

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2023-534882
(P2023-534882A)

(43)公表日 令和5年8月14日(2023.8.14)

(51)国際特許分類		F I		テーマコード(参考)	
C 1 2 Q	1/6837(2018.01)	C 1 2 Q	1/6837	Z Z N A	4 B 0 2 9
C 1 2 Q	1/6869(2018.01)	C 1 2 Q	1/6869	Z	4 B 0 6 3
C 1 2 Q	1/34 (2006.01)	C 1 2 Q	1/34		
C 1 2 Q	1/25 (2006.01)	C 1 2 Q	1/25		
C 1 2 M	1/00 (2006.01)	C 1 2 M	1/00	A	

審査請求 未請求 予備審査請求 未請求 (全69頁) 最終頁に続く

(21)出願番号	特願2023-521274(P2023-521274)	(71)出願人	515236259 ザ・ブロード・インスティテュート・インコーポレイテッド アメリカ合衆国・マサチューセッツ・02142・ケンブリッジ・メイン・ストリート・415
(86)(22)出願日	令和3年6月14日(2021.6.14)	(71)出願人	596060697 マサチューセッツ インスティテュート オブ テクノロジー アメリカ合衆国マサチューセッツ州02139ケンブリッジ, マサチューセッツ・アヴェニュー・77
(85)翻訳文提出日	令和5年2月14日(2023.2.14)	(71)出願人	592017633 ザ ジェネラル ホスピタル コーポレーション
(86)国際出願番号	PCT/US2021/037226		
(87)国際公開番号	WO2021/257453		
(87)国際公開日	令和3年12月23日(2021.12.23)		
(31)優先権主張番号	63/039,004		
(32)優先日	令和2年6月15日(2020.6.15)		
(33)優先権主張国・地域又は機関	米国(US)		
(81)指定国・地域	AP(BW,GH,GM,KE,LR,LS,MW,MZ,NA,RW,SD,SL,ST,SZ,TZ,UG,ZM,ZW),EA(AM,AZ,BY,KG,KZ,RU,TJ,TM),EP(AL,AT,BE,BG,CH,CY,CZ,DE,DK,EE,ES,FI,FR,GB,GR,HR,HU,IE,IS,IT,LT,LU,LV,MC,		

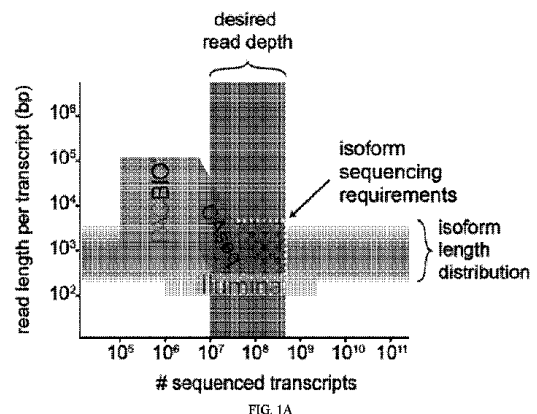
最終頁に続く

(54)【発明の名称】 キメラアンプリコンアレイ配列決定

(57)【要約】

本開示は、核酸配列決定のための組成物及び方法に関し、具体的には、少なくともある態様では、入力配列のキメラアレイを提供することによって、既知のロングレンジ配列決定プラットフォームの有効性、スループット及び/又は収率を増強するための方法及び組成物を提供する。そのようなコンポーネント核酸配列要素のアレイは、バイアスの導入を最小限に抑える方法によって調製することができる。本キメラアンプリコン配列決定プロセスを用いるミトコンドリア系統追跡のための方法と同様に、例えば患者試料からアイソフォーム配列決定情報を得るための現在の方法の適用も具体的に提供される。アレイ核酸配列の処理及び解釈のための方法及びシステムも提供される。

【選択図】 図 1 A



【特許請求の範囲】

【請求項 1】

アレイ核酸配列を調製するための方法であって、前記方法が、

- i) 各入力核酸配列が約 30 キロベース長以下である複数の入力核酸配列を得ること、
 - ii) 1 つ又は複数のアダプタ配列を前記複数の入力核酸配列に付加し、それによって適合核酸配列の集団を生成すること、
 - iii) 前記適合核酸配列の集団を、前記適合核酸配列の集団内の各適合核酸配列の少なくとも 1 つの末端に一本鎖末端を生成することができる酵素と接触させ、それにより、一本鎖末端を有する核酸配列の集団を形成すること、及び
 - iv) 前記一本鎖末端を有する核酸配列の集団をリガーゼと接触させること、
- を含み、
それにより、アレイ核酸配列を形成する、方法。

10

【請求項 2】

前記 1 つ又は複数のアダプタ配列のうち少なくとも 1 つが、1 つの鎖上に内部 d U を含む、請求項 1 に記載の方法。

【請求項 3】

前記アレイ核酸配列が、少なくとも 20 キロベース、任意選択的に少なくとも 50 キロベース、任意選択的に約 100 kb 以上の長さを有する、請求項 1 に記載の方法。

【請求項 4】

前記複数の入力核酸配列が、約 0.5 kb ~ 20 kb の長さである、請求項 1 に記載の方法。

20

【請求項 5】

前記複数の入力核酸配列が、1 つ又は複数の cDNA ライブラリ、任意選択で 1 つ又は複数の単一細胞若しくは空間 cDNA ライブラリから得られる、請求項 1 に記載の方法。

【請求項 6】

工程 (ii) が、前記複数の核酸配列を前記対になった増幅プライマーと接触させること、ここで、前記対になった増幅プライマー内の少なくとも 1 つのプライマーが 1 つの鎖上に内部 d U を含むアダプタ配列を含み、及び、少なくとも 1 ラウンドの増幅を実施し、それにより、適合核酸配列の集団を生成することを含む、請求項 1 に記載の方法。

【請求項 7】

前記対になった増幅プライマー内の少なくとも 1 つのプライマーがビオチン化されており、任意選択で、アダプタ配列テールアンプリコンのためのビオチン媒介選択が行われる、請求項 6 に記載の方法。

30

【請求項 8】

工程 (iii) が、前記適合核酸配列の集団をウラシル DNA グリコシラーゼ及びエンドヌクレアーゼ V I I I と接触させ、それにより一本鎖末端を有する核酸配列の集団を形成することを更に含む、請求項 2、6 又は 7 のいずれか一項に記載の方法。

【請求項 9】

前記アダプタ配列が 5 ~ 30 塩基対の長さを含み (標的核酸配列を除く)、任意選択で、前記アダプタ配列が 6 ~ 25 塩基対の長さであり、任意選択で、前記アダプタ配列が構造 5' - N 6 - 1 6 __ d U __ t a r g e t - D N A - 3' を有する、請求項 1 に記載の方法。

40

【請求項 10】

1 つの鎖上に内部 d U を含む前記アダプタ配列が、配列番号 1 ~ 18 からなる群より選択される配列を含む、請求項 1 に記載の方法。

【請求項 11】

アダプタ配列を有する複数の核酸配列について、各アダプタ配列が、アダプタ配列を有する前記複数の核酸配列のうち少なくとも 1 つの他のものと相補的な 1 つ又は 2 つの指定された配列を有し、それにより、前記複数のアダプタ配列が相補的なアダプタ配列の集団を形成し、任意選択で、前記相補的なアダプタ配列の集団の各相補的なアダプタ配列が

50

、前記相補的なアダプタ配列の集団の互いの相補的なアダプタ配列に対して最小の類似性を有し、任意選択で、前記相補的なアダプタ配列の集団の各相補的なアダプタ配列が、前記相補的なアダプタ配列の集団の他の全ての相補的なアダプタ配列から少なくとも11ハミング距離単位離れている、請求項1に記載の方法。

【請求項12】

以下の：前記複数の入力核酸配列；前記適合核酸配列の集団；及び/又は前記一本鎖末端を有する核酸配列の集団のうちの1つ又は複数がサイズ選択され、任意選択で前記サイズ選択が電気泳動を介して、任意選択でアガロースゲル上で行われる、請求項1に記載の方法。

【請求項13】

前記アレイ核酸配列の配列情報が、任意選択でロングリード配列決定プラットフォームを使用して得られる、請求項1に記載の方法。

【請求項14】

ハプロタイプフェージングの配列情報が前記アレイ核酸配列にわたって得られる、請求項13に記載の方法。

【請求項15】

形成される前記アレイ核酸配列が、5つ以上の入力核酸配列、任意選択的に6つ以上、任意選択的に7つ以上、任意選択的に8つ以上、任意選択的に9つ以上、任意選択的に10以上、任意選択的に11以上、任意選択的に12以上、任意選択的に13以上、任意選択的に14以上、任意選択的に15以上、任意選択的に16以上、任意選択的に17以上、任意選択的に18以上、任意選択的に19以上、任意選択的に20以上を含む、請求項1に記載の方法。

【請求項16】

標的化アイソフォーム配列決定情報が、前記複数の入力核酸配列を得る工程(i)中に遺伝子パネルの標的化を介して得られる、請求項13に記載の方法。

【請求項17】

前記複数の入力核酸配列が、免疫応答経路のためのcDNAを含む、請求項1に記載の方法。

【請求項18】

前記複数の入力核酸配列がミトコンドリアDNAから得られ、任意選択で、前記アレイ核酸配列の配列決定がミトコンドリアDNA系統追跡に使用される、請求項1に記載の方法。

【請求項19】

前記適合核酸配列の集団が、ギブソンアセンブリを介して連結される、請求項1に記載の方法。

【請求項20】

前記アレイ核酸配列が線状アレイである、請求項1に記載の方法。

【請求項21】

前記アレイ核酸配列が環状アレイである、請求項1に記載の方法。

【請求項22】

核酸配列の線状アレイのアレイを調製するための方法であって、前記方法が、

i) 請求項20に記載の方法によって入力核酸配列の第1の集団から第1の線状アレイを調製すること、

ii) 請求項20に記載の方法によって、入力核酸配列の第2の集団から第2の線状アレイを調製すること、ここで、前記第1の線状アレイ及び前記第2の線状アレイは各々、適合する相補的フランキング配列を有し、

iii) 前記第1の線状アレイ及び前記第2の線状アレイを溶液中で組み合わせること、及び

iv) 溶液中の前記第1の線状アレイ及び前記第2の線状アレイをリガーゼと接触させること、

10

20

30

40

50

を含み、

それにより、核酸配列の線状アレイのアレイを形成する、方法。

【請求項 23】

前記第 1 の線状アレイ若しくは前記第 2 の線状アレイ、又はその両方が線状アレイのアレイを含む、請求項 22 に記載の方法。

【請求項 24】

v) 請求項 20 に記載の方法によって入力核酸配列の第 3 の集団から第 3 の線状アレイを調製すること、ここで、前記線状アレイのアレイ及び第 3 の線状アレイは各々、適合する相補的フランキング配列を有し、

v i) 前記線状アレイのアレイ及び前記第 3 の線状アレイを溶液中で組み合わせること 10

v i i) 溶液中の前記線状アレイのアレイ及び前記第 3 の線状アレイをリガーゼと接触させ、それにより、核酸配列の線状アレイのより大きなアレイを形成すること、

を更に含み、任意選択で、工程 (v) ~ (v i i) が繰り返されて、第 4 の線状アレイ、第 5 の線状アレイ、及び / 又はより多くの線状アレイが、線状アレイのより大きなアレイに組み込まれる、請求項 22 又は請求項 23 に記載の方法。

