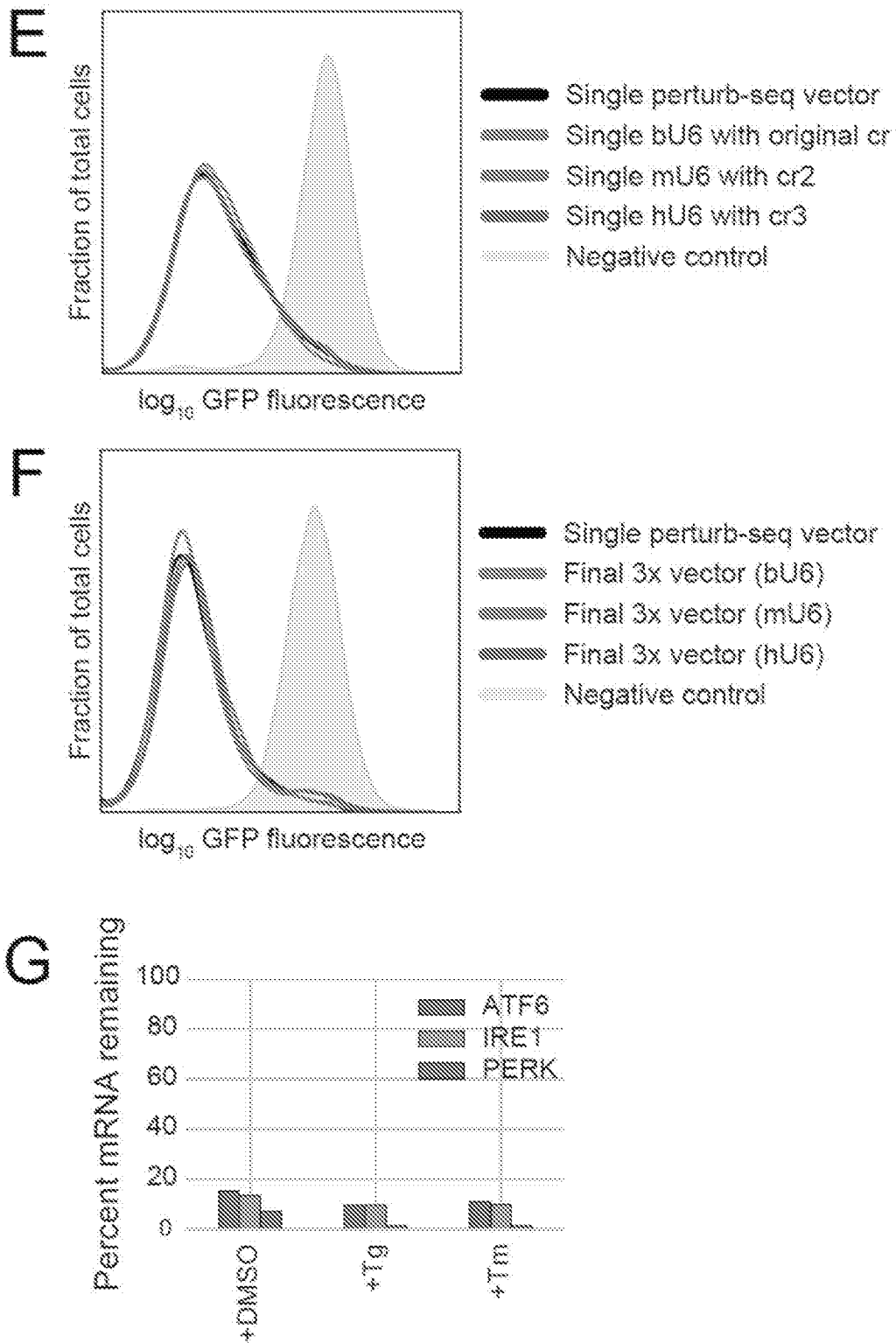


FIG. 83A-C

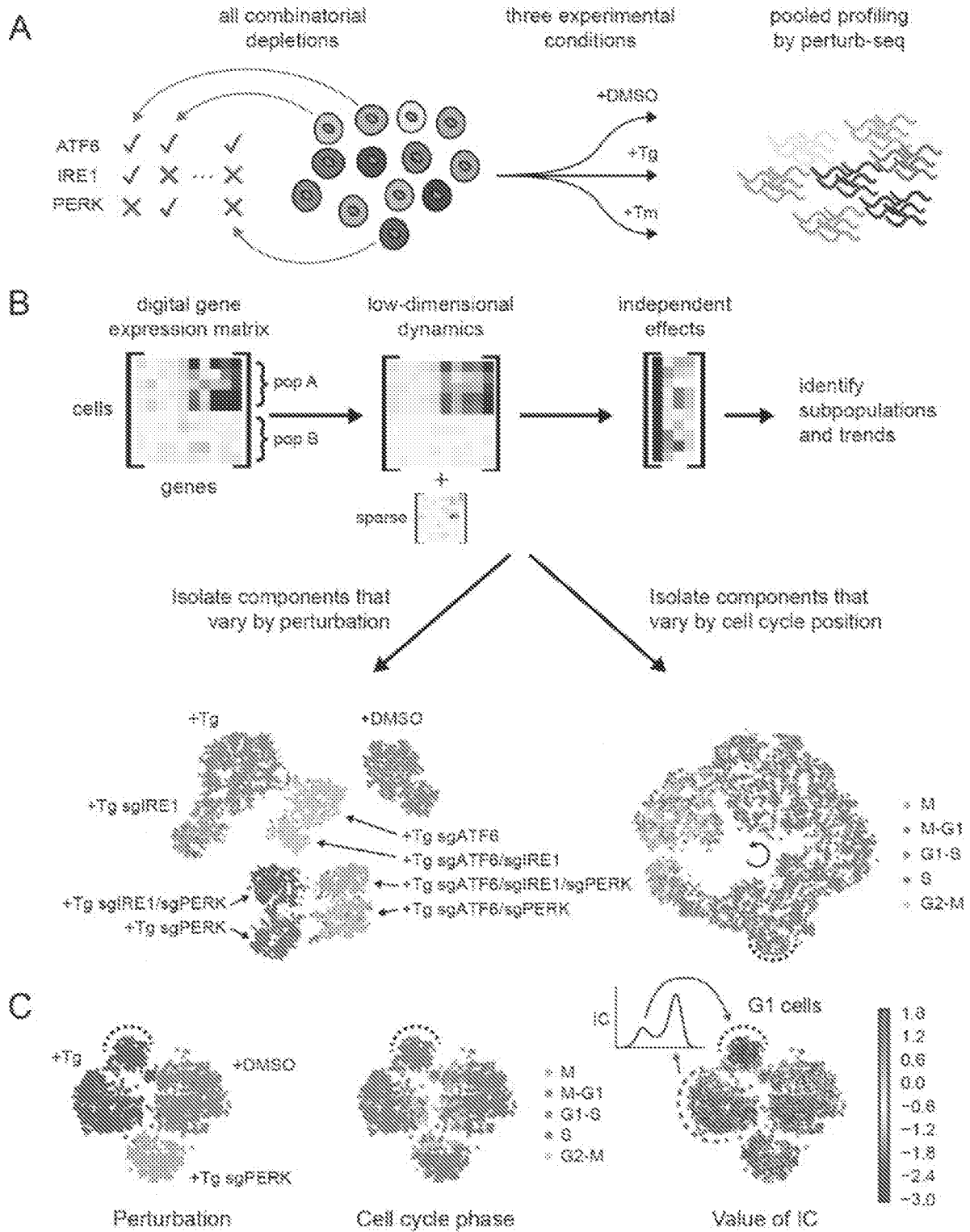


FIG. 83D

D

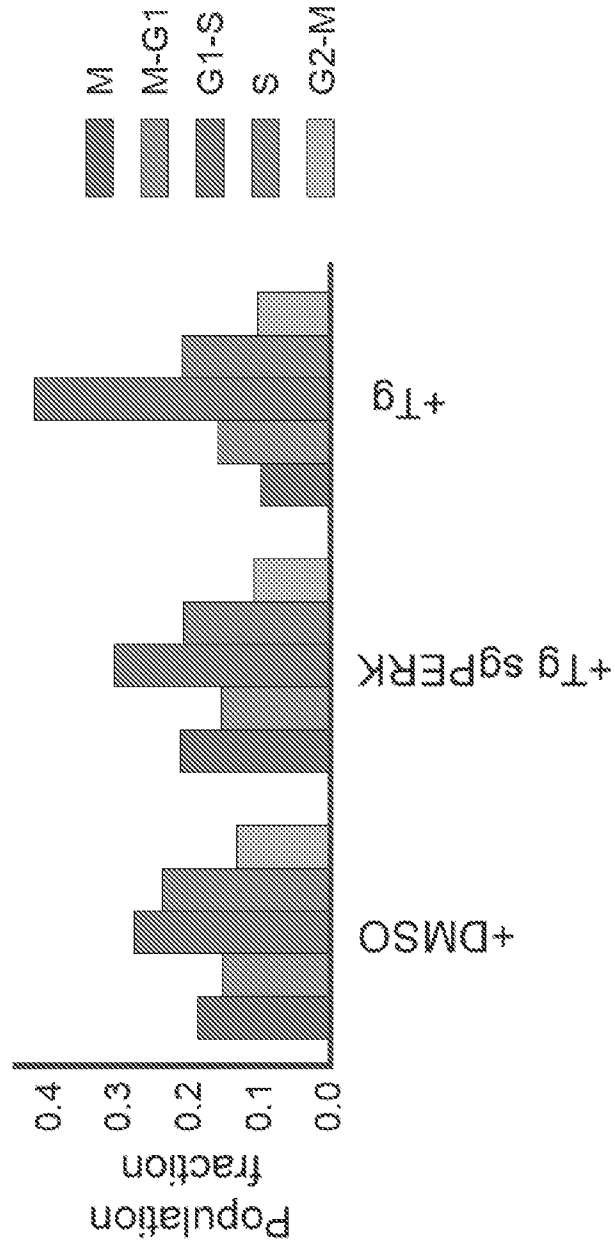


FIG. 83E

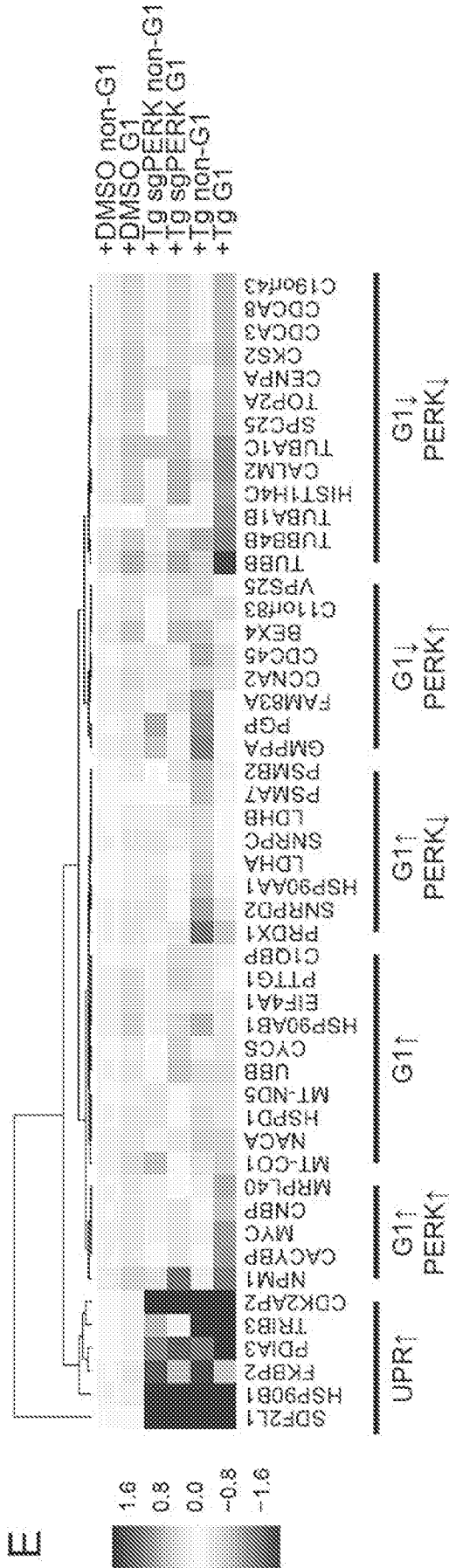