【請求項 25】

入力 c D N A 配列の集団からアイソフォーム配列決定情報を得るための方法であって、前記方法が、

i) 複数の入力 c D N A 配列を得ること、 20

i i) 前記複数の入力 c D N A 配列を対になった増幅プライマーと接触させること、ここで、前記対になった増幅プライマー内の少なくとも 1 つのプライマーは、1 つの鎖上に内部 d U を含むアダプタ配列を含み、少なくとも 1 ラウンドの増幅を行い、それにより、適合 c D N A 配列の集団を生成させ、

i i i) 前記適合 c D N A 配列の集団をウラシル D N A グリコシラーゼ及びエンドヌクレアーゼ V I I I と接触させ、それにより、一本鎖末端を有する適合 c D N A 配列の集団を形成すること、

i v) 前記一本鎖末端を有する適合 c D N A 配列の集団をリガーゼと接触させ、それによって線状アレイ核酸配列を形成すること、

v) 前記線状アレイ核酸配列から、任意選択でロングリード配列決定によって配列情報 30

を得ること、及び

v i) 前記線状アレイ核酸配列から得られた前記配列情報を分析して、アイソフォーム配列決定情報を得ること、

を含み、

それにより、前記入力 c D N A 配列の集団からアイソフォーム配列決定情報を得る、方法。

【請求項 26】

入力ミトコンドリア c D N A 配列の集団からミトコンドリア系統追跡を行うための方法であって、前記方法が、

i) 複数の入力ミトコンドリア c D N A 配列を得ること、 40

i i) 前記複数の入力ミトコンドリア c D N A 配列を対になった増幅プライマーと接触させること、ここで、前記対になった増幅プライマー内の少なくとも 1 つのプライマーが、1 つの鎖上に内部 d U を含むアダプタ配列を含み、少なくとも 1 ラウンドの増幅を行い、それにより、適合ミトコンドリア c D N A 配列の集団を生成し、

i i i) 前記適合ミトコンドリア c D N A 配列の集団をウラシル D N A グリコシラーゼ及びエンドヌクレアーゼ V I I I と接触させ、それにより、一本鎖末端を有する適合ミトコンドリア c D N A 配列の集団を形成すること、

i v) 前記一本鎖末端を有する適合ミトコンドリア c D N A 配列の集団をリガーゼと接触させ、それによりアレイ核酸配列を形成すること、

v) 配列情報を前記アレイ核酸配列から、任意選択でロングリード配列決定によって得 50

ること、及び

v i) 前記アレイ核酸配列から得られた前記配列情報を分析してミトコンドリア系統を追跡すること、

を含み、

それにより、前記入力ミトコンドリア c D N A 配列の集団からミトコンドリア系統追跡を行う、方法。

【請求項 27】

アレイ核酸配列を調製するための方法であって、前記方法が、

i) 複数の入力核酸配列を得ること、ここで、前記複数の入力配列内の各入力核酸配列が約 300 キロベース以下の長さであり、

i i) 前記複数の入力核酸配列を対になった増幅プライマーと接触させること、ここで、前記対になった増幅プライマー内の少なくとも 1 つのプライマーが 1 つの鎖上に内部 d U を含むアダプタ配列を含み、少なくとも 1 ラウンドの増幅を実施し、それにより、適合核酸配列の集団を生成し、

i i i) 前記適合核酸配列の集団をウラシル D N A グリコシラーゼ及びエンドヌクレアーゼ V I I I と接触させ、それにより一本鎖末端を有する適合核酸配列の集団を形成すること、及び

i v) 前記一本鎖末端を有する適合核酸配列の集団をリガーゼと接触させること、を含み、

それにより、アレイ核酸配列を形成する、方法。

【請求項 28】

アレイ核酸配列を調製するための方法であって、前記方法が、

i) 複数の入力核酸配列を得ること、ここで、前記複数の入力配列内の各入力核酸配列が約 300 キロベース以下の長さであり、

i i) 前記複数の入力核酸配列を、一本の鎖上の内部 d U を含むアダプタ配列及びリガーゼと接触させ、それにより、適合核酸配列の集団を生成すること、

i i i) 前記適合核酸配列の集団をウラシル D N A グリコシラーゼ及びエンドヌクレアーゼ V I I I と接触させ、それにより一本鎖末端を有する適合核酸配列の集団を形成すること、及び

i v) 前記一本鎖末端を有する適合核酸配列の集団をリガーゼと接触させること、を含み、

それにより、線状アレイ核酸配列を形成する、方法。

【請求項 29】

前記複数の入力配列内の各入力核酸配列は、約 30 キロベース以下の長さである、請求項 27 又は 28 に記載の方法。

【請求項 30】

複数の核酸配列を含み、前記複数の核酸配列の少なくとも 2 つが、配列番号 1 ~ 18 からなる群から選択されるアダプタ配列を含む、組成物。

【請求項 31】

配列番号 1 ~ 18 からなる群から選択される複数のアダプタ配列及びその使用説明書を含むキット。

【請求項 32】

核酸配列リードの集団の個々の核酸配列リード内の別個の配列要素を同定するための方法であって、前記個々の核酸配列リードは、配列要素の線状アレイを有し、

前記配列要素の線状アレイの各々は、高複雑度のライブラリから引き出された 2 つ以上の核酸配列要素を含み、高複雑度のライブラリから引き出された各核酸配列要素は、低複雑度のライブラリから引き出された 1 つ若しくは複数の予想される核酸配列、又は低複雑度のライブラリから引き出された 1 つ若しくは複数の予想される核酸配列及び配列リード終端に隣接し、前記方法は、

(a) 前記核酸配列リードの集団の配列データに 1 つ又は複数の統計的アノテーション

10

20

30

40

50

モデルを適用して、前記核酸配列リードの集団内で、高複雑度のライブラリから引き出された個々の核酸配列要素の領域及び低複雑度のライブラリから引き出された核酸配列の領域を予測すること、ここで、前記1つ又は複数の統計的アノテーションモデルが、

i) 核酸配列リード全体に散在する1つ又は複数の予想される核酸配列を認識するための生成統計的アライメントモデル、

ii) 既知ではない、又は高複雑度の配列の辞書から引き出された配列を認識するためのランダム統計アライメントモデル、を含み、

予測された転位部位は、各モデルの末端に配置され、前記生成統計的アライメントモデルの内部位置内では許容されず、

(b) 複数の核酸配列リードに対して工程(a)を繰り返し、それにより、前記1つ又は複数の統計的モデルを前記複数の核酸配列リードの各核酸配列リードに順相補性配向及び逆相補性配向の両方で適用し、最大対数尤度値を有するモデルを同定することによって、選択された最大事後状態経路の最終的リード当たりのモデル(maximum a posteriori state path Final per-read model) 1) 選択を決定すること、及び

(c) 前記複数の核酸配列リードの各核酸配列リードを、工程(b)の最大事後状態経路の最終的リード当たりのモデル選択によって同定される転位部位によって区画された別個の配列要素にセグメント化すること、

を含み、

それにより、前記核酸配列リードの集団内の別個の配列要素を同定する、方法。 20

【請求項32】

前記高複雑度のライブラリが、1,000を超える異なる要素、任意選択で10,000を超える異なる要素を含むか、又は含む可能性がある、請求項32に記載の方法。

【請求項33】

前記高複雑度のライブラリ及び/又は知られていない若しくは高複雑度の配列の辞書から引き出された配列が、cDNA転写物配列、バーコード配列及び固有の分子識別子からなる群から選択される要素を含む、請求項32に記載の方法。

【請求項34】

前記低複雑度のライブラリが、100個以下の異なる配列、任意選択的に50個以下の異なる配列、任意選択的に25個以下の異なる配列、任意選択的に15個以下の異なる配列を含む、請求項32に記載の方法。 30

【請求項35】

前記低複雑度のライブラリがアダプタ及び/又はリンカー配列を含む、請求項35に記載の方法。

【請求項36】

前記先験的に予想される核酸配列が、アダプタ及び/又はリンカー配列を含む、請求項32に記載の方法。

【請求項37】

先験的に知られていない配列、又は高複雑度の配列の辞書から引き出された配列が、cDNA配列、バーコード配列及び固有の分子識別子配列からなる群から選択される1つ又は複数の配列を含み、任意選択で前記バーコード配列が単一細胞バーコード配列を含む、請求項32に記載の方法。 40

【請求項38】

複数の核酸配列リードの個々の配列リード内の別個の配列要素を同定し、配列要素データを保存するためのシステムであって、前記システムは、

ネットワークと通信するための1つ又は複数のネットワークインターフェース；

前記ネットワークインターフェースに結合される、1つ又は複数のプロセスを実行するように構成されたプロセッサ；及び

前記プロセッサによって実行可能なプロセスを保存するように構成された非一時的メモリを含み、

前記プロセスは、実行されると、

(a) 配列要素の線状アレイを有する個々の核酸配列リードを含む複数の核酸配列リードを取得する、

ここで、配列要素の線状アレイを有する各リードが高複雑度のライブラリから引き出された2つ以上の個々の核酸配列要素を含み、高複雑度のライブラリから引き出された各核酸配列要素が低複雑度の1つ若しくは複数の予想される核酸配列、又は低複雑度の1つ若しくは複数の予想される核酸配列及び配列リード終端のいずれかに隣接している、

(b) 核酸配列リード内で、高複雑度のライブラリから引き出された個々の核酸配列要素の前記複数の領域及び低複雑度のライブラリから引き出された核酸配列の領域を予測するために、1つ又は複数の統計的アノテーションモデルを前記複数の核酸配列リードの配列データに適用する、

10

ここで、前記1つ又は複数の統計的アノテーションモデルは、

i) 核酸配列リード全体に散在する1つ又は複数の予想される核酸配列を認識するための生成統計的アライメントモデル、及び

i i) 既知ではない配列、又は高複雑度の配列の辞書から引き出された配列を認識するためのランダム統計アライメントモデルを含み、

前記生成統計的アライメントモデルでは、予測された転位部位は各モデルの末端に配置され、かつ内部位置内では許容されない、

(c) 複数の核酸配列リードに対して工程 (a) を繰り返し、それにより、前記1つ又は複数の統計的モデルを順相補配向及び逆相補配向の両方で前記複数の核酸配列リードの各核酸配列リードに適用し、最大対数尤度値を有するモデルを同定することによって選択された最終リード当たりのモデル選択により、各モデルの最大事後状態経路を決定し、それにより、前記核酸配列リード内の既知のセグメントを標識する、

20

(d) 前記複数の核酸配列リードの各核酸配列リードを、工程 (c) の最大事後状態経路の最終的リード当たりのモデルによって同定される転位部位によって区画された標識された既知のセグメントの個別の配列要素にセグメント化し、

それにより、前記複数の核酸配列リード内の別個の配列要素を同定する、及び

(e) 前記複数の核酸配列リード内で同定された前記別個の配列要素を配列要素データファイルに保存する、

ように構成される、システム。

30

【請求項39】

前記高複雑度のライブラリが、1,000を超える異なる要素、任意選択で10,000を超える異なる要素を含むか、又は含む可能性がある、請求項39に記載のシステム。

【請求項40】

前記高複雑度のライブラリ及び/又は先験的に知られていない配列若しくは高複雑度の配列の辞書から引き出された配列が、cDNA転写物配列、バーコード配列及び固有の分子識別子からなる群から選択される要素を含む、請求項39に記載のシステム。

【請求項41】

前記低複雑度のライブラリが、100個以下の異なる配列、任意選択的に50個以下の異なる配列、任意選択的に25個以下の異なる配列、任意選択的に15個以下の異なる配列を含む、請求項39に記載のシステム。

40

【請求項42】

前記低複雑度のライブラリがアダプタ及び/又はリンカー配列を含む、請求項42に記載のシステム。

【請求項43】

前記先験的に予想される核酸配列が、アダプタ及び/又はリンカー配列を含む、請求項39に記載のシステム。

【請求項44】

先験的に知られていない配列、又は高複雑度の配列の辞書から引き出された配列が、cDNA配列、バーコード配列及び固有の分子識別子配列からなる群から選択される1つ又

50

は複数の配列を含み、任意選択で前記バーコード配列が単一細胞バーコード配列を含む、請求項 39 に記載のシステム。

【請求項 45】

複数の核酸配列リードの個々の配列リードを、低品質として同定し、除去し、配列データを保存するためのシステムであって、前記システムは、

ネットワークと通信するための 1 つ又は複数のネットワークインターフェース；

前記ネットワークインターフェースに結合される、1 つ又は複数のプロセスを実行するように構成されたプロセッサ；及び

前記プロセッサによって実行可能なプロセスを保存するように構成された非一時的メモリを含み、

前記プロセスは、実行されると、

i) 複数の核酸配列リードの個々の配列リードに対して請求項 39 に記載の工程 (a) ~ (e) を実施する、

ii) ライブラリ調製より予想される順序で起こらない別個の配列要素を含む任意のリードを低品質として同定し、除去する、

ここで、最初の別個の配列要素の後で開始するが、残りの別個の配列要素が順番であるリード、及び最後の別個の配列要素の前で終わるが、前のセクションが全て順番であるリード、並びにこれらの場合の組み合わせは除去されず、

iii) 低品質リードが除去された前記複数の核酸配列リードを配列データファイルに保存する、ように構成される、システム。

【請求項 46】

サーキュラーコンセンサスシーケンシングソフトウェアによって高品質であると同定された 1 つ又は複数の核酸配列リードが低品質であると同定され、除去される、請求項 46 に記載のシステム。

【請求項 47】

更なる分析のために十分に高品質の個々の配列リードを同定し、複数の核酸配列リードの個々の配列リードを配列データに付加し、配列データを保存するためのシステムであって、前記システムは、

ネットワークと通信するための 1 つ又は複数のネットワークインターフェース；

前記ネットワークインターフェースに結合される、1 つ又は複数のプロセスを実行するように構成されたプロセッサ；及び

前記プロセッサによって実行可能なプロセスを保存するように構成された非一時的メモリ、を含み、

前記プロセスは、実行される場合、

i) 複数の核酸配列リードの各々における各ヌクレオチドである複数の核酸配列リードの個々の配列リードに対して請求項 39 に記載の工程 (a) ~ (e) を実施する、

ii) 更なる分析のための十分に高い品質に関して、最初の予想されるセグメントの後で開始するが残りのセクションが順番であるリード、及び最後の予想されるセグメントの前で終わるが前のセクションが順番であるリード、並びにこれらの場合の任意の組み合わせを含む、ライブラリ調製により出現すると予想される順序で標識されたセクションを含む任意のリードを同定する、及び

iii) 更なる分析のために十分に高品質であると同定された前記核酸配列リードを配列データファイルに保存する、ように構成される、システム。

【請求項 48】

サーキュラーコンセンサスシーケンシングソフトウェアによって低品質であると同定された 1 つ又は複数の核酸配列リードが、更なる分析に対して十分に高品質であると同定される、請求項 48 に記載のシステム。

【請求項 49】

請求項 46 に記載の低品質として同定されたリード、又は請求項 48 に記載の高品質として同定されたリードの品質を概算し、推定品質スコアをデータに付加し、データを保存

10

20

30

40

50

するためのシステムであって、

前記システムは、

ネットワークと通信するための1つ又は複数のネットワークインターフェース；

前記ネットワークインターフェースに結合される、1つ又は複数のプロセスを実行するように構成されたプロセッサ；及び

前記プロセッサによって実行可能なプロセスを保存するように構成された非一時的メモリ、を含み、

前記プロセスは、実行されると、

(i) 請求項 4 6 に記載の低品質として同定された各リード又は請求項 4 8 に記載の高品質として同定された各リードでの各別個の配列要素について、別個の配列要素内のヌクレオチドと個別の配列要素に対する予想される配列との間の観察されたアライメントスコアを計算し、個別の配列要素内のヌクレオチドと個別の配列要素に対する予想される配列との間の最良の可能なアライメントスコアを計算する；

10

(i i) 各セクションの品質スコアを取得するために、任意選択的に、工程 (i) で計算された前記アライメントスコアを前記最良の可能なアライメントスコアで除算する、

(i i i) 全体的な観察されたアライメントスコアを得るために工程 (i) で計算された全ての観察されたアライメントスコアを合計し、全体的な最良のアライメントスコアを得るために工程 (i) で計算された全ての最良の可能なアライメントスコアを合計し、及び、前記全体的な観察されたアライメントスコアと前記全体的な可能な最良のアライメントスコアとの比を得ることによって、前記核酸配列リードの推定品質スコアを計算する；

20

及び

(i v) 前記核酸配列リードについての前記推定された品質スコアをデータファイルに保存する、ように構成される、システム。

【請求項 5 0】

前記観察されたアライメントスコアが、工程 (i) において、動的プログラミングアルゴリズムを直接的に使用して、又は直接的に前記別個の配列要素と前記予想される配列との間のレーベンシュタイン距離を計算し、前記予想される配列の長さからその距離を減算することによって計算され、任意選択で、前記動的プログラミングアルゴリズムが、Smith-Watermanアルゴリズム、Needleman-Wunschアルゴリズム、及びペア隠れマルコフモデルアルゴリズムからなる群から選択される、請求項 5 0 に記載のシステム。

30

【請求項 5 1】

前記最良の可能なアライメントスコアが、前記予想される配列とそれ自体との間のアライメントスコアを計算することによって得られる、請求項 5 0 に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

関連出願の相互参照

本出願は、「キメラアンプリコンアレイ配列決定 (Chimeric Amplicon Array Sequencing)」と題する、2020年6月15日に提出された米国仮特許出願第63/039,004号の利益を主張する。上記出願の全内容は、参照により本明細書に組み込まれる。

40

【0002】

連邦政府による資金提供を受けた研究に関する記載

本発明は、国立衛生研究所によって授与された助成金番号U19AI082630の下で政府の支援を受けてなされた。政府は、本発明に一定の権利を有する。

【0003】

本発明は、一般に、核酸配列決定のための方法及び組成物、特に配列決定のための核酸集団の調製に関する。

【背景技術】

50

【 0 0 0 4 】

次世代DNA配列決定の出現は生物学的研究に革命をもたらしたが、現在の配列決定プラットフォームによって解決が依然として不十分である多数の重要な遺伝的特徴が存在する。例えば、mRNA成熟中にエクソンの差次的スプライシングを介して遺伝子機能の深く本質的な多様化を可能にするコア生物学的プロセスである選択的スプライシングは、公知の単一細胞配列決定法によって十分に捕捉されていない。腫瘍のクローン進化研究のために、単一細胞のマーカ対立遺伝子からクローン関係を導き出す能力は、頑強な配列決定カバレッジを必要とし、単一細胞遺伝子発現ワークフローでもこれまで達成出来ていない試みを必要とする。更に、潜在する遺伝的障害に起因する疾患には、診断及び病因の解明の両方のためにゲノム組成を忠実に再構築する能力が必要とされる。特に、接合後の変異の結果であり、重度の神経障害に寄与することが知られている体細胞モザイク現象を特徴付けることは、多数の個々の細胞のサンプリングを必要とし、これは現在の方法では扱いにくい作業である。以前に記載されたアプローチではこれらの重要な特徴を解決することができないことは、複雑な生物学的系を忠実に特徴付ける当技術分野の能力が著しく不足していることを強調している。これらの制限は、既知のアプローチが現在の配列決定技術でロングレンジDNA情報を効率的に捕捉することができないことから生じる。したがって、現在のロングリード配列決定プラットフォームでのロングレンジDNA情報の捕捉を最適化することができるアプローチが必要とされている。

10

【 発明の概要 】

【 課題を解決するための手段 】

20

【 0 0 0 5 】

本開示は、少なくとも部分的には、特に、ロングリード配列決定プラットフォームを使用してキメラ核酸に対して核酸配列決定を行うための組成物及び方法に関する。ある態様において、本開示は、ハイスループット構築のための方法及び組成物、並びにロングリード配列決定プラットフォームへの適用のための、（本明細書において「キメラアレイ配列決定」又は「Caseq」と呼ばれるプロセスを介した）核酸のキメラアレイの使用を提供する。そのようなキメラアレイは、以前は不明瞭であった遺伝的特徴の解明、例えば選択的スプライシングの検出；腫瘍クローン進化等のクローン進化の改善された検出；例えば、疾患診断及び疾患病因の解明のための、ゲノム組成の忠実な再構成；体細胞モザイク現象の特徴付け；及びより一般的には改良されたゲノムハプロタイプ評価を可能にする。

30

【 0 0 0 6 】

本開示は、そのロングリードプラットフォームの固有の特徴を利用して、複数の共通配列決定ライブラリの出力を増強するための一般化可能なワークフローを提供する。ロングリードシーケンサは、非常に大きな配列決定出力を有するが（例えば、PacBio（登録商標）Sequel IIは約300GBである）、ラン当たりのリードの総数は限られている（例えば、PacBio（登録商標）Sequel IIは約4Mである）。出力を最大化するために、より小さい断片のライブラリをアレイにアセンブルし、ロングリードシーケンサで効率的に配列決定し、配列決定されたライブラリメンバーの数をアレイ中の断片の数に対して線形に増加させることができる。したがって、本開示のある態様は、単一細胞の遺伝子発現試料からのハイスループット完全転写物配列決定を可能にするという本開示の主な利点を有する、高効率ロングリード配列決定のためのアレイのアセンブリのための合理化され、一般化可能な方法を詳述する。

40

【 0 0 0 7 】

一態様では、本開示は、アレイ核酸配列を調製する方法を提供し、方法は、（i）それぞれが約300キロベース長以下（任意選択で30キロベース長以下）である、複数の入力核酸配列を取得すること、（ii）1つ又は複数のアダプタ配列を複数の核酸配列に付着させ、それにより、適合（adapted）核酸配列の集団を作製すること、（iii）適合核酸配列の集団を、適合核酸配列の集団内の各二本鎖適合核酸配列の少なくとも1つの末端に一本鎖末端を生成することができる酵素と接触させ、それにより一本鎖末端を有する核酸配列の集団を形成すること、及び（iv）一本鎖末端を有する核酸配列の集団