FIG. 84A-C

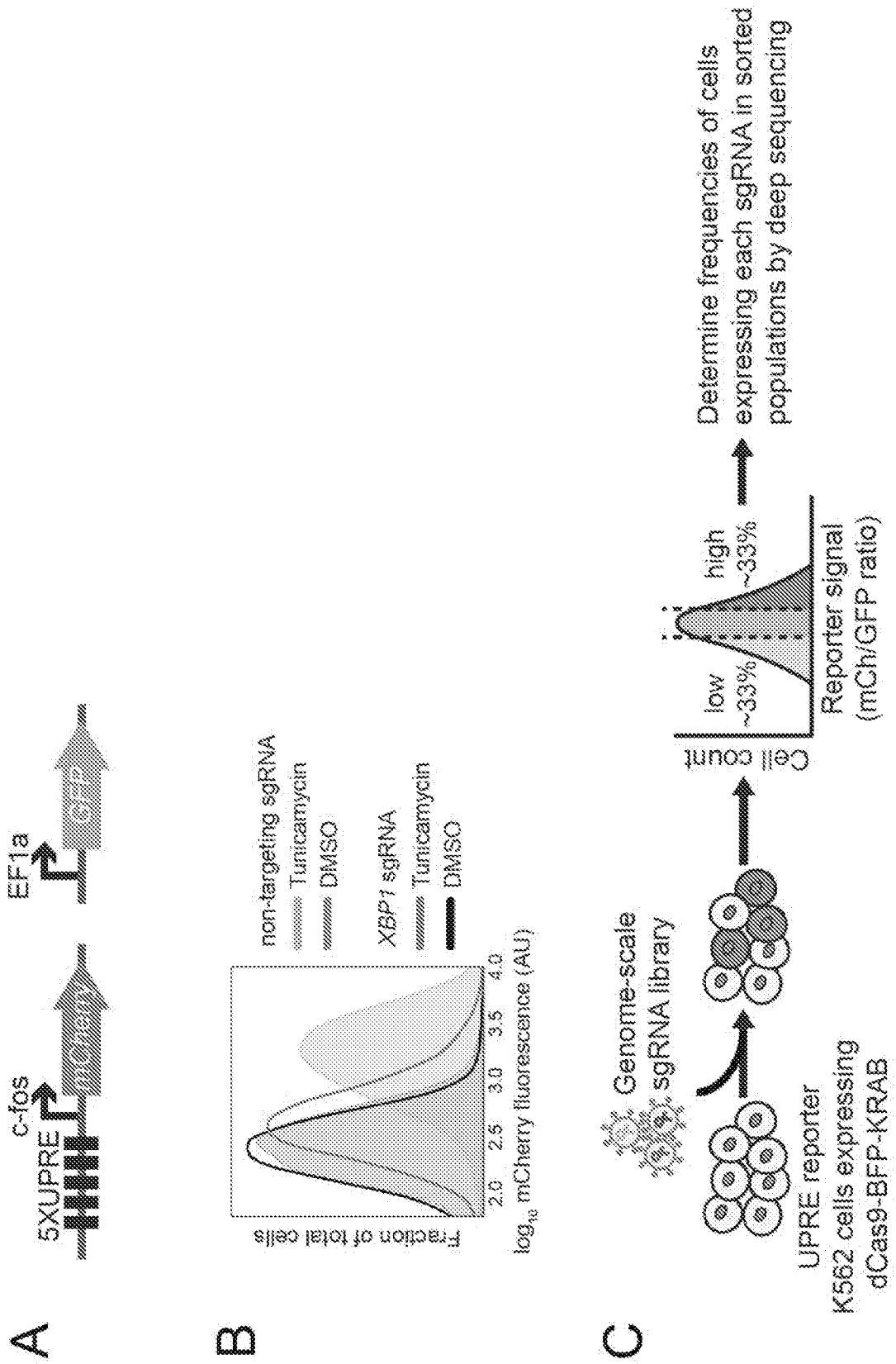


FIG. 84D-E

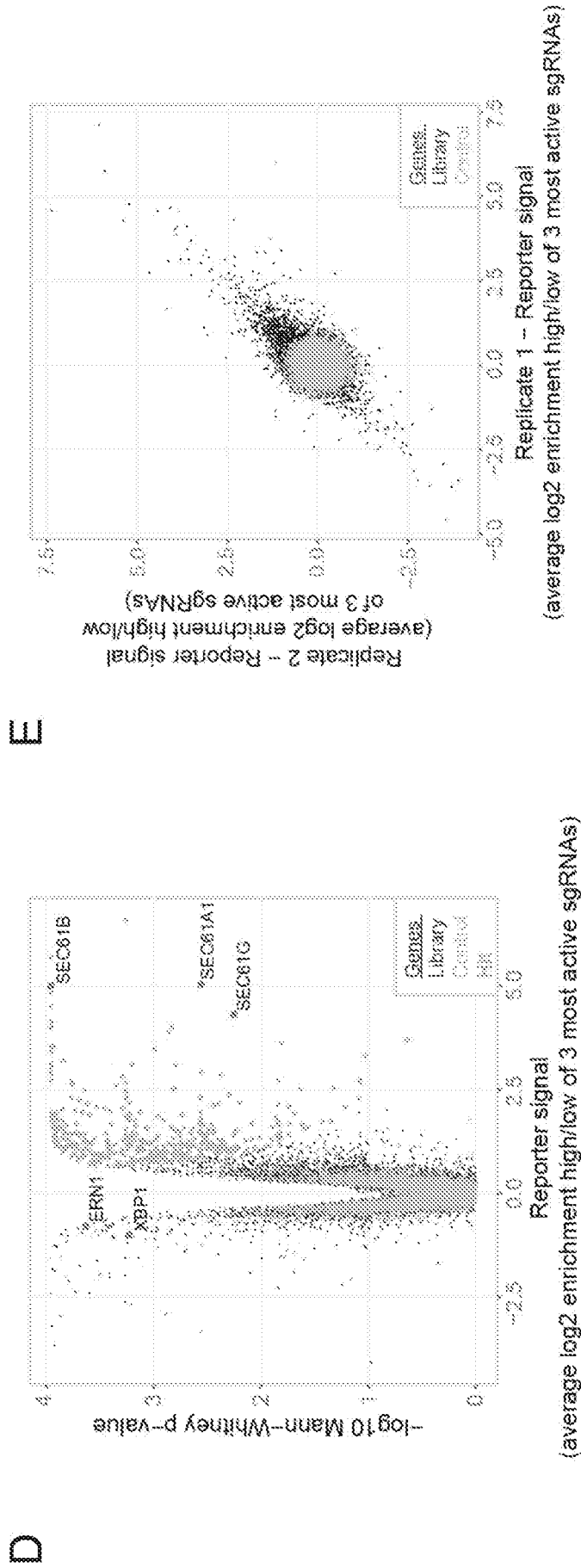


FIG. 84G

G

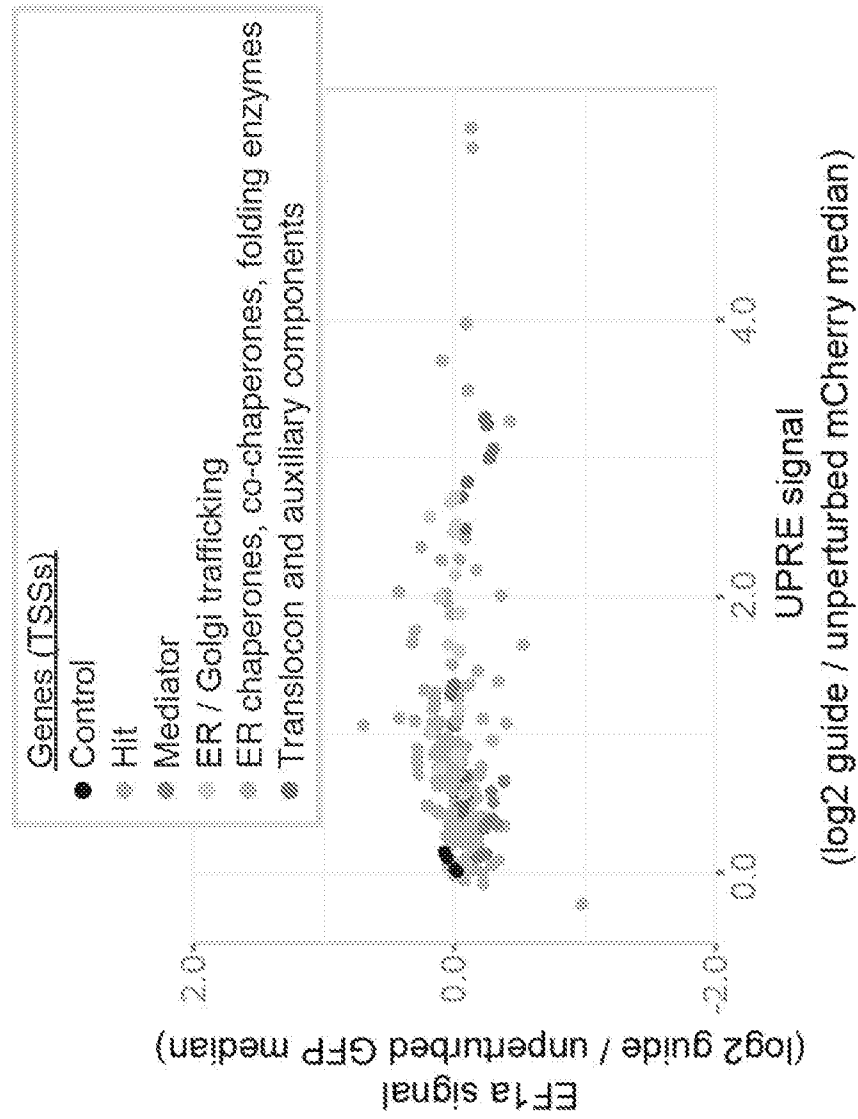


FIG. 85B-D