50

をリガーゼと接触させること、を含み、それによりアレイ核酸配列を形成する。

【0008】

いくつかの実施形態において、アダプタ配列の少なくとも1つは、1つの鎖上に内部dUを含む。

【0009】

実施形態では、アレイ核酸配列は、少なくとも20キロベースの長さを有する。任意選択で、アレイ核酸配列は、少なくとも50キロベースの長さを有する。関連する実施形態では、アレイ核酸配列は、約100キロベース以上の長さを有する。

【0010】

一実施形態では、複数の入力核酸配列は、約0.5 kb ~ 20 kbの長さである。

10

【0011】

ある実施形態において、複数の入力核酸配列は、1つ又は複数のcDNAライブラリから得られる。任意選択で、複数の入力核酸配列は、1つ又は複数の単一細胞又は空間cDNAライブラリから得られる。

【0012】

実施形態では、工程(i i)は、複数の核酸配列を対になった増幅プライマーと接触させること、この際、対になった増幅プライマーの少なくとも1つが1つの鎖上の内部dUを含むアダプタ配列を含む、及び、少なくとも1ラウンドの増幅を実行すること、を含み、それにより適合核酸配列の集団を生成する。

【0013】

いくつかの実施形態では、増幅プライマーの各対の少なくとも1つがビオチン化されている。任意選択で、アダプタ配列テールアンプリコンのためのビオチン媒介選択が行われる。

20

【0014】

実施形態では、工程(i i i)は、適合核酸配列の集団をウラシルDNAグリコシラーゼ及びエンドヌクレアーゼV I I Iと接触させ、それにより一本鎖末端を有する核酸配列の集団を形成することを含む。

【0015】

いくつかの実施形態では、アダプタ配列は、5 ~ 30塩基対の長さを含む(標的核酸配列を除く)。任意選択で、アダプタ配列は6 ~ 25塩基対の長さである。任意選択で、アダプタ配列は、構造5' - N 6 - 1 6 _ d U _ t a r g e t - D N A - 3'を有する。

30

【0016】

実施形態では、一方の鎖に内部dUを有するアダプタ配列は配列番号: 1 ~ 18の配列を含む。

【0017】

いくつかの実施形態では、アダプタ配列を有する複数の核酸配列について、各アダプタ配列は、アダプタ配列を有する複数の核酸配列のうち少なくとも1つの他のものと相補的な1つ又は2つの指定配列を有し、それにより、複数のアダプタ配列は相補的なアダプタ配列の集団を形成する。任意選択で、相補的なアダプタ配列の集団の各相補的なアダプタ配列は、相補的なアダプタ配列の集団の互いに相補的なアダプタ配列に対して最小の類似性を有する。関連する実施形態では、相補的なアダプタ配列の集団の各相補的なアダプタ配列は、相補的なアダプタ配列の集団の他の全ての相補的なアダプタ配列から少なくとも11ハミング距離単位離れている。

40

【0018】

ある実施形態では、以下の1つ又は複数がサイズ選択される: 複数の入力核酸配列; 適合核酸配列の集団; 及び/又は一本鎖末端を有する核酸配列の集団。任意選択で、サイズ選択は電気泳動を介して行われる。関連する実施形態では、サイズ選択は、アガロースゲルを使用して行われる。

【0019】

一定の実施形態では、アレイ核酸配列の配列情報が得られる。任意選択で、アレイ核酸

50

配列の配列情報は、ロングリード配列決定プラットフォームを使用して得られる。

【0020】

関連する実施形態では、ハプロタイプフェージングの配列情報がアレイ核酸配列にわたって得られる。

【0021】

別の実施形態では、形成されるアレイ核酸配列は、5つ以上の入力核酸配列を含む。任意選択で、形成されるアレイ核酸配列は、6個以上、7個以上、8個以上、9個以上、10個以上、11個以上、12個以上、13個以上、14個以上、15個以上、16個以上、17個以上、18個以上、19個以上、又は20個以上の入力核酸配列を含む。

【0022】

ある実施形態において、標的化アイソフォーム配列決定情報は、複数の入力核酸配列を得る工程(i)の間に遺伝子パネルの標的化を介して得られる。

【0023】

実施形態では、複数の入力核酸配列は、免疫応答経路のためのcDNAを含む。

【0024】

いくつかの実施形態では、複数の入力核酸配列は、ミトコンドリアDNAから得られる。任意選択で、アレイ核酸配列の配列決定は、ミトコンドリアDNA系統追跡に使用される。

【0025】

ある実施形態では、適合核酸配列の集団は、ギブソンアセンブリを介して結合される。

【0026】

いくつかの実施形態では、アレイ核酸配列は線状アレイである。

【0027】

ある実施形態において、アレイ核酸配列は、環状アレイである。

【0028】

本開示の更なる態様は、入力cDNA配列の集団からアイソフォーム配列決定情報を得るための方法を提供し、方法は、(i)複数の入力cDNA配列を得ること、(ii)複数のcDNA配列を対になった増幅プライマーと接触させ、それにより、適合cDNA配列の集団を生成させること、この際、対になった増幅プライマーのうち少なくとも1つは1つの鎖上に内部dUを含むアダプタ配列を提示し、少なくとも1回の増幅を行い、(iii)適合cDNA配列の集団をウラシルDNAグリコシラーゼ及びエンドヌクレアーゼVIIIIと接触させ、それにより一本鎖末端を有する適合cDNA配列の集団を形成すること、(iv)一本鎖末端を有する適合cDNA配列の集団をリガーゼと接触させ、それにより線状アレイ核酸配列を形成すること、(v)線状アレイ核酸配列から配列情報を得ること(任意選択で、配列は、ロングリード配列決定によって得られる)、及び(vi)線状アレイ核酸配列から得られた配列情報を分析して、アイソフォーム配列決定情報を得ること、を含み、それにより、入力cDNA配列の集団からアイソフォーム配列決定情報を得る。

【0029】

本開示の別の態様は、入力ミトコンドリアcDNA配列の集団からミトコンドリア系統追跡を行うための方法を提供し、方法は、(i)複数の入力ミトコンドリアcDNA配列を得ること、(ii)複数のミトコンドリアcDNA配列を対になった増幅プライマーと接触させ、それによって適合ミトコンドリアcDNA配列の集団を生成させること、この際、対になった増幅プライマーのうち少なくとも1つが1つの鎖上に内部dUを含むアダプタ配列を含み、少なくとも1ラウンドの増幅を行い、(iii)適合ミトコンドリアcDNA配列の集団をウラシルDNAグリコシラーゼ及びエンドヌクレアーゼVIIIIと接触させ、それにより一本鎖末端を有する適合ミトコンドリアcDNA配列の集団を形成すること、(iv)一本鎖末端を有する適合ミトコンドリアcDNA配列の集団をリガーゼと接触させ、それによりアレイ核酸配列を形成すること、(v)アレイ核酸配列から配列情報を取得すること(任意選択で、配列は、ロングリード配列決定によって得られる)

10

20

30

40

50

、及び(v i)線状アレイ核酸配列から得られた配列情報を分析してミトコンドリア系統を追跡すること、を含み、それによって入力ミトコンドリアcDNA配列の集団に対してミトコンドリア系統追跡を実施する。本開示の更なる態様は、核酸配列の線状アレイのアレイを調製する方法を提供し、方法は、(i)本明細書に開示されるC A s e q方法によって入力核酸配列の第1の集団から第1の線状アレイを調製すること、(i i)本明細書に開示されるC A s e q法によって入力核酸配列の第2の集団から第2の線状アレイを調製すること、この際、第1の線状アレイ及び第2の線状アレイがそれぞれ適合する相補的フランキング配列を有し、(i i i)第1の線状アレイ及び第2の線状アレイを溶液中で組み合わせること、及び(i v)溶液中の第1の線状アレイ及び第2の線状アレイをリガーゼと接触させること、を含み、それにより、核酸配列の線状アレイのアレイを形成する

10

【0030】

ある実施形態では、第1の線状アレイ若しくは第2の線状アレイ、又はその両方は、線状アレイのアレイを含む。

【0031】

いくつかの実施形態では、方法は更に、(v)本明細書に開示されるC A s e q法によって入力核酸配列の第3の集団から第3の線状アレイを調製すること、この際、線状アレイ及び第3の線状アレイのアレイはそれぞれ、適合する相補的フランキング配列を有する、(v i)線状アレイ及び第3の線状アレイのアレイを溶液中で組み合わせること、及び、(v i i)溶液中の線状アレイのアレイ及び第3の線状アレイをリガーゼと接触させること、を含み、それにより、核酸配列の線状アレイのより大きなアレイを形成する。任意選択的に、工程(v)~(v i i)は、第4の線状アレイ、第5の線状アレイ、及び/又はより多くの線状アレイを線状アレイのより大きなアレイに組み込むために繰り返される。

20

【0032】

本開示の別の態様は、アレイ核酸配列を調製する方法を提供し、方法は、(a)複数の入力核酸配列を得ること、この際、各入力配列は、約300キロベース以下の長さであり、(b)複数の核酸配列を、一本の鎖上の内部dUを含むアダプタ配列及びリガーゼと接触させ、それによって適合核酸配列の集団を生成すること、(c)適合核酸配列の集団をウラシルDNAグリコシラーゼ及びエンドヌクレアーゼV I I Iと接触させ、それにより一本鎖末端を有する核酸配列の集団を形成すること、及び(d)一本鎖末端を有する核酸配列の集団をリガーゼと接触させること、を含み、それによりアレイ核酸配列を形成する。

30

【0033】

更なる態様では、本開示は、アレイ核酸配列を調製するための方法を提供し、方法は、(i)複数の入力核酸配列を得ること、この際、各入力配列が約300キロベース以下の長さである；(i i)複数の核酸配列を、1つの鎖上に内部dUを有するアダプタ配列と接触させ、少なくとも1回の増幅を行い、それにより適合核酸配列の集団を生成すること；(i i i)適合核酸配列の集団をウラシルDNAグリコシラーゼ及びエンドヌクレアーゼV I I Iと接触させ、それにより一本鎖末端を有する核酸配列の集団を形成すること；及び(i v)一本鎖末端を有する核酸配列の集団をリガーゼと接触させること、を含み、それにより線状アレイ核酸配列を形成する。

40

【0034】

実施形態では、複数の入力配列内の各入力核酸配列は、約30キロベース以下の長さである。

【0035】

本開示の更なる態様は、複数の核酸配列を含む組成物を提供し、複数の核酸配列の少なくとも2つは、配列番号1~18から選択されるアダプタ配列を含む。

【0036】

本開示の別の態様は、配列番号1~18から選択される複数のアダプタ配列、及びその

50

使用説明書を含むキットを提供する。

【0037】

本開示の更なる態様は、核酸配列リードの集団の個々の核酸配列リード内の別個の配列要素を同定するための方法を提供し、個々の核酸配列リードは、配列要素の線状アレイを有し、配列要素の線状アレイの各々は、高複雑度のライブラリから引き出された2つ以上の核酸配列要素を含み、高複雑度のライブラリから引き出された各核酸配列要素は、低複雑度のライブラリから引き出された1つ若しくは複数の予想される核酸配列の、又は低複雑度のライブラリから引き出された1つ若しくは複数の予想される核酸配列及び配列リード末端のいずれかに隣接 (flanked) し、前記方法は：(a) 核酸配列リードの集団の配列データに1つ又は複数の統計的アノテーションモデルを適用して、高複雑度のライブラリから引き出された個々の核酸配列要素の領域及び低複雑度のライブラリから引き出された核酸配列リードの領域を集団内で予測すること、この際、前記1つ又は複数の統計的アノテーションモデルは、i) 核酸配列リード全体に散在する1つ又は複数の予想される核酸配列を認識するための生成統計的アライメントモデル (generative statistical alignment model)、及び、ii) 既知ではない配列又は高複雑度の配列の辞書から引き出された配列を認識するためのランダム統計的アライメントモデル (random statistical alignment model) を含み、予測された転位部位は各モデルの末端に配置され、生成統計的アライメントモデルの内部位置内では許容されず；(b) 複数の核酸配列リードに対して工程(a)を繰り返し、それによって前記1つ又は複数の統計的モデルを複数の核酸配列リードの各核酸配列リードに順相補性配向及び逆相補性配向の両方で適用し、最大対数尤度値を有するモデルを同定することによって選択された最大事後状態経路の最終的リード当たりのモデル (maximum a posteriori state path Final per-read model) 選択を決定すること；及び、(c) 複数の核酸配列リードの各核酸配列リードを、工程(b)の最大事後状態経路の最終的リード当たりのモデル選択によって同定される転位部位によって区画された別個の配列要素にセグメント化すること、を含み、それにより核酸配列リードの集団内の別個の配列要素を同定する。

10

20

【0038】

一実施形態では、高複雑度のライブラリは、1,000を超える異なる要素を含むか、又は潜在的に含む。任意選択的に、高複雑度のライブラリは、10,000を超える異なる要素を含むか、又は潜在的に含む。

30

【0039】

別の実施形態では、高複雑度のライブラリ及び/又は先験的に知られていない配列、又は高複雑度の配列の辞書から引き出された配列は、cDNA転写物配列、バーコード配列、及び/又は固有の分子識別子である要素を含む。

【0040】

ある実施形態において、低複雑度のライブラリは、100個以下の異なる配列を含む。任意選択で、低複雑度のライブラリは、50個以下の異なる配列を含む。任意選択で、低複雑度のライブラリは、25個以下の異なる配列を含む。任意選択で、低複雑度のライブラリは、15個以下の異なる配列を含む。

40

【0041】

いくつかの実施形態において、低複雑度のライブラリは、アダプタ及び/又はリンカー配列を含む。

【0042】

実施形態では、先験的に予想される核酸配列は、アダプタ及び/又はリンカー配列を含む。

【0043】

ある実施形態において、先験的に知られていない配列又は高複雑度の配列の辞書から引き出された配列は、以下のタイプの配列：cDNA配列、バーコード配列及び/又は固有の分子識別子 (unique molecular identifier) 配列のうち

50

の1つ又は複数を含む。任意選択的に、バーコード配列は、単一細胞バーコード配列を含む。

【0044】

本開示の別の態様は、複数の核酸配列リードの個々の配列リード内の別個の配列要素を同定し、配列要素データを保存するためのシステムを提供し、システムは、ネットワークと通信するための1つ又は複数のネットワークインターフェース；ネットワークインターフェースに結合され、1つ又は複数のプロセスを実行するように構成されたプロセッサ；及び、プロセッサによって実行可能なプロセスを保存するように構成された非一時的メモリを含み、プロセスは、実行されると、(a)配列要素の線状アレイを有する個々の核酸配列リードを含む複数の核酸配列リードを得るように、この際、配列要素の線状アレイを有する各リードは、高複雑度のライブラリから引き出された2つ以上の個々の核酸配列要素を含み、高複雑度のライブラリから引き出された各核酸配列要素は、低複雑度の1つ若しくは複数の予想される核酸配列に、又は低複雑度の1つ若しくは複数の予想される核酸配列及び配列リード終端のいずれかに隣接する；(b)高複雑度のライブラリから引き出された個々の核酸配列要素の複数の領域及び低複雑度のライブラリから引き出された核酸配列の領域の核酸配列リード内で予測するために、1つ又は複数の統計的アノテーションモデルを複数の核酸配列リードの配列データに適用するように、この際、1つ又は複数の統計的アノテーションモデルは、i)核酸配列リード全体に散在する1つ又は複数の予想される核酸配列を認識するための生成統計的アライメントモデル、及び、ii)既知ではない配列又は高複雑度の配列の辞書から引き出された配列を認識するためのランダム統計的アライメントモデル含み、予測された転位部位が各モデルの末端に配置され、生成統計的アライメントモデルの内部位置内では許容されない；(c)複数の核酸配列リードに対して工程(a)を繰り返し、それにより、1つ又は複数の統計的モデルを複数の核酸配列リードの各核酸配列リードに順相補性配向及び逆相補性配向の両方で適用し、モデルを最大対数尤度値で同定することによって選択された最大事後状態経路の最終的リード当たりのモデル選択を決定し、それにより、核酸配列リード内の既知のセグメントを標識するように；及び、(d)複数の核酸配列リードの各核酸配列リードを、工程(c)の最大事後状態経路の最終的リード当たりのモデル選択によって同定される転位部位によって区画された(標識された既知のセグメントの)個別の配列要素にセグメント化し、それにより、複数の核酸配列リード内の個別の配列要素を同定するように；及び、(e)複数の核酸配列リード内で同定された別個の配列要素を配列要素データファイルに保存するように、構成される。

10

20

30

【0045】

本開示の更なる態様は、複数の核酸配列リードの個々の配列リードを低品質として識別し、除去し、配列データを保存するためのシステムを提供し、システムは、ネットワークと通信するための1つ又は複数のネットワークインターフェース；ネットワークインターフェースに結合され、1つ又は複数のプロセスを実行するように構成されたプロセッサ；及び、プロセッサによって実行可能なプロセスを保存するように構成された非一時的メモリ、を含み、プロセスが実行される場合、i)複数の核酸配列リードの個々の配列リードに対して上記の工程(a)~(e)を実施するように；ii)ライブラリ調製により予想される順序では生じない別個の配列要素を有する任意のリードを同定及び除去するように、この際、最初の別個の配列要素の後に始まるが、残りの別個の配列要素が順番になっているリード、及び最終的に予想される別個の配列要素の前に終わるが、前のセクションが全て順番になっているリード、並びにこれらの場合の組み合わせは除去されず；及び、iii)低品質リードが除去された複数の核酸配列リードを配列データファイルに保存するように、構成される。

40

【0046】

ある実施形態において、Circular Consensus Sequencingソフトウェアが高品質であると識別した個々の配列リードは、この方法によって低品質であると識別される。

50

【 0 0 4 7 】

本開示の別の態様は、更なる分析のために十分に高い品質の個々の配列リードを同定し、複数の核酸配列リードの個々の配列リードを配列データに付加し、配列データを保存するためのシステムを提供し、システムは：ネットワークと通信するための1つ又は複数のネットワークインターフェース；ネットワークインターフェースに結合され、1つ又は複数のプロセスを実行するように構成されたプロセッサ；及び、プロセッサによって実行可能なプロセスを保存するように構成された非一時的メモリ、を含み、プロセスは実行される場合、i) 複数の核酸配列リードの個々の配列リードに対して上記の工程(a)~(e)を実施し、最初に予想される別個の配列要素の後に始まるが、残りの別個の配列要素が順番になっているリード、及び最後に予想される別個の配列要素の前に終わるが、以前の別個の配列要素が順番になっているリード、並びにこれらの場合の任意の組合わせを含む、ライブラリ調製より出現すると予想される順序で別個の配列要素を有する任意のリードを、更なる分析のために十分に高品質であると識別するように；及び、v) 更なる分析のために十分に高品質であると識別された核酸配列リードを配列データファイルに保存するように、構成される。

10

【 0 0 4 8 】

ある実施形態において、Circular Consensus Sequencingソフトウェアが低品質であると識別した個々の配列リードは、この方法によって高品質であると識別される。

【 0 0 4 9 】

本開示の最終態様は、新たに識別された高品質及び低品質リードの品質を概算し、推定品質スコアをデータに追加し、データを保存するためのシステムを提供し、システムは、ネットワークと通信するための1つ又は複数のネットワークインターフェース；ネットワークインターフェースに結合され、1つ又は複数のプロセスを実行するように構成されたプロセッサ；及び、プロセッサによって実行可能なプロセスを保存するように構成された非一時的メモリ、を含み、プロセスは、実行されると：(i) 各新しく識別された高品質又は低品質のリード内の各別個の配列要素について、別個の配列要素内のヌクレオチド間の観察されたアライメントスコア及び別個の配列要素に対する予想される配列を計算し、別個の配列要素内のヌクレオチドと別個の配列要素に対する予想される配列のヌクレオチドとの間の最良のアライメントスコアを計算するように；(ii) 任意選択的に、工程(i)で計算されたアライメントスコアを最良のアライメントスコアで除算して、各セクションの品質スコアを得るように；及び、(iii) 工程(i)で計算された全ての観察されたアライメントスコアを合計して、全体的な観察されたアライメントスコアを得、工程(i)で計算された全ての最良の可能なアライメントスコアを合計して、全体的な最良のアライメントスコアを得；全体的な観察されたアライメントスコアと全体的な最良のアライメントスコアとの比を得ることによって、核酸配列リードの推定品質スコアを計算するように；及び、(iv) 核酸配列リードについての推定品質スコアをデータファイルに保存するように、構成される。

20

30

【 0 0 5 0 】

ある実施形態では、アライメントスコアは、動的プログラミングアルゴリズムを直接使用して、又は別個の配列要素と予想される配列との間のレーベンシュタイン距離を計算し、その距離を予想される配列の長さから減算することによって直接、工程(a)で計算される。任意選択で、動的プログラミングアルゴリズムは、Smith-Waterman(ローカル)アルゴリズム、Needleman-Wunsch(グローバル)アルゴリズム、又は類似/同等のアライメントアルゴリズム(例えば、ペア隠れマルコフモデル(Pair Hidden Markov Model))のうちの1つ又は複数を含む。

40

【 0 0 5 1 】

いくつかの実施形態において、最良のアライメントスコアは、予想される配列とそれ自体との間のアライメントスコアを計算することによって得られる。

【 0 0 5 2 】

50

定義

本明細書で使用される場合、特に明記されない限り、又は文脈から明らかでない限り、「約」という用語は、当技術分野における通常の許容範囲内、例えば平均の2標準偏差以内であると理解される。「約」は、記載された値の10%、9%、8%、7%、6%、5%、4%、3%、2%、1%、0.5%、0.1%、0.05%、又は0.01%以内と理解することができる。

【0053】

ある実施形態では、「およそ」又は「約」という用語は、特に明記しない限り、又は文脈から明らかでない限り（そのような数が可能な値の100%を超える場合を除き）、記載された基準値のいずれかの方向（より大きい又はより小さい）において25%、20%、19%、18%、17%、16%、15%、14%、13%、12%、11%、10%、9%、8%、7%、6%、5%、4%、3%、2%、1%、又はそれ未満に入る値の範囲を指す。

10

【0054】

文脈から明らかでない限り、本明細書で提供される全ての数値は、「約」という用語によって修飾される。

【0055】

「対照」又は「参照」とは、比較の基準を意味する。対照試料を選択及び試験する方法は、当業者の能力の範囲内である。統計学的有意性の決定は、当業者の能力の範囲内であり、例えば、陽性結果を構成する平均からの標準偏差の数である。