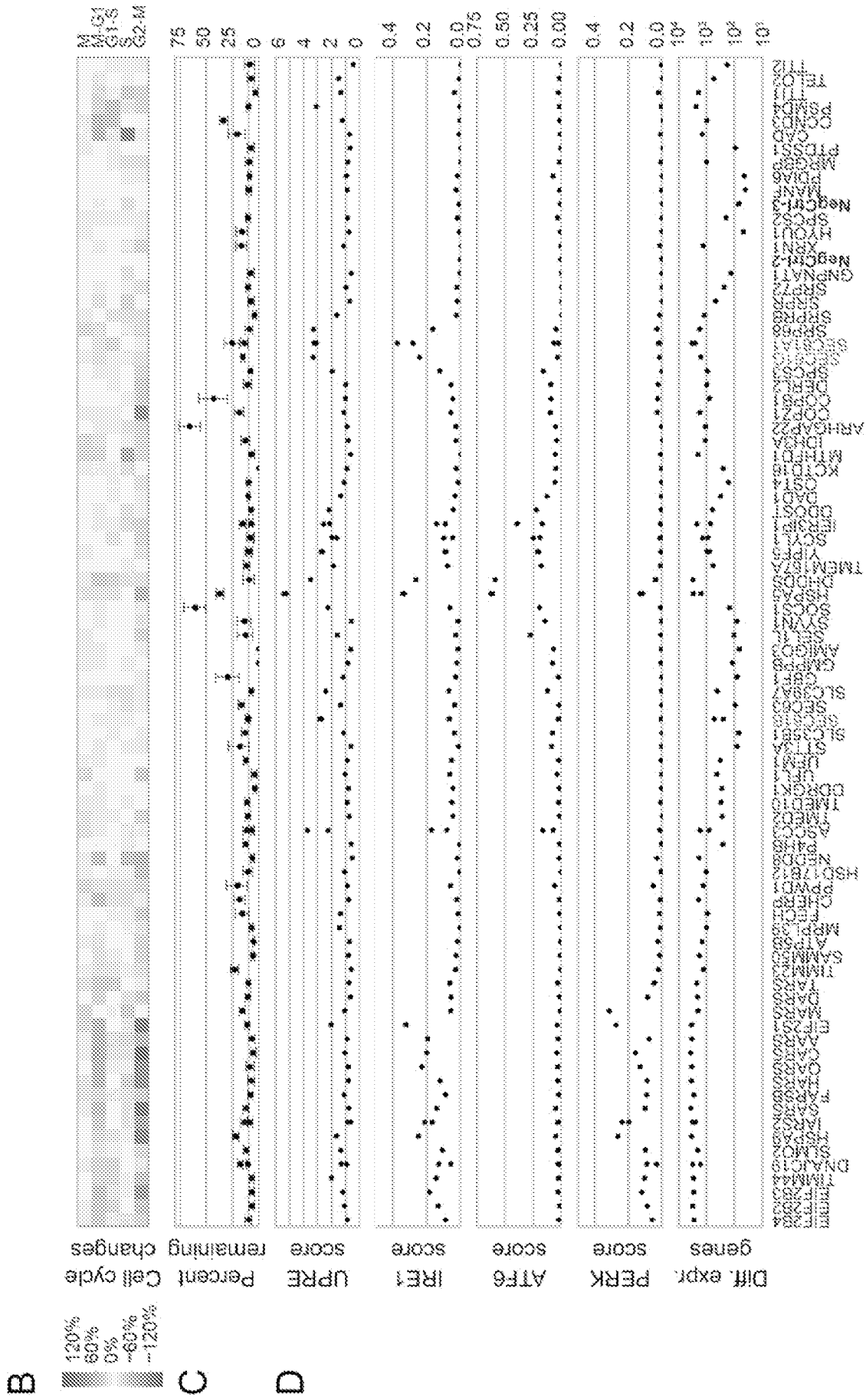


FIG. 86A-C

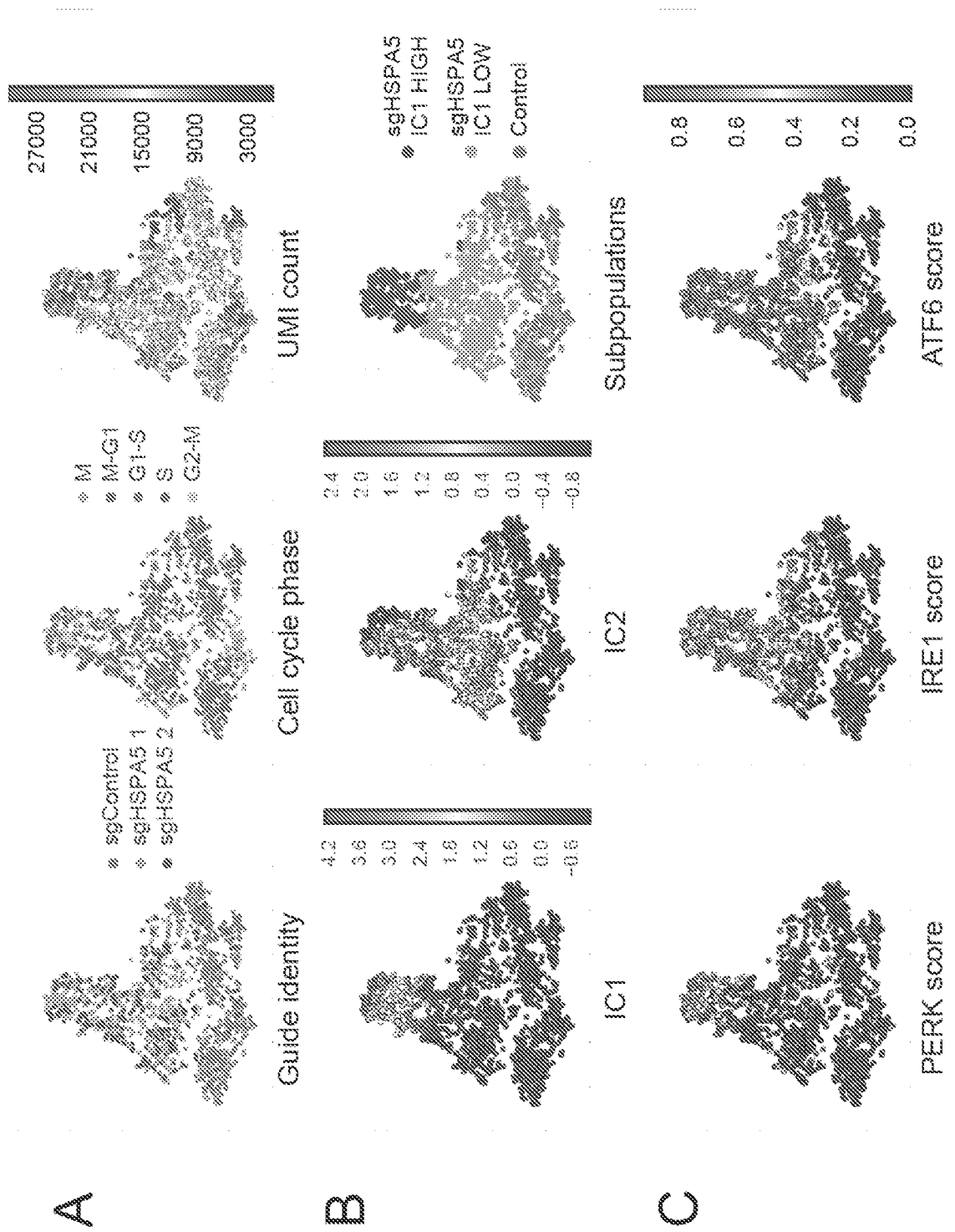


FIG. 86D-F

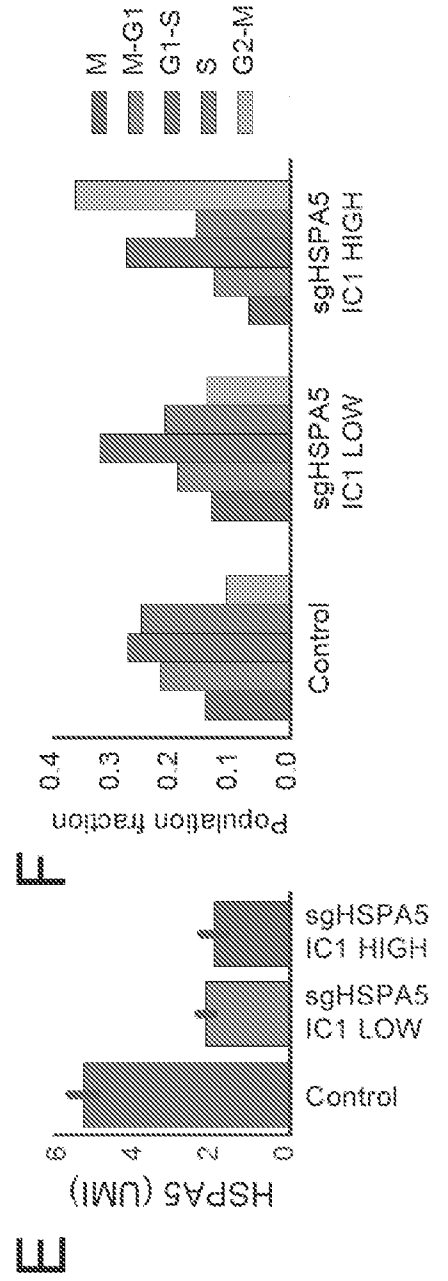
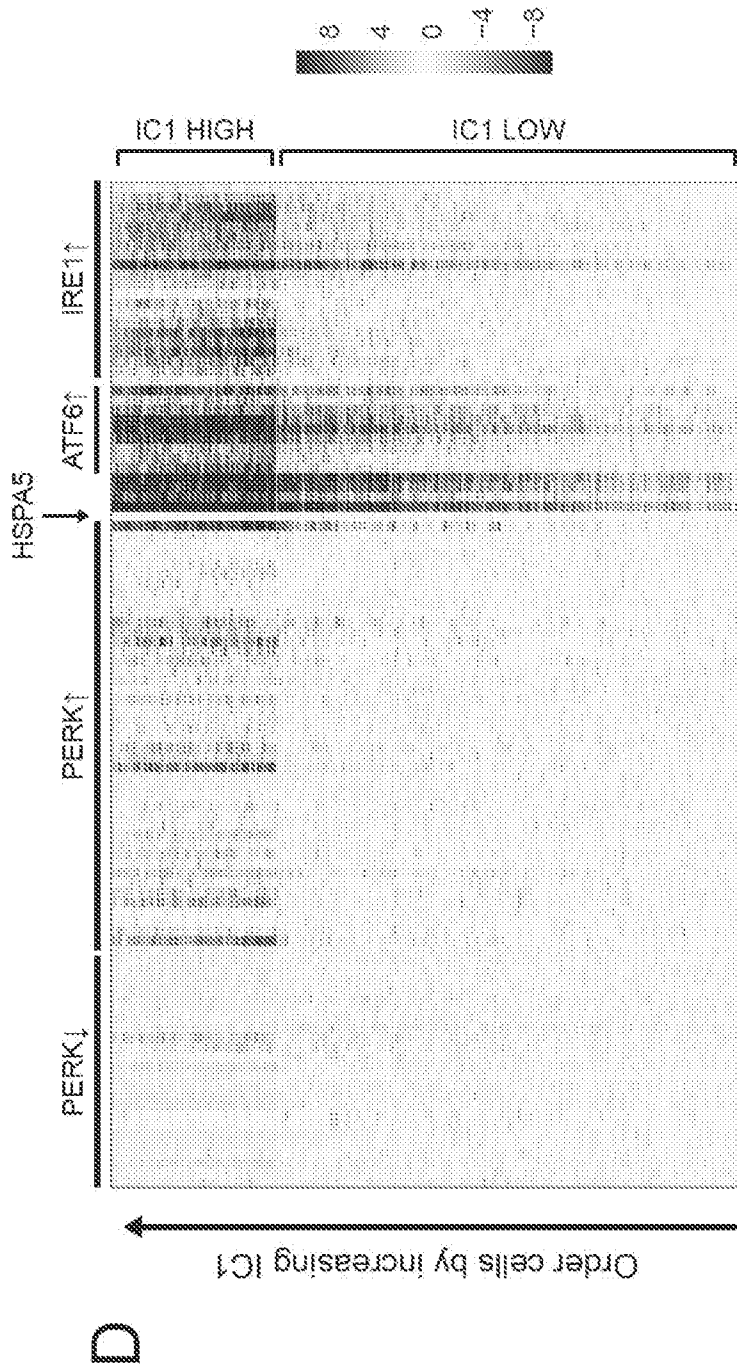


FIG. 86G-H

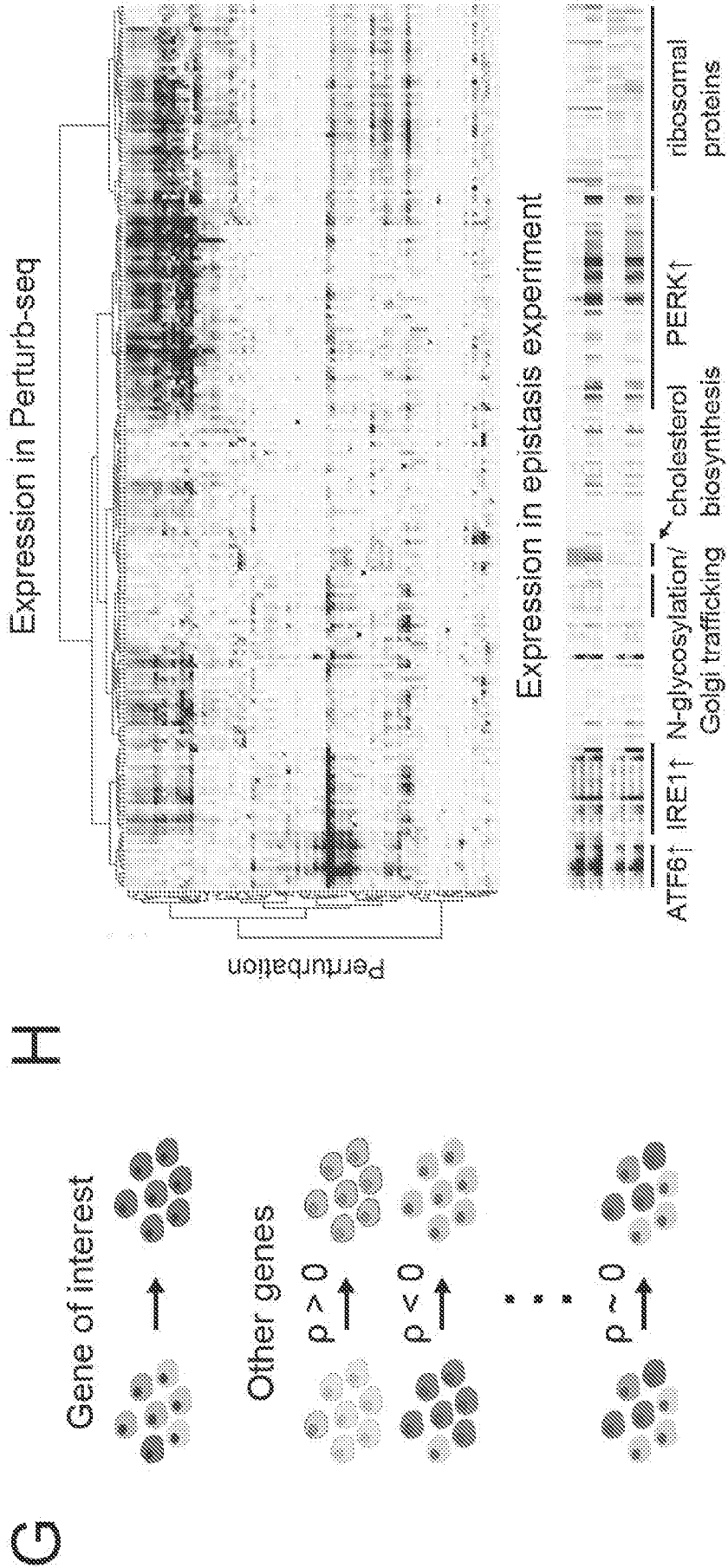


FIG. 86I-J

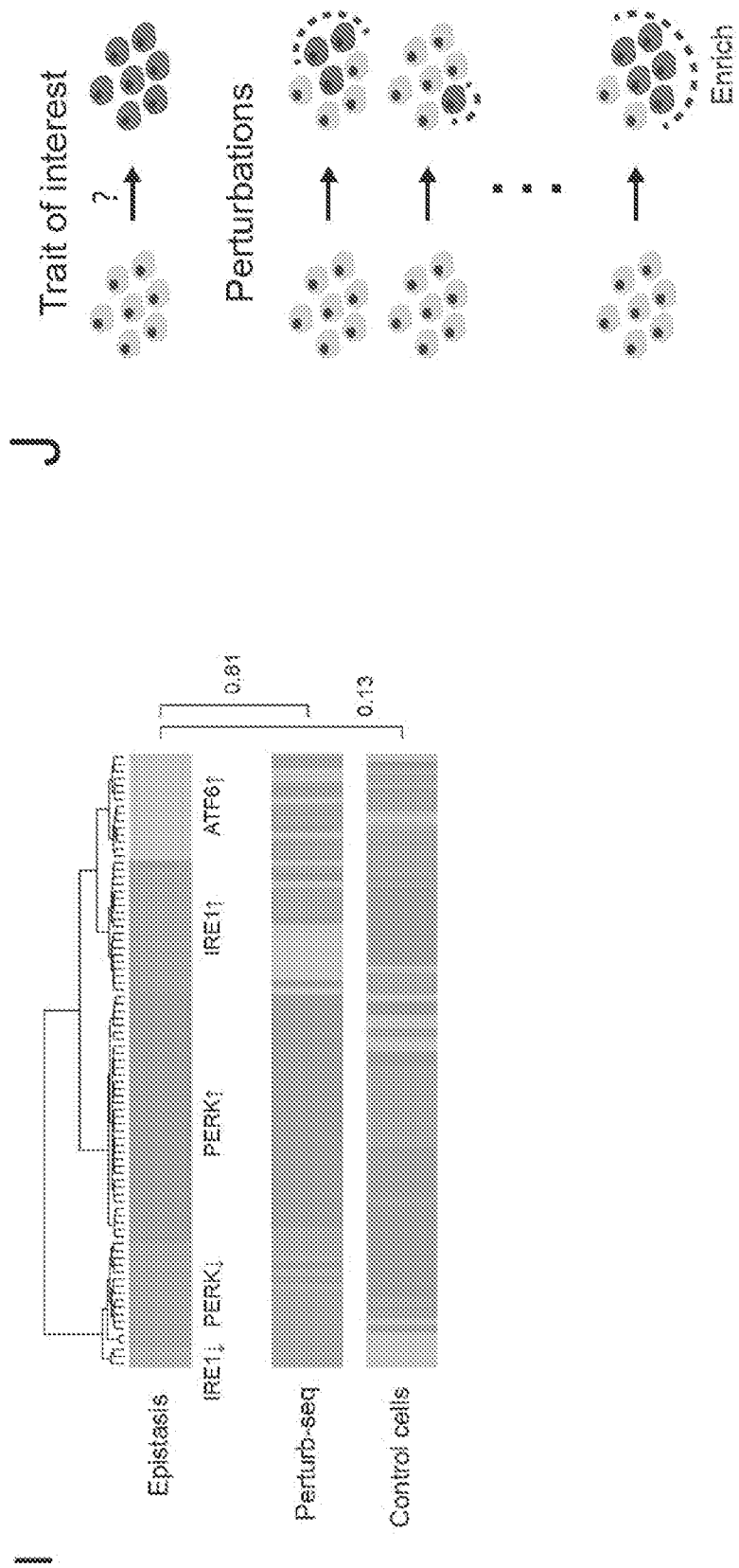


FIG. 86K-L

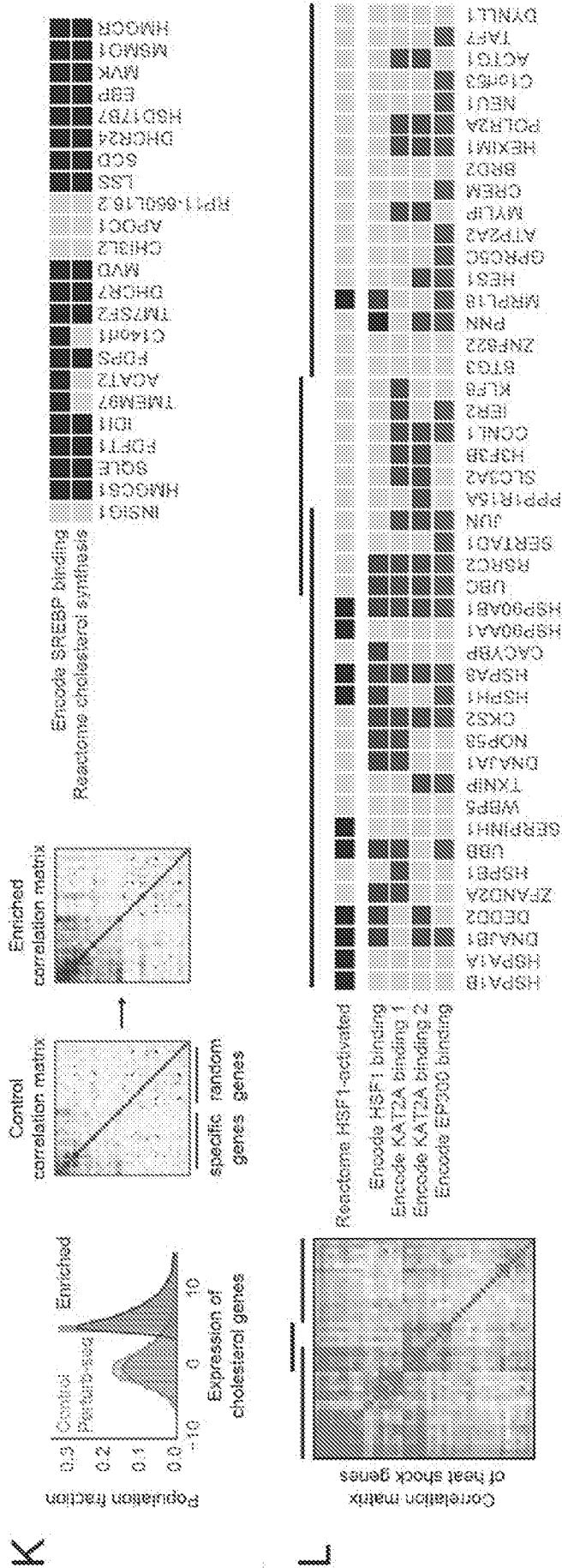


FIG. 87A-C

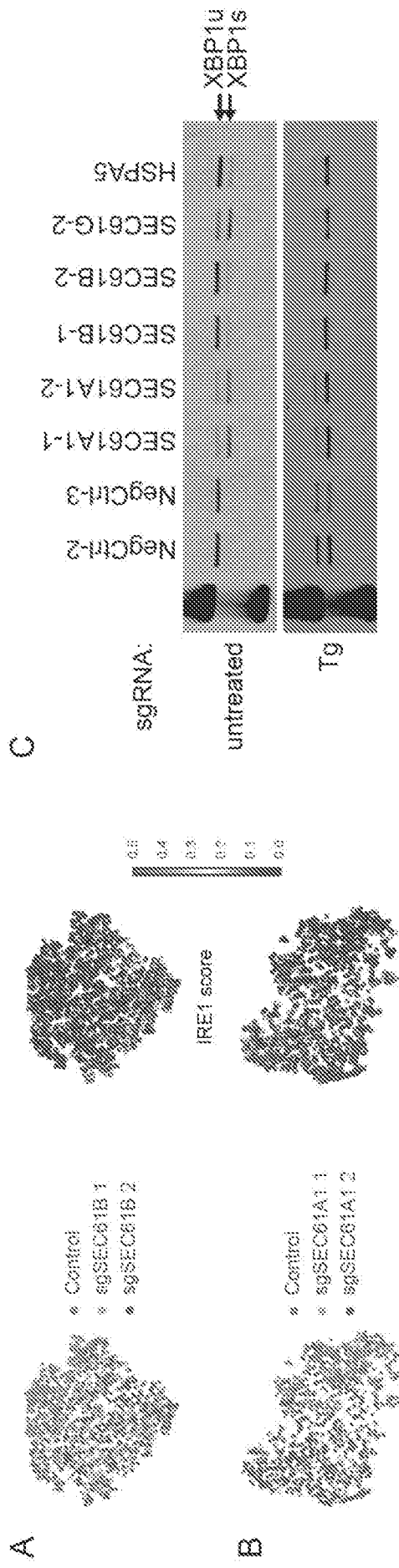
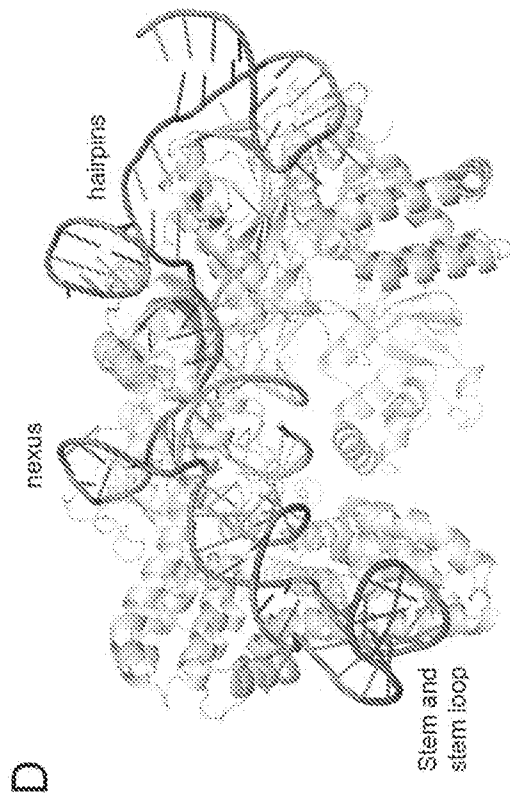
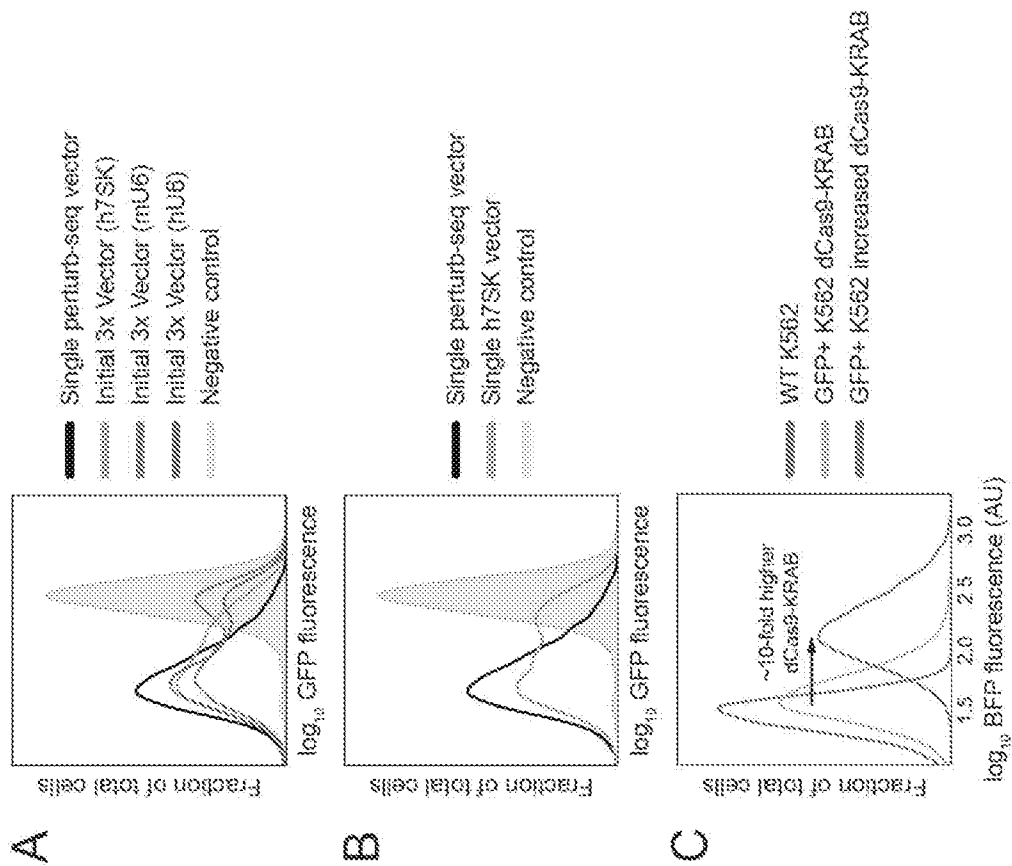


FIG. 88A-E



E

promoter	% knockdown (GFP)
mU6 (perturb-seq vector, exp 1)	95.2
hU6-2	95.4
hU6-3	93.4
sU6-1	94.9
sU6-2	95.1
mU6 (perturb-seq vector, exp 2)	94.4
bulU6	94.8
pU6	94.4