20

【0056】

本明細書で使用される場合、「異なる」という用語は、核酸に関して使用される場合、核酸が互いに同じではないヌクレオチド配列を有することを意味する。2つ以上の核酸は、それらの全長に沿って異なるヌクレオチド配列を有することができる。あるいは、2つ以上の核酸は、それらの長さのかなりの部分に沿って異なるヌクレオチド配列を有することができる。例えば、2つ以上の核酸は、2つ以上の分子について異なる標的ヌクレオチド配列部分を有することができるが、2つ以上の分子上で同じであるユニバーサル配列部分も有することができる。

【0057】

本明細書で使用される場合、「各」という用語は、アイテムの集合に関して使用される場合、集合内の個々のアイテムを識別することを意図しているが、必ずしも集合内の全てのアイテムを指すとは限らない。明示的な開示又は文脈がそうでないことを明確に指示する場合、例外が発生する可能性がある。

30

【0058】

本明細書で使用される場合、単一細胞核酸配列決定は、試料中の細胞又は他の種類の核酸の配列を測定し、その細胞及び/又は試料核酸が得られた個々の細胞及び/又は供給源を同定する方法を指す。同様に、単一細胞RNA配列決定は、細胞RNA（任意選択で、転写物）の配列を測定し、その細胞RNAが得られた個々の細胞を同定する方法を指す。

【0059】

本明細書で使用される場合、「アンプリコン」という用語は、核酸に関して使用される場合、核酸を複製する産物を意味し、産物は、核酸のヌクレオチド配列の少なくとも一部と同じ又は相補的なヌクレオチド配列を有する。アンプリコンは、例えば、ポリメラーゼ伸長、ポリメラーゼ連鎖反応（PCR）、ローリングサークル増幅（RCA）、多重置換増幅（MDA）、ライゲーション伸長、又はライゲーション連鎖反応を含む、核酸又はそのアンプリコンを鋳型として使用する様々な増幅方法のいずれかによって産生され得る。アンプリコンは、特定のヌクレオチド配列の単一コピー（例えば、PCR産物）又はヌクレオチド配列の複数コピー（例えば、RCAのコンカテマー生成物）を有する核酸分子であり得る。標的核酸の第1アンプリコンは、典型的には相補的コピーである。後続のアンプリコンは、第1のアンプリコンの生成後に、標的核酸又は第1のアンプリコンから作製されるコピーである。後続のアンプリコンは、標的核酸に実質的に相補的であるか又は標

40

50

的核酸と実質的に同一である配列を有することができる。

【0060】

本明細書で使用される場合、「アレイ」という用語は、相対的な位置によって互いに区別することができる特徴又は部位の集団を指す。アレイの異なる部位にある異なる分子は、アレイ内の部位の位置に応じて互いに区別することができる。アレイの個々の部位は、特定の種類の1つ又は複数の分子を含むことができる。例えば、部位は、特定の配列を有する単一の核酸分子を含むことができ、又は部位は、いくつかの核酸分子を含むことができる。ある実施形態では、「線状アレイ (linear array)」という用語は、より大きな線状の核酸分子に沿ったアレイの別個の位置における配列要素の線状の集合体を指すために使用される。

10

【0061】

本明細書で使用される場合、「バーコード配列」という用語は、核酸、核酸の特徴（例えば、同一性）、又は核酸に対して行われた操作を識別するために使用することができる核酸中の一連のヌクレオチドを意味することを意図している。バーコード配列は、天然に存在する配列又はバーコード化核酸が得られた生物に天然には存在しない配列であり得る。バーコード配列は、集団中の単一の核酸種に固有であり得るか、又はバーコード配列は、集団中のいくつかの異なる核酸種によって共有され得る。更なる例として、集団中の各核酸プローブは、集団中の他の全ての核酸プローブとは異なるバーコード配列を含むことができる。あるいは、集団中の各核酸プローブは、集団中のいくつか又はほとんどの他の核酸プローブからの異なるバーコード配列を含み得る。例えば、集団中の各プローブは、共通のバーコードを有するプローブがそれらの長さに沿った他の配列領域において互いに異なる場合であっても、集団中のいくつかの異なるプローブについて存在するバーコードを有することができる。特定の実施形態では、生物学的検体（例えば、組織試料）と共に使用される1つ又は複数のバーコード配列は、生物学的検体のゲノム、トランスクリプトーム又は他の核酸には存在しない。例えば、バーコード配列は、特定の生物学的検体中の核酸配列に対して80%、70%、60%、50%又は40%未満の配列同一性を有し得る。

20

【0062】

本明細書で使用される場合、「伸長する」という用語は、核酸に関して使用される場合、核酸への少なくとも1つのヌクレオチド又はオリゴヌクレオチドの付加を意味すること
を意図している。特定の実施形態では、1つ又は複数のヌクレオチドを、例えばポリメ
ラーゼ触媒作用（例えば、DNAポリメラーゼ、RNAポリメラーゼ又は逆転写酵素）を介
して核酸の3'末端に付加することができる。化学的又は酵素的方法を使用して、核酸の
3'又は5'末端に1つ又は複数のヌクレオチドを付加することができる。1つ又は複数の
オリゴヌクレオチドを、例えば、化学的又は酵素的（例えばリガーゼ触媒）方法によっ
て、核酸の3'末端又は5'末端に付加することができる。核酸は、鋳型指向的に伸長するこ
とができ、それによって伸長産物は、伸長される核酸にハイブリダイズする鋳型核酸に相
補的である。

30

【0063】

本明細書で使用される場合、「逆転写酵素」という用語は、RNA鋳型から相補的DN
A (cDNA) を生成するために使用される酵素を指す。当技術分野で一般的に使用され
る逆転写酵素 (RT) には、非鎖置換転写酵素 RTX、及びウイルス逆転写酵素 M - ML
V が含まれる。

40

【0064】

本明細書で使用される場合、「増幅する」、「増幅」又は「増幅反応」及びそれらの派
生語は、一般に、核酸分子の少なくとも一部が少なくとも1つの更なる核酸分子に複製又
はコピーされる任意の作用又はプロセスを指す。追加の核酸分子は、任意選択で、鋳型核
酸分子の少なくとも一部と実質的に同一又は実質的に相補的な配列を含む。鋳型核酸分子
は一本鎖又は二本鎖であり得、追加の核酸分子は独立して一本鎖又は二本鎖であり得る。
増幅は、任意選択で、核酸分子の線状の又は指数関数的複製を含む。いくつかの実施形態

50

において、そのような増幅は、等温条件を用いて行うことができ、他の実施形態では、そのような増幅は熱サイクリングを含むことができる。いくつかの実施形態において、増幅は、単一の増幅反応における複数の標的配列の同時増幅を含む多重増幅である。増幅反応は、当業者に公知の増幅プロセスのいずれかを含むことができる。いくつかの実施形態では、増幅反応は、1つ又は複数の核酸配列を増幅するポリメラーゼ連鎖反応（PCR）を含む。そのような増幅は、線状又は指数関数的であり得る。いくつかの実施形態では、増幅条件は、等温条件を含むことができ、あるいは熱サイクリング条件、又は等温条件と熱サイクリング条件との組み合わせを含むことができる。いくつかの実施形態では、1つ又は複数の核酸配列を増幅するのに適した条件は、ポリメラーゼ連鎖反応（PCR）条件を含む。典型的には、増幅条件は、ユニバーサル配列に隣接する1つ又は複数の標的配列等の核酸を増幅するために、又は1つ又は複数のアダプタに連結された増幅標的配列を増幅するために十分な反応混合物を指す。一般に、増幅条件は、増幅又は核酸合成のための触媒、例えばポリメラーゼ；増幅される核酸に対してある程度の相補性を有するプライマー；及び、核酸にハイブリダイズするとプライマーの伸長を促進するための、ヌクレオチド、例えばデオキシリボヌクレオチド三リン酸及びリボヌクレオチド三リン酸を含む。増幅条件は、プライマーの核酸へのハイブリダイゼーション又はアニーリング、プライマーの伸長、及び伸長されたプライマーが増幅を受ける核酸配列から分離される変性工程を必要とし得る。本明細書で使用される場合、「ポリメラーゼ連鎖反応」（「PCR」）という用語は、目的のポリヌクレオチドのセグメントの濃度を増加させるための方法を記載している、Mullisの米国特許第4,683,195号及び同第4,683,202号の方法を指す。本明細書で使用される場合、「増幅標的配列」及びその派生語は、一般に、標的的特異的プライマー及び本明細書で提供される方法を使用して標的配列を増幅することによって生成される核酸配列を指す。増幅された標的配列は、標的配列に関して同じセンス（すなわち、プラス鎖）又はアンチセンス（すなわち、マイナス鎖）のいずれかであり得る。

10

20

【0065】

本明細書で使用される場合、「サーキュラーコンセンサスシーケンシング（Circular Consensus Sequencing）ソフトウェア低品質リード」という用語は、サーキュラーコンセンサスシーケンシングソフトウェアが0.99未満のリード品質スコアを割り当てる配列決定リード、又はサーキュラーコンセンサスシーケンシングソフトウェアが「ZMWパスフィルタ」以外のカテゴリにリードを割り当てるリードを指す。

30

【0066】

本明細書で使用される場合、「サーキュラーコンセンサスシーケンシングソフトウェア高品質リード」という用語は、サーキュラーコンセンサスシーケンシングソフトウェアが「ZMWパスフィルタ」カテゴリにリードを割り当てる配列リードを指す。ある実施形態では、CCSソフトウェア高品質リードは、CCSソフトウェアが0.99以上のリード品質スコアを割り当てたリードである。

【0067】

本明細書で使用される場合、「高複雑度のライブラリ」という用語は、特定のライブラリ要素が所与の位置に存在するかどうかの先験的な予測を統計的に不確実にする（例えば、所与の場所における特定のライブラリ要素の可能性は1%未満、所与の場所における特定のライブラリ要素の可能性は0.1%未満等である）のに十分な数の異なる要素（異なる配列、サイズ、長さ等を有する要素）を含むか、又は潜在的に含むライブラリを指す。ある実施形態では、「高複雑度のライブラリ」は、100を超える別個の要素、任意選択で1000を超える別個の要素、任意選択的に10,000を超える別個の要素、及び/又は任意選択的に100,000を超える別個の要素を含むか、又は潜在的に含む。実施形態では、「高複雑度のライブラリ」はcDNA配列ライブラリ、任意選択でゲノムcDNA配列ライブラリを指す。いくつかの実施形態において、「高複雑度のライブラリ」は、後の処理工程（例えば、バーコード配列（任意選択で、単一細胞バーコード配列、ピー

40

50

ズバーコード配列等)、固有の分子識別子等)において異なる検討に値するほど大きな配列の辞書から引き出されたライブラリを指す。

【0068】

本明細書で使用される場合、「低複雑度のライブラリ」という用語は、特定のライブラリ要素が所与の位置に存在するかどうかの先験的予測を、限られた統計的不確実性のみで可能にするために(例えば、特定のライブラリ要素が所与の場所で発生する可能性は1%超、特定のライブラリ要素が所与の場所で発生する可能性は5%超、特定のライブラリ要素が所与の場所で発生する可能性は20%超等である)、十分に少数の別個の要素(異なる配列、サイズ、長さ等を有する要素)を含むか、又は潜在的に含むライブラリを指す。ある実施形態では、「低複雑度のライブラリ」は、100個未満の異なる要素、任意選択的に50個未満の異なる要素、任意選択的に30個未満の異なる要素、及び/又は任意選択的に15個未満の異なる要素を含むか、又は潜在的に含む。実施形態では、「低複雑度のライブラリ」は、リンカー及び/又はアダプタ配列ライブラリを指す。