FIG. 88F

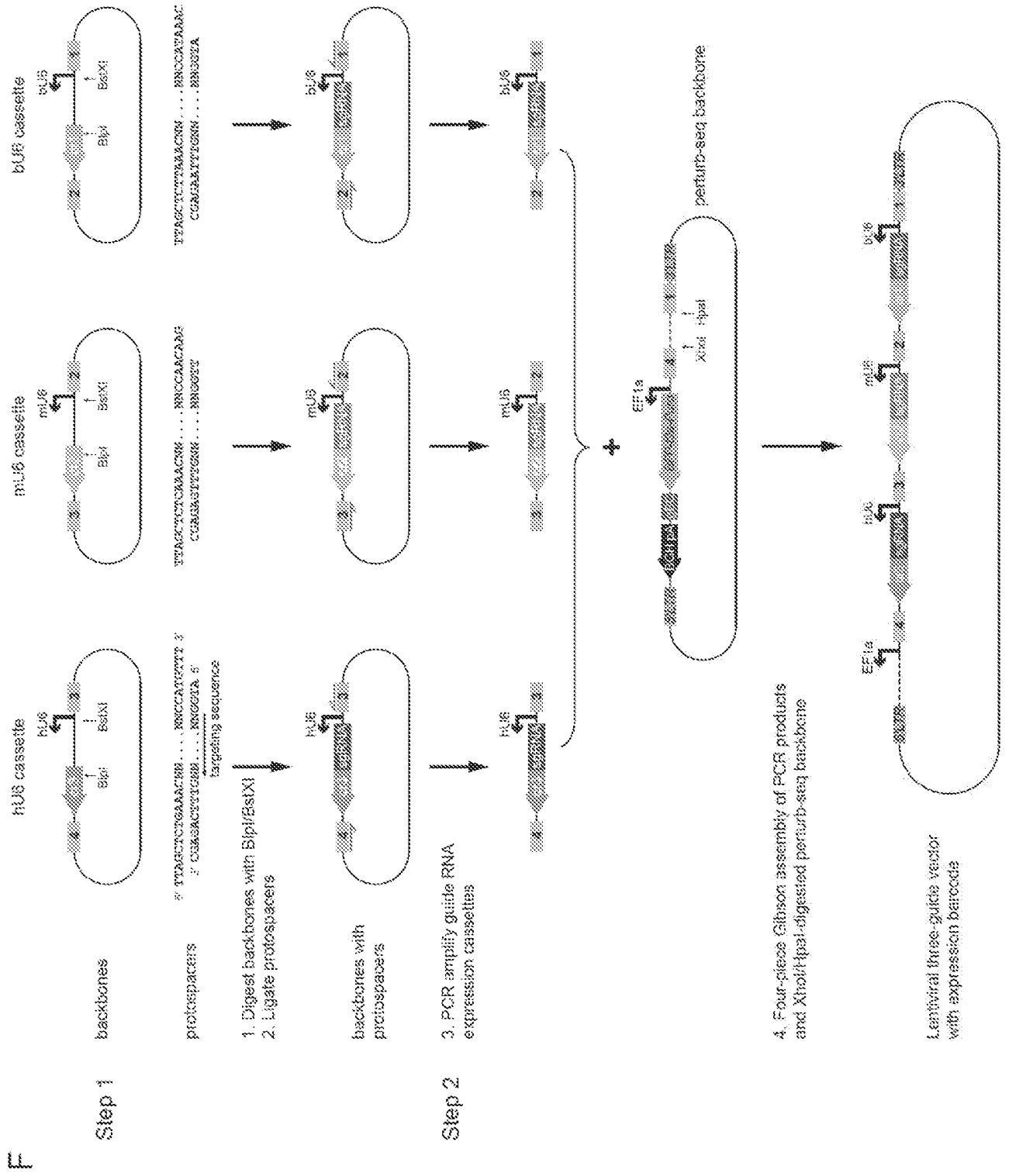


FIG. 89A

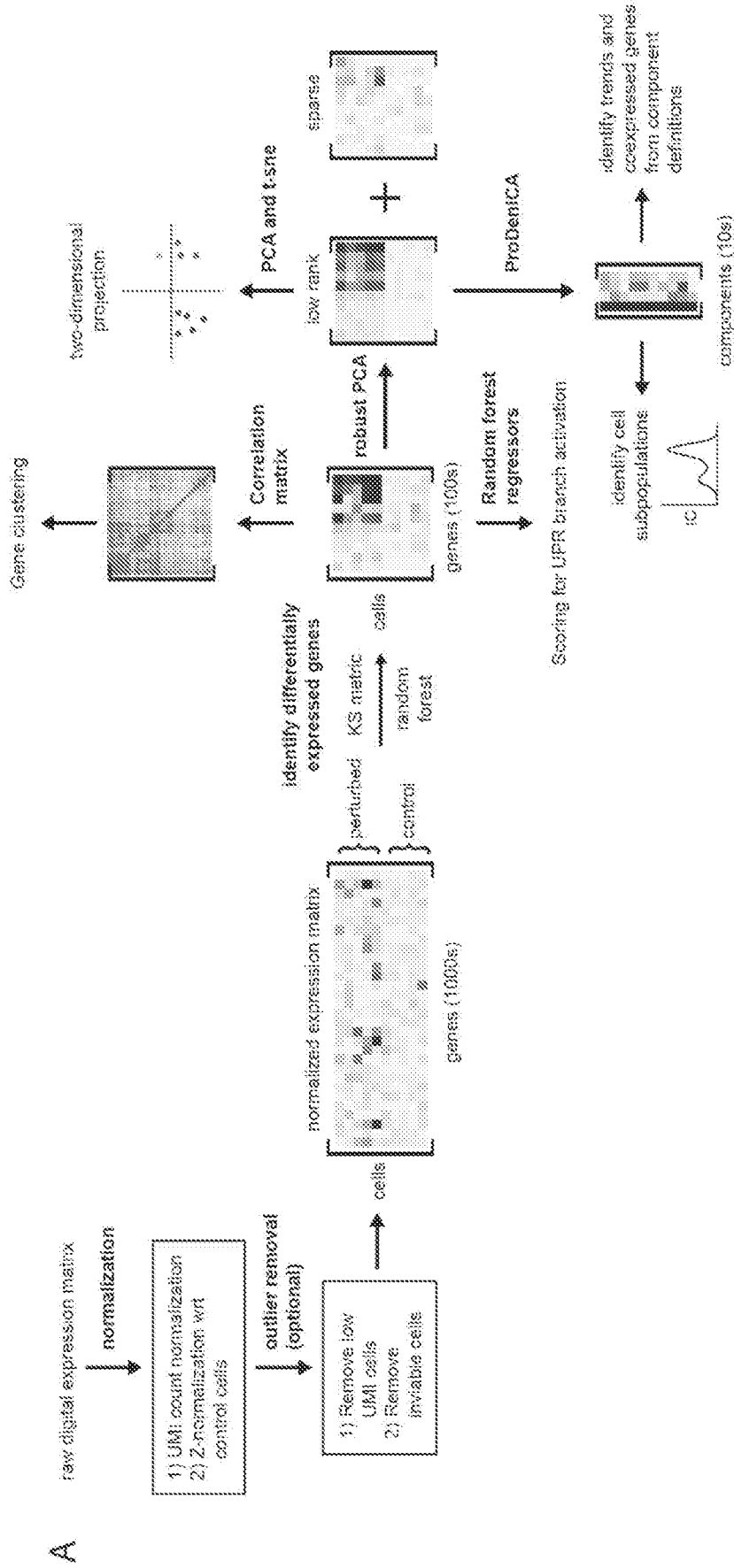


FIG. 89B

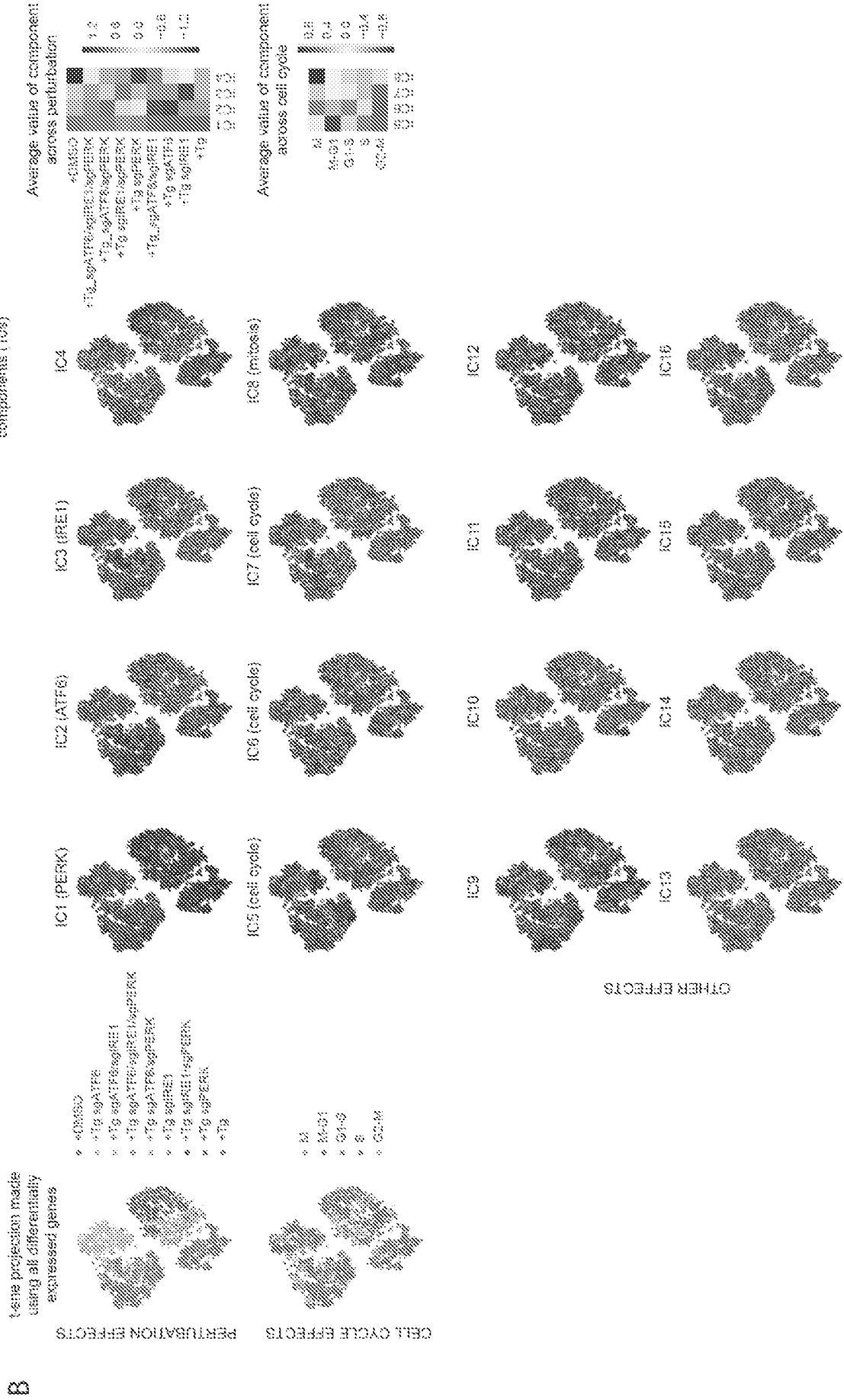


FIG. 90A-B

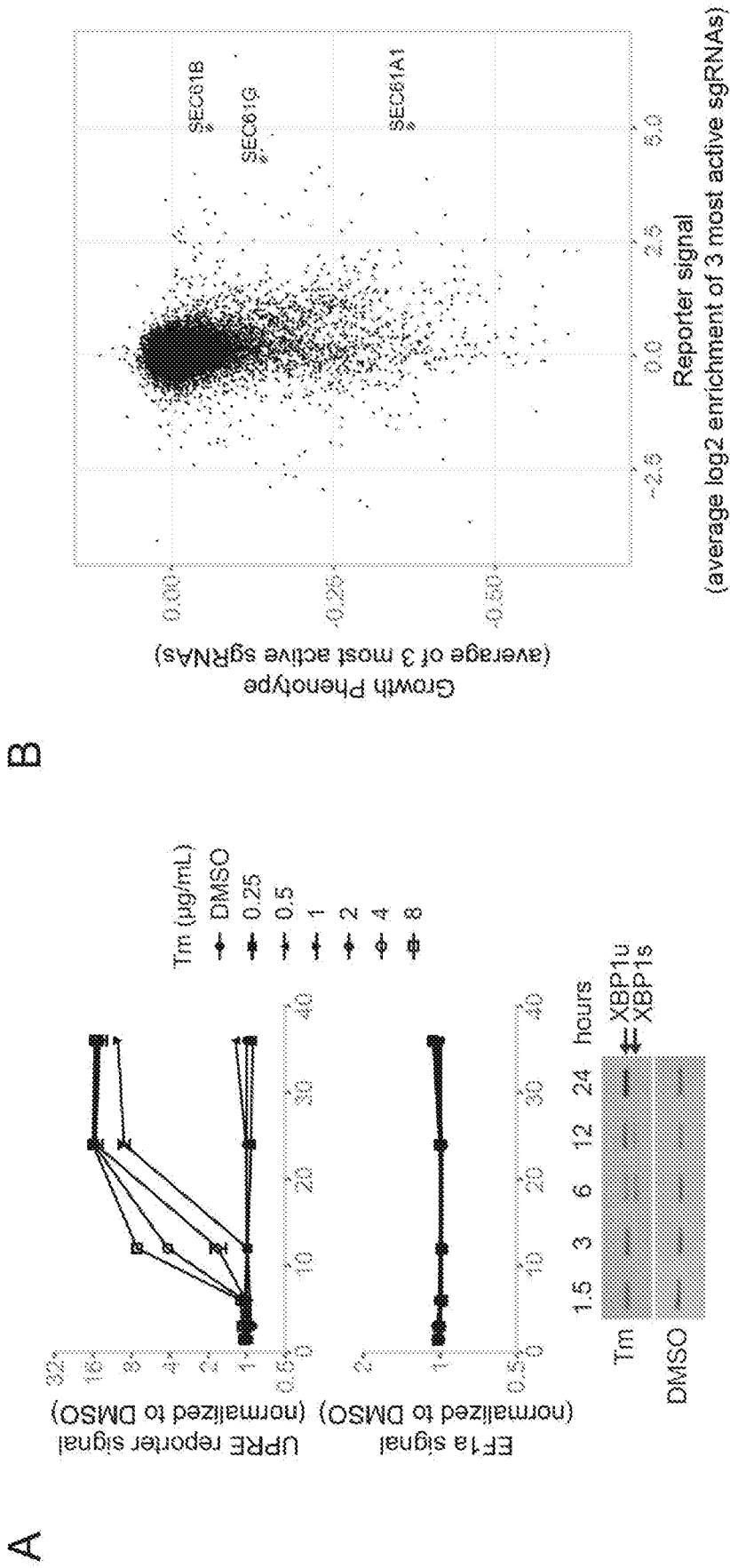
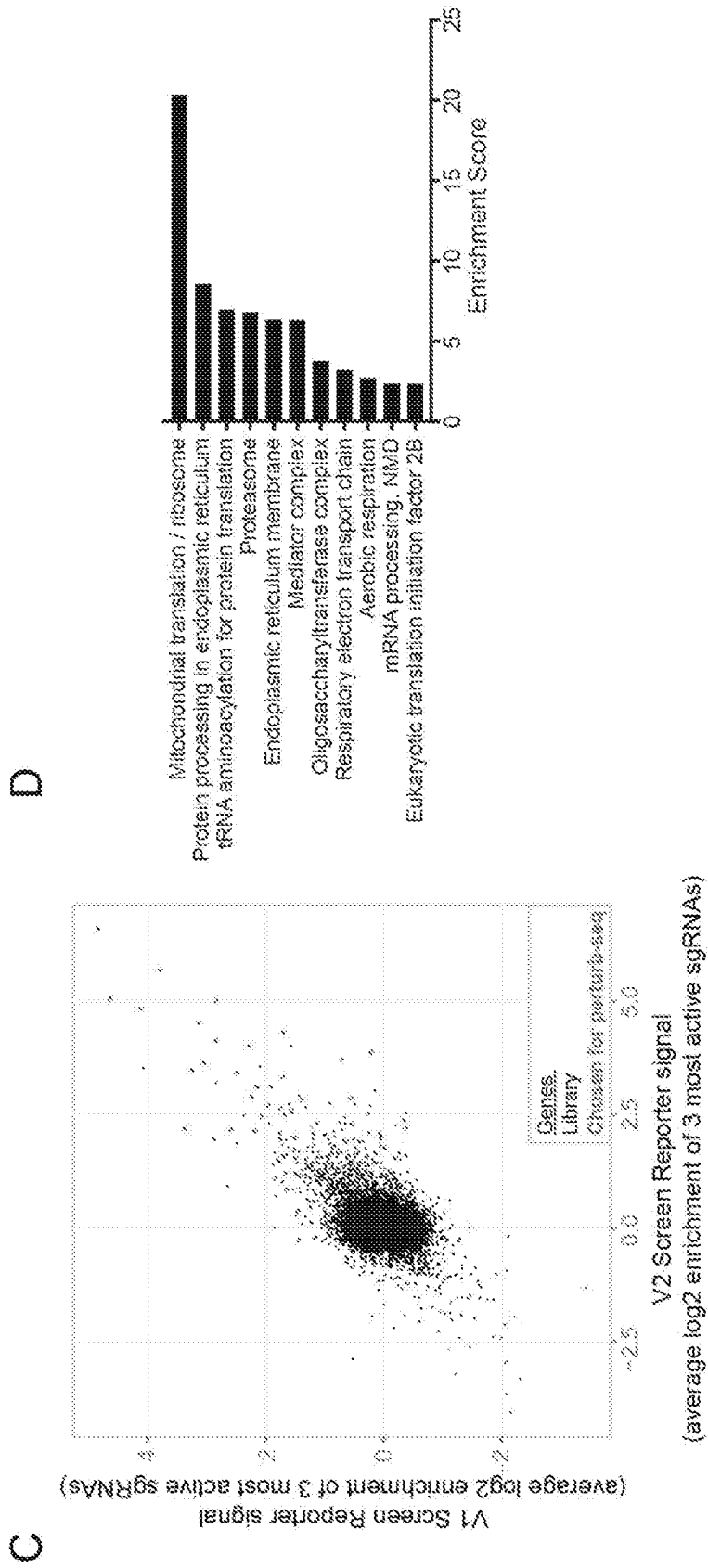


FIG. 90C-D



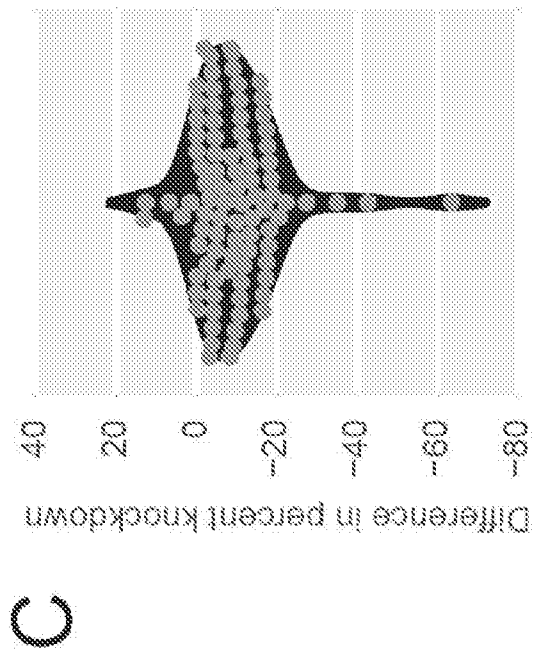
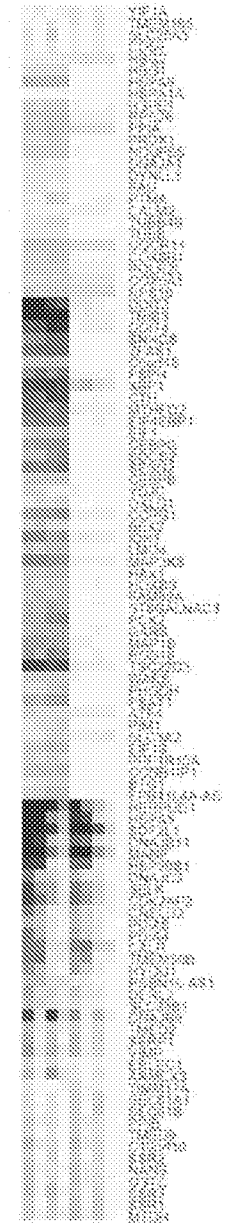
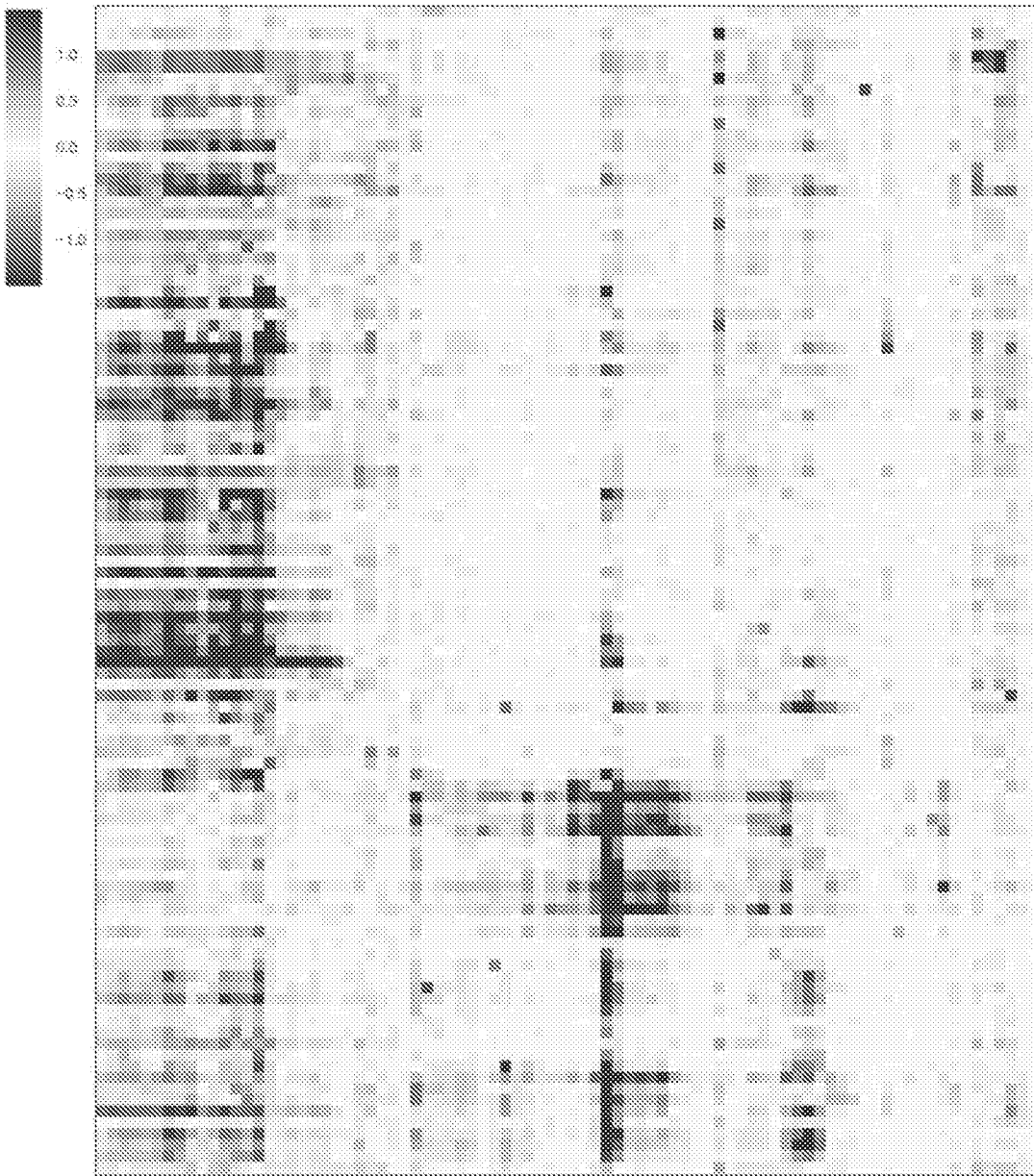