10

【0069】

本明細書中で使用されるとき、用語「ライゲーションすること」、「ライゲーション」及びそれらの派生語は、一般に、2つ以上の分子を互いに共有結合的に連結するためのプロセス、例えば、2つ以上の核酸分子を互いに共有結合的に連結するためのプロセスのことを指す。いくつかの実施形態において、ライゲーションは、核酸の隣接するヌクレオチド間にニックをつなぐことを含む。いくつかの実施形態では、ライゲーションは、第1の核酸分子の末端と第2の核酸分子の末端との間に共有結合を形成することを含む。いくつかの実施形態では、ライゲーションは、1つの核酸の5'リン酸基と第2の核酸の3'ヒドロキシル基との間に共有結合を形成し、それによってライゲーションされた核酸分子を形成することを含み得る。一般に、本開示の目的のために、ライブラリ配列(任意選択で増幅されたライブラリ配列)をアダプタ配列にライゲーションして(又はそうでなければプライマー媒介増幅を介して付着させて)アダプタ連結配列を生成することができ、次いで、これを更に操作して、異なる配列要素を線状アレイ核酸に連結することができる。

20

【0070】

本明細書で使用される場合、「リガーゼ」及びその派生語は、一般に、2つの基質分子のライゲーションを触媒することができる任意の薬剤を指す。いくつかの実施形態において、リガーゼは、核酸の隣接するヌクレオチド間のニックの連結を触媒することができる酵素を含む。いくつかの実施形態では、リガーゼは、1つの核酸分子の5'リン酸と別の核酸分子の3'ヒドロキシルとの間の共有結合の形成を触媒し、それによってライゲーションされた核酸分子を形成することができる酵素を含む。適切なリガーゼには、T4 DNAリガーゼ、T7 DNAリガーゼ、Taq DNAリガーゼ、及び大腸菌(E. coli) DNAリガーゼが含まれ得るが、これらに限定されない。

30

【0071】

本明細書で使用される場合、「ライゲーション条件」及びその派生語は、一般に、2つの分子を互いに連結するのに適した条件を指す。

【0072】

本明細書中で使用されるとき、用語「次世代配列決定」又は「NGS」とは、従来の配列決定方法(例えば、標準的なサンガー又はマクサム-ギルバート配列決定法)を使用したときには前例のない速度でポリヌクレオチドを配列決定する能力を有する配列決定技術のことを指し得る。これらの前例のない速度は、数千から数百万の配列決定反応を並行して実行し、読み出すことによって達成される。NGS配列決定プラットフォームとしては、限定されないが、以下が挙げられる: Massively Parallel Signature Sequencing (Lynx Therapeutics); 454 pyro-sequencing (454 Life Sciences/Roche Diagnostics); solid-phase, reversible dye-terminator sequencing (Solexa/Illumina (商標)); SOLiD (商標)技術 (Applied Biosystems); Ion s

40

50

emiconductor sequencing (Ion Torrent (商標)) ; 及び DNA nanoball sequencing (Complete Genomics)。ある NGS プラットフォームの説明は、以下に見出すことができる: Sheendure, et al., ' ' Next-generation DNA sequencing, ' ' Nature, 2008, vol. 26, No. 10, 135 - 145 ; Mardis, ' ' The impact of next-generation sequencing technology on genetics, ' ' Trends in Genetics, 2007, vol. 24, No. 3, pp. 133 - 141 ; Su, et al., ' ' Next-generation sequencing and its applications in molecular diagnostics ' ' Expert Rev Mol Diagn, 2011, 11 (3) : 333 - 43 ; 及び Zhang et al., ' ' The impact of next-generation sequencing on genomics ' ' , J Genet Genomics, 2011, 38 (3) : 95 - 109。

【0073】

本明細書で使用される場合、「核酸」及び「ヌクレオチド」という用語は、当技術分野でのそれらの使用と一致し、天然に存在する種又はその機能的類似体を含むことを意図している。核酸の特に有用な機能的類似体は、配列特異的な様式で核酸にハイブリダイズすることができるか、又は特定のヌクレオチド配列の複製のための鋳型として使用することができる。

【0074】

天然に存在する核酸は、一般に、ホスホジエステル結合を含む骨格を有する。類似体構造は、当技術分野で公知の様々なもののいずれかを含む代替の骨格連結を有することができる。天然に存在する核酸は、一般に、デオキシリボース糖（例えば、デオキシリボ核酸 (DNA) に見られる）又はリボース糖（例えば、リボ核酸 (RNA) に見られる）を有する。核酸は、当技術分野で公知のこれらの糖部分の様々な類似体のいずれかを含むヌクレオチドを含むことができる。核酸は、天然又は非天然ヌクレオチドを含むことができる。これに関して、天然デオキシリボ核酸は、アデニン、チミン、シトシン又はグアニンからなる群から選択される1つ又は複数の塩基を有することができ、リボ核酸は、ウラシル、アデニン、シトシン又はグアニンからなる群から選択される1つ又は複数の塩基を有することができる。核酸又はヌクレオチドに含めることができる有用な非天然塩基は、当技術分野で公知である。「プローブ」又は「標的」という用語は、核酸又は核酸の配列に関して使用される場合、本明細書に記載の方法又は組成物の文脈における核酸又は配列の意味的識別子として意図され、核酸又は配列の構造又は機能を、他に明示的に示されるものを超えて必ずしも限定しない。

【0075】

本明細書で使用される場合、「プライマー」という用語及びその派生語は、一般に、目的の標的配列にハイブリダイズすることができる任意の核酸を指す。典型的には、プライマーは、ポリメラーゼによってヌクレオチドを重合することができるか、又はインデックス等のヌクレオチド配列をライゲーションすることができる基質として機能するが、いくつかの実施形態では、プライマーは、合成された核酸鎖に組み込まれ、別のプライマーがハイブリダイズして、合成された核酸分子に相補的な新しい鎖の合成を開始することができる部位を提供することができる。プライマーは、ヌクレオチド又はその類似体の任意の組み合わせを含むことができる。いくつかの実施形態において、プライマーは、一本鎖オリゴヌクレオチド又はポリヌクレオチドである。「ポリヌクレオチド」及び「オリゴヌクレオチド」という用語は、任意の長さのヌクレオチドのポリマー形態を指すために本明細書で互換的に使用され、リボヌクレオチド、デオキシリボヌクレオチド、それらの類似体、又はそれらの混合物を含み得る。この用語は、等価物として、DNA、RNA 又は cDNA のいずれかの類似体及び二本鎖ポリヌクレオチドを含むと理解されるべきである。本明細書で使用される用語はまた、例えば逆転写酵素の作用によって RNA 鋳型から産生さ

れる相補的又はコピーDNAであるcDNAを包含する。この用語は、分子の一次構造のみを指す。

【0076】

例として与えられるが、説明されるある実施形態のみに本開示を限定することを意図するものではない以下の詳細な説明は、添付の図面と併せて最もよく理解され得る。

【図面の簡単な説明】

【0077】

【図1】図1A~1Cは、アイソフォーム配列決定を効果的に実行するための核酸リード長及びスループット要件を実証し、本明細書に開示される「Caseq」アプローチを提示するグラフを示す。図1Aは、以前に記載された配列決定アプローチがアイソフォーム配列決定領域にギャップを残したことを実証するプロットを示す。具体的には、組み合わせたハイスループット(>20Mリード)及び中間リード長(0.5~5kb)配列決定アプローチは存在せず、本Caseqアプローチは対処するために本明細書で提供されている。図1Bは、本明細書に開示される線状核酸アレイが、ロングリードプラットフォーム上で配列決定され、配列決定されたDNA分子の全出力をアレイあたりの断片の数に多重化して、それらの個々の全長DNA断片に逆多重化され得ることを示す(現在のグラフに示すように3倍であるが、有効配列出力の10倍以上の多重化を容易に達成することができる)。図1Cは、デオキシウラシル(dU)消化を用いて断片の協調的アセンブリを駆動する技術によって、アレイへのDNAアンプリコンの制御された不偏ライゲーションがどのようにして本明細書において達成されたかの描写を示す。例示されるように、DNAライブラリを、5'「相補配列」とそれに続くdUとを含有するプライマーを用いて増幅する。増幅後、dU含有アンプリコンをウラシルDNAグリコシラーゼ及びエンドヌクレアーゼVIIで消化すると、dUが除去され、DNAの残りの上流鎖が融解し、それによって一本鎖「相補配列」が露出する。次いで、これらのdU消化アンプリコンは、相補的な「相補配列」を含むアンプリコンとハイブリダイズして、標的化されたアセンブリを駆動することができる。アレイ長は、生成される「オーバーラップ配列」断片の数によって単純に変調される。

10

20

【図2】図2A及び2Bは、本開示のCaseqプロセスを使用して、1.2kbの平均断片サイズを有するcDNAライブラリからの8断片多重化アセンブリについて得られた結果を示す。図2Aは、例示されたようなCaseqプロセスが、表示されたcDNAサイズ分布(開始、ライゲーション及び配列決定/逆多重化cDNA)に従って、ライゲーション時に約10kbの多重化断片をもたらしたことを示す。図2Bは、SequelIIで配列決定された多重化ライブラリについて得られた結果を示す図であり、これは、逆多重化後に約23Mの転写物を伴う合計約2.5Mのリードをもたらし、これは、これまでに知られているアプローチを超えるスループットの約9倍の増加を表したことを示す。逆多重化されたリードの分析により、元のcDNAライブラリと同様のサイズ分布が確認された(図2Aに見られるように)。

30

【図3】図3A及び図3Bは、本開示のキメラアレイの完全な配列内容を、そのようなキメラアレイに存在する構造を利用する様式で解明することに関連する、ヒトゲノムにわたる遺伝子及び転写物の長さの分布を示す。図3Aは、ヒトゲノムにわたるタンパク質コード遺伝子転写物(左側の緑色の点)及び遺伝子(黒色の点、右の分布)についてのカウント及び長さの分布を示す。ヒトタンパク質をコードする遺伝子転写物の大部分は10kb未満の長さであり、事実上全てのタンパク質をコードする転写物は100kb未満の長さであるが、遺伝子のかなり大部分は10kbの長さを超え、かなりの数の遺伝子は100kbの長さを超え、数は1Mbの長さを超える。図3Bは、長さが増加するにつれて累積頻度をより明確に示すように表される、タンパク質コード遺伝子転写産物長(左側の緑色の点)及び遺伝子(黒色の点、右の分布)のヒトゲノムにおける累積分布(頻度)を示す。80%のヒトタンパク質コード遺伝子転写物は、5000塩基未満を含有すると特異的に認められた。

40

【図4】図4は、それぞれSpike-In RNAバリエーション(SIRV)で実施した

50

場合の、長鎖リード配列分析のための既存の「Smart-seq3」プロセスと本開示のキメラアンプリコンアレイ配列決定分析との混同行列の比較を示す。SIRVは、ヒト遺伝子と同様に選択的にスプライシングされる7つのSIRV遺伝子(SIRV1~SIRV7)に分けられる。各遺伝子の転写物群は、四角で囲まれた領域によって示される。影付きの四角は、データ間の類似性を示す。対角線(左上から右下)は、SIRV転写物の自己類似性を示す。Smart-seq3で生成されたデータは、各SIRV遺伝子の個々の転写物を区別することが困難であることが観察されたが、本開示のキメラアンプリコンアレイ配列決定法及び分析によって生成されたデータは、配列決定されたSIRV転写物にほぼ完全にマッピングして戻された。

【図5】図5は、本開示のキメラアンプリコンアレイ配列決定方法及びヒトT細胞試料に対して行われた分析の全体的な収率のサンキーダイアグラムを示す。本開示の計算的逆多重化方法及び低品質リード再生方法を組み合わせたライブラリ調製は、既存のCCS補正HiFiリードプロセス(すなわち、「Smart-seq3」)のみを使用する方法と比較して、データ収率の全体的な21.85倍の増加をもたらした。

【図6】図6は、本開示のキメラアンプリコンアレイ配列決定法を用いて調製したヒトT細胞試料中のアダプタライゲーションのヒートマップを示す。カウントは、各列に示されるオーバーハングアダプタから各行に示されるオーバーハングアダプタまでのライゲーションの数を示す。逆方向に補完された配列は、'記号で示されている。この特定のライブラリでは、アレイサイズは15であり、予想されるライゲーション順序はA->B->C->D->E->F->G->H->I->J->K->L->M->N->O->Pであった。対角線に沿った高いカウント(1つ下にシフト)は、調製されたライブラリ全体にわたって予想されるライゲーションの極めて高い割合を示す。中央の切れ目は、プロットが向きを切り替える場所である(逆相補ライゲーションを別々に示すため)。「ホット対角」上にない正方形のほとんどのカウントは0であり、予想外に検出されたライゲーションを示す正方形の最大のカウントでさえ、「ホット対角」のカウントよりも最大で3桁小さい。

【図7】図7は、予想されるライゲーション順序A->B->C->D->E->F->G->H->I->J->K->L->M->N->O->Pでの長さ15アレイライブラリ調製の上位20のライゲーションプロファイル(有病率による)を示す。逆の相補的アダプタは'記号で示されている。これらのデータは、本明細書に現在開示されているキメラアレイの分析方法によって未だフィルタリングされていない。

【図8】図8は、2つのヒトT細胞試料にわたる、直接配列決定と、本開示のキメラアンプリコンアレイの配列決定方法及び分析の使用との比較を示す。

【図9】図9A及び図9Bは、本開示の方法によって調製及び分析されたキメラアンプリコンアレイについての高品質及び低品質のアダプタライゲーションのヒートマップをそれぞれ示す。図9Aは、本開示のキメラアンプリコンアレイ配列決定法を用いて調製したヒトT細胞試料中の高品質アダプタライゲーションのヒートマップを示す。カウントは、各列に示されるオーバーハングアダプタから各行に示されるオーバーハングアダプタまでのライゲーションの数を示す。逆方向に補完された配列は、'記号で示されている。この特定のライブラリでは、アレイサイズは15であり、予想されるライゲーション順序はA->B->C->D->E->F->G->H->I->J->K->L->M->N->O->Pであった。高品質のデータを、本開示のキメラアンプリコンアレイ配列決定分析プロセス(「Longbow」と呼ばれる)によって決定した。図9Bは、本開示のキメラアンプリコンアレイ配列決定法を用いて調製したヒトT細胞試料中の低品質アダプタライゲーションのヒートマップを示す。カウントは、各列に示されるオーバーハングアダプタから各行に示されるオーバーハングアダプタまでのライゲーションの数を示す。逆方向に補完された配列は、'記号で示されている。この特定のライブラリでは、アレイサイズは15であり、予想されるライゲーション順序はA->B->C->D->E->F->G->H->I->J->K->L->M->N->O->Pであった。低品質のデータを、本開示のキメラアンプリコンアレイ配列決定分析プロセス(「Longbow」)に

10

20

30

40

50

よって決定した。対角線上に生じないライゲーションは多数存在するが、低品質のデータであってもほぼ全てのライゲーションが予想通りに生じた。

【図10】図10A～図10Dは、COVID-19患者と健康な対照(HC)との間で行われた比較から得られた転写物データのクラスタリング評価を提示するt分布型確率の近傍埋め込み(t-distributed Stochastic Neighbor Embedding: t-SNE)プロットを示し、これにより、健康な患者と軽度及び重度のCOVID-19を有する患者との間の単球区画における著しい転写の違いが識別された。t-SNEプロットは、健康な人及びCOVID-19患者の血液試料の評価から得られ、本明細書に開示のCaseqプロセスを介して得られた遺伝子アイソフォーム情報をショートリードデジタル遺伝子発現データにどのように補足できるかを実証している。図10Aは、表現型によってクラスター化されたt-SNE分析プロットを示す。図10Bは、試料によってクラスタリングされたt-SNE分析プロットを示す。図10Cは、ライデンクラスタリング(leiden clustering)を使用して実行されたt-SNE分析のプロットを示す。図10Dは、細胞タイプによってクラスター化されたt-SNE分析プロットを示す。

【図11】図11A～図11Cは、末梢血単核球(PBMC)試料から得られた結果を示す。図11Aは、免疫細胞型を同定するために使用される、PBMC試料からの標準的なショートリード遺伝子発現データのクラスタリングの結果を示す。図11Bは、同じ試料からの遺伝子(ショートリード)及びアイソフォーム(ロングリード)発現データの統合を示す。図11Cは、図11Bに示される遺伝子(ショートリード)及びイソ型(ロングリード)の発現データの統合により、カノニカルCD45(PTPRC)アイソフォームの細胞型特異的発現が明らかにされたことを示す。

【図12】図12は、本開示のシステムを示す。

【図13】図13は、本開示の1つ又は複数の実施形態による最大の状態経路を決定するための例示的な手順を示す。

【発明を実施するための形態】

【0078】

本開示は、少なくとも部分的には、核酸配列の入力集団に見出され得る不偏であり、及び/又はバイアスを最小化する方法で、ロングリード配列決定プラットフォームのスループット及び/又は収率を増強するための方法及び組成物に関する。したがって、ある態様において、特に、ロングリード配列決定プラットフォームを使用してキメラ核酸に対して核酸配列決定を行うための方法が提供される。ある実施形態において、本方法の核酸の線状キメラアレイは、ロングリード配列決定プラットフォームへの適用に有用である。そのような線状キメラアレイは、以前は不明瞭であった遺伝的特徴の解明、例えば選択的スプライシングの検出；腫瘍クローン進化等のクローン進化の改善された検出；例えば、疾患診断及び疾患病因の解明のための、ゲノム組成の忠実な再構成；体細胞モザイク現象の特徴付け；及びより一般的には改良されたゲノムハプロタイプ評価を可能にする。

【0079】

本開示は、特に、ロングリードプラットフォームの固有の特徴を利用して、複数の共通配列決定ライブラリの出力を増強するための一般化可能なワークフローを提供する。ロングリードシーケンサは、非常に大きな配列決定出力を有するが(例えば、PacBio(登録商標)Sequel IIは約300GBである)、ラン当たりのリードの総数は限られている(例えば、PacBio(登録商標)Sequel IIは約4Mである)。出力を最大化するために、より小さい断片のライブラリをアレイにアセンブルし、ロングリードシーケンサで効率的に配列決定し、配列決定されたライブラリメンバーの数をアレイ中の断片の数に対して線状に増加させることができる。したがって、本開示のある態様は、単一細胞遺伝子発現試料からのハイスループット完全転写物配列決定を可能にするという本開示の主な利点を有する、高効率ロングリード配列決定のためのアレイの組立てのための合理化され、一般化可能な方法を詳述する。

【0080】

10

20

30

40

50

近年、単一細胞遺伝子発現研究の劇的な増加が見られているが、そのような研究の注目すべき欠点は、これまで、そのような試みにおいてアイソフォーム組成又は遺伝的変異を解決することができなかったことである。ハイスループット単一細胞配列決定/発現分析における完全長転写物情報の捕捉における制限は、これらのワークフローにおけるハイスループットショートリード配列決定への依存を反映している。ショートリードアプローチは、転写産物の5'末端又は3'末端からの小さな約100bpのスナップショットを効果的に配列決定し、 1×10^8 を超える転写産物から遺伝子カウントを効率的に取得するのに十分であるが、遺伝子アイソフォーム組成又は遺伝的変異を捕捉するには短すぎる(約5kb以上のリード長を必要とする)。ロングリード配列決定技術における印象的な最近の進歩があるが、それらのスループットは、単一細胞遺伝子発現試料から完全長転写物を適切にサンプリングするには依然として不十分である。したがって、ある態様において、本明細書中に提供されるのは、これらの制限を克服するための合理化された方法であり、この方法は、ある態様では、ロングリード配列決定プラットフォームのための核酸配列の精密に設計された線状アレイを作製することに依存し、それにより、本方法は、単一細胞遺伝子発現試料からのハイスループット完全転写物配列決定を可能にする。

10

20

30

40

50

【0081】

上記のように、PacBio(登録商標)及びOxford Nanopore Technologies(「Nanopore」)によって製造された2つの先駆的なロングリード配列決定プラットフォームにおける最近の著しい進歩は、ロングリード配列決定のリード長、スループット、及び精度を劇的に増加させ、単一細胞アイソフォーム配列決定の目標をほぼ手の届くところに置いた。最近の試みはこの2つのロングリード配列決定プラットフォーム(1~3)を活用してきたが、それらのワークフローは、大量のアーチファクト及びスループットの欠如に関連する著しい制限を受ける。これらの非効率性の合計は結果として、トランスクリプトーム内容物のスパース(sparse)サンプリングをもたらしており、これは今日まで、ロングリード配列決定分析の能力を厳しく制限してきた。例えば、Nanoporeアイソフォーム配列決定法であるR2C2(Rolling Circle Amplification to Concatemeric Consensus)は、フィルタを通過する転写物の52%しか達成しないことが観察されており、これはNanoporeフローセルあたり約300,000個の配列決定された転写物(約790ドル)に相当する(2)。PacBio(登録商標)方法、ScISOR-seqも同様にアーチファクトによって制限されており、リードの約36%のみがフィルタを通過し、PacBio(登録商標)1Mフローセルあたり約360,000個の全長転写物(約\$640)になる(1)。これらの欠点は、既知の配列決定技術(図1A)間にこれまで存在していたギャップ、具体的にはハイスループット(>20Mリード)及び中間リード長(0.5~5kb)配列決定が存在しないことを強調している。本開示のある態様は、配列決定アーチファクトを>90%(図1A)減少させながら、ロングリード配列決定プラットフォームのスループットを10倍超増加させることができる方法、キメラアレイ配列決定(CAseq)を提供する。

【0082】

本明細書に開示されるCAseq方法は、これらのプラットフォームの固有の特徴に対処することによってロングリードシーケンサの分子配列決定出力を増強する特殊な多重化ワークフローである。特定のリード