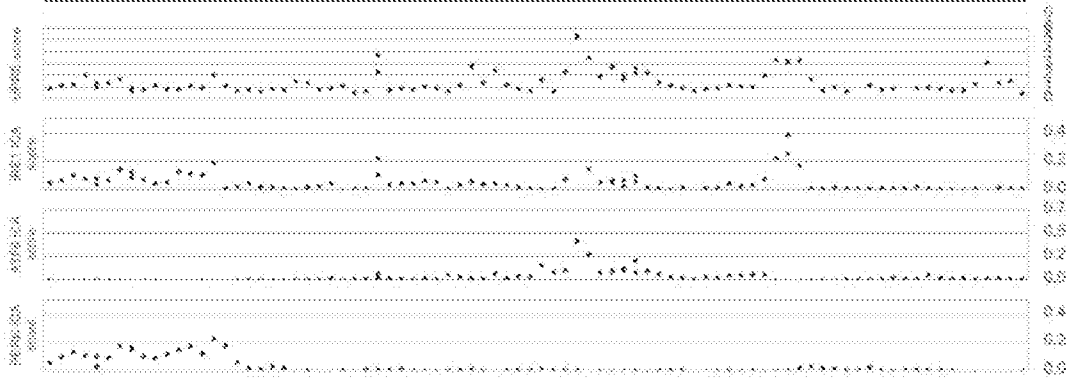
FIG. 91C

FIG. 91D-E

D



E



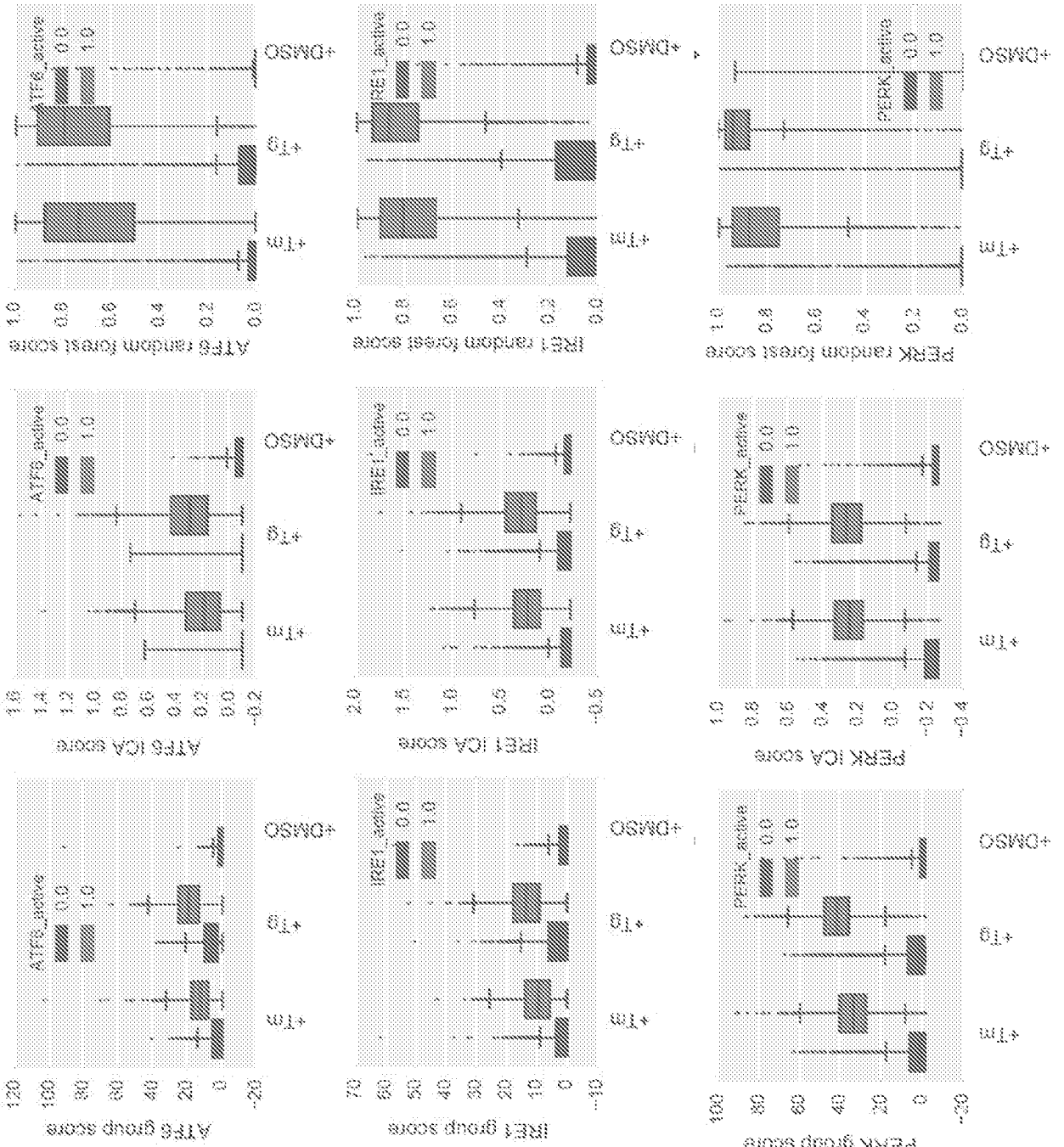


FIG. 91F
F

FIG. 92A

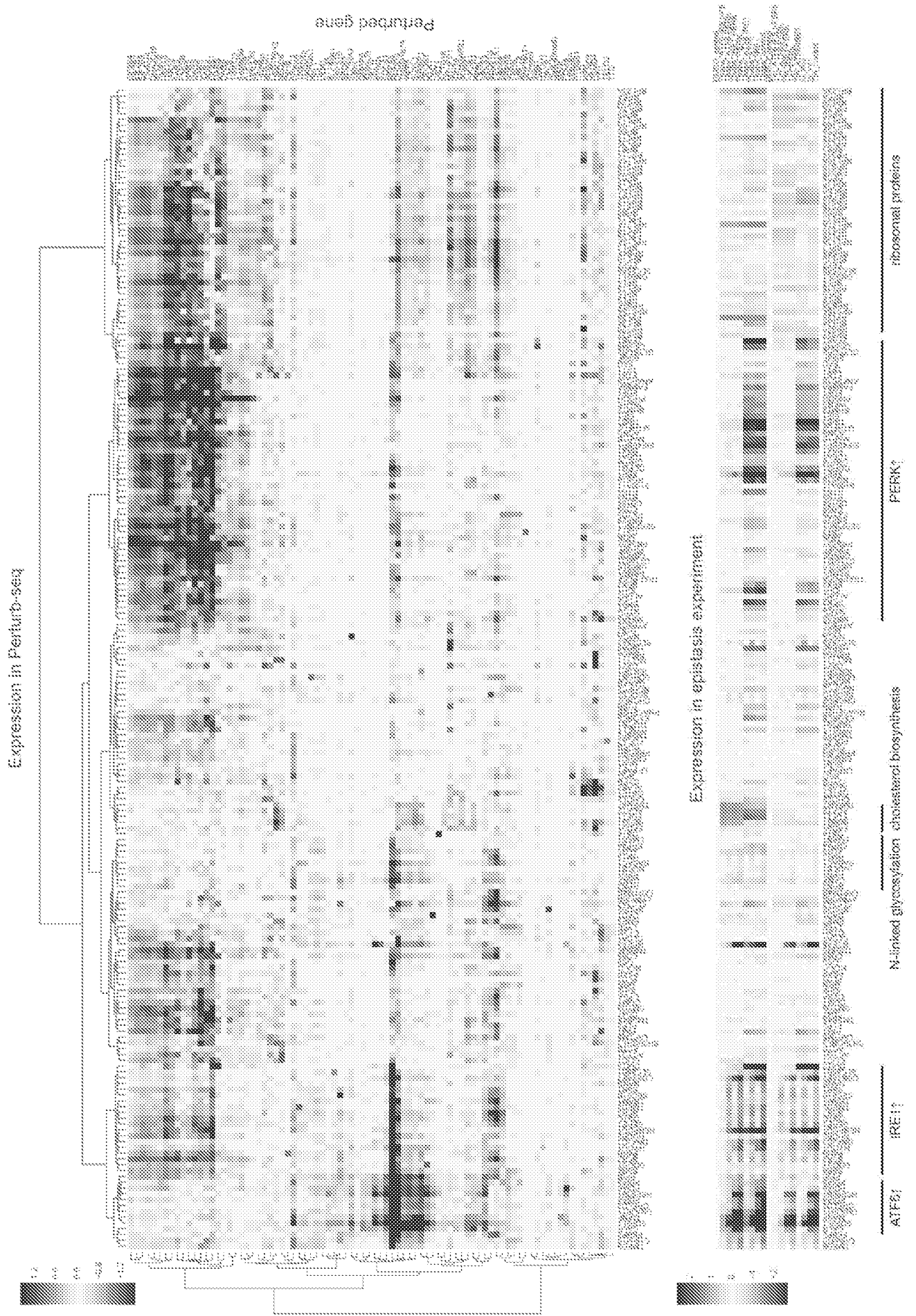
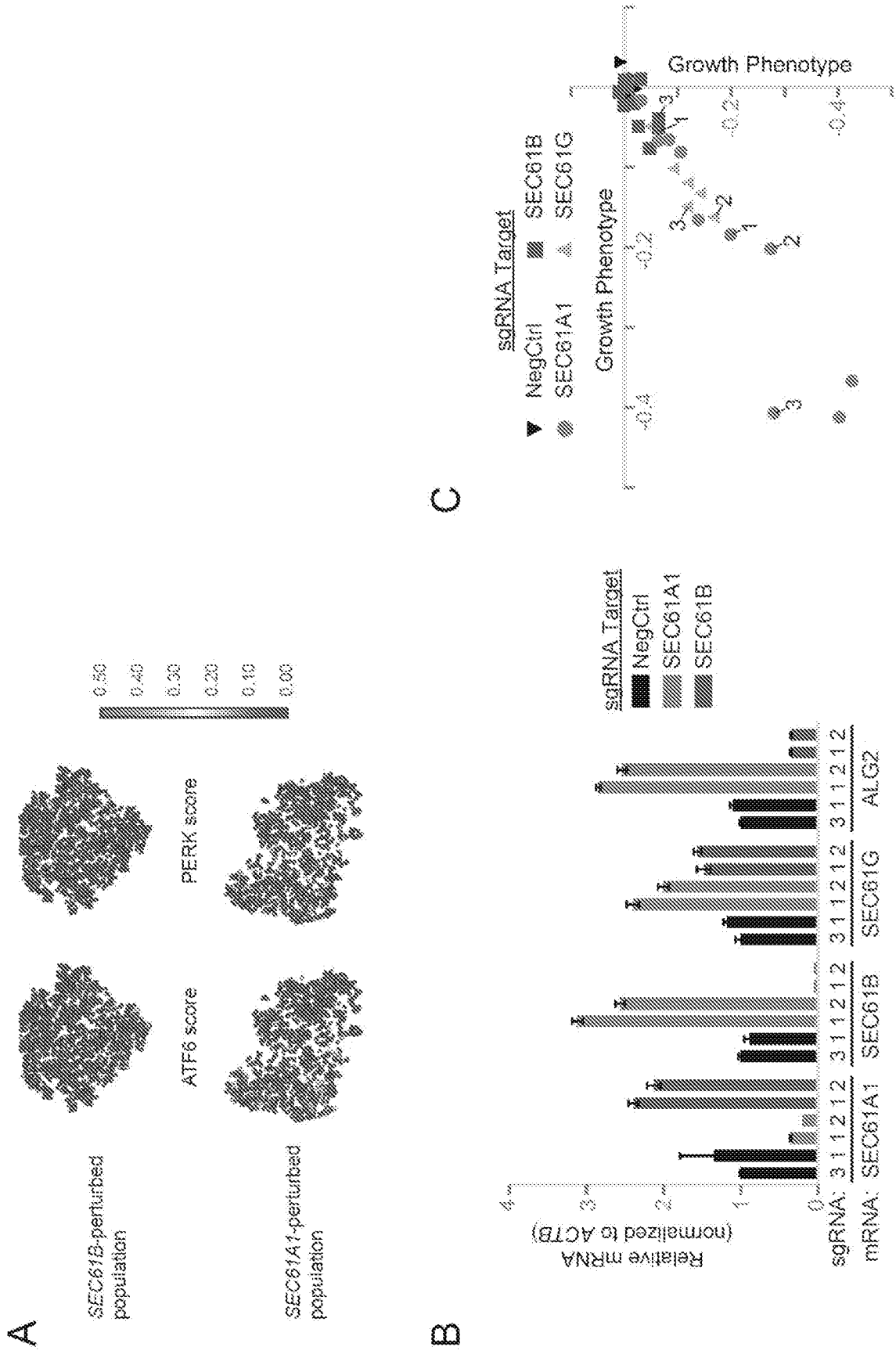


FIG. 93A-C



INTERNATIONAL SEARCH REPORT

International application No.
PCT/US2016/059233

Box No. II Observations where certain claims were found unsearchable (Continuation of item 2 of first sheet)

This international search report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. Claims Nos.: 70
because they relate to subject matter not required to be searched by this Authority, namely:
see FURTHER INFORMATION sheet PCT/ISA/210
2. Claims Nos.:
because they relate to parts of the international application that do not comply with the prescribed requirements to such an extent that no meaningful international search can be carried out, specifically:
3. Claims Nos.:
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

Box No. III Observations where unity of invention is lacking (Continuation of item 3 of first sheet)

This International Searching Authority found multiple inventions in this international application, as follows:

1. As all required additional search fees were timely paid by the applicant, this international search report covers all searchable claims.
2. As all searchable claims could be searched without effort justifying an additional fees, this Authority did not invite payment of additional fees.
3. As only some of the required additional search fees were timely paid by the applicant, this international search report covers only those claims for which fees were paid, specifically claims Nos.:
4. No required additional search fees were timely paid by the applicant. Consequently, this international search report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

Remark on Protest

- The additional search fees were accompanied by the applicant's protest and, where applicable, the payment of a protest fee.
- The additional search fees were accompanied by the applicant's protest but the applicable protest fee was not paid within the time limit specified in the invitation.
- No protest accompanied the payment of additional search fees.

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2016/059233

A. CLASSIFICATION OF SUBJECT MATTER
INV. G06F19/12
ADD.
According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED
Minimum documentation searched (classification system followed by classification symbols)
G06F
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)
EPO-Internal, WPI Data

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	<p>PARNAS OREN ET AL: "A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks", CELL, vol. 162, no. 3, 30 July 2015 (2015-07-30) , pages 675-686, XP029248090, ISSN: 0092-8674, DOI: 10.1016/J.CELL.2015.06.059 abstract page 675 - page 677, left-hand column, paragraph 2 page 679, right-hand column, paragraph 2 - page 681, left-hand column, paragraph 2 page 684, left-hand column, paragraph 4 - page 685, left-hand column, paragraph 2 ----- -/--</p>	<p>1-69, 71-83</p>

Further documents are listed in the continuation of Box C.

See patent family annex.

* Special categories of cited documents :

- "A" document defining the general state of the art which is not considered to be of particular relevance
- "E" earlier application or patent but published on or after the international filing date
- "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- "O" document referring to an oral disclosure, use, exhibition or other means
- "P" document published prior to the international filing date but later than the priority date claimed

- "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
- "&" document member of the same patent family

Date of the actual completion of the international search 20 January 2017	Date of mailing of the international search report 31/01/2017
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016	Authorized officer Kürten, Ivayla

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2016/059233

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	<p>WO 2015/130968 A2 (BROAD INST INC [US]; BRIGHAM & WOMENS HOSPITAL [US]; HARVARD COLLEGE []) 3 September 2015 (2015-09-03) abstract paragraph [0226] - paragraph [0230] paragraph [0234] - paragraph [0269] paragraph [0277] - paragraph [0302]</p> <p>-----</p>	1-69, 71-83
X	<p>KAREN SACHS ET AL: "Causal Protein-Signaling Networks Derived from Multiparameter Single-Cell Data", SCIENCE, AMERICAN ASSOCIATION FOR THE ADVANCEMENT OF SCIENCE, vol. 308, no. 5721, 22 April 2005 (2005-04-22), pages 523-529, XP007901146, ISSN: 0036-8075, DOI: 10.1126/SCIENCE.1105809 abstract page 523, right-hand column, paragraph 2 - page 528, left-hand column</p> <p>-----</p>	1-69, 71-83
X	<p>WO 2014/145631 A1 (BROAD INST INC [US]; HARVARD COLLEGE [US]) 18 September 2014 (2014-09-18) abstract paragraph [0005] - paragraph [0009] paragraph [0067] - paragraph [0069] paragraph [0079] - paragraph [0090] paragraph [0161] - paragraph [0181] paragraph [0191] - paragraph [0216]</p> <p>-----</p>	1-69, 71-83
A	<p>OLIVER STEGLE ET AL: "Computational and analytical challenges in single-cell transcriptomics", NATURE REVIEWS GENETICS, vol. 16, no. 3, 28 January 2015 (2015-01-28), pages 133-145, XP055336846, GB ISSN: 1471-0056, DOI: 10.1038/nrg3833 abstract page 134, left-hand column, paragraph 2 - right-hand column, paragraph 2; figure 1 page 138, right-hand column, paragraph 4 - page 142, left-hand column, paragraph 1</p> <p>-----</p> <p style="text-align: center;">-/--</p>	1-69, 71-83

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2016/059233

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	PE'ER ET AL: "Inferring subnetworks from perturbed expression profiles", BIOINFORMATICS, OXFORD UNIVERSITY PRESS, SURREY, GB, vol. 17, no. SUPPL. 1, 1 June 2001 (2001-06-01), pages S215-S224, XP002963503, ISSN: 1367-4803 abstract page S215 - page S217 -----	1-69, 71-83

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2016/059233

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
WO 2015130968 A2	03-09-2015	AU 2015222944 A1	08-09-2016
		CA 2940653 A1	03-09-2015
		EP 3110971 A2	04-01-2017
		US 2016377631 A1	29-12-2016
		WO 2015130968 A2	03-09-2015

WO 2014145631 A1	18-09-2014	AU 2014232794 A1	08-10-2015
		CA 2904099 A1	18-09-2014
		CN 105683392 A	15-06-2016
		EP 2971116 A1	20-01-2016
		JP 2016515383 A	30-05-2016
		KR 20150139537 A	11-12-2015
		US 2015368719 A1	24-12-2015
		WO 2014145631 A1	18-09-2014

FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210

Continuation of Box II.1

Claims Nos.: 70

Claim 70 discloses a method for treatment of the human or animal body ("the population of cells is adoptively transferred to a subject"). Hence, under Article 17(2)(a)(i) and Rule 39.1(iv) PCT, this International Searching Authority is not required to search the subject-matter of said claims. Consequently, no opinion with regard to novelty, inventive step or industrial applicability can be given for the subject-matter of said claim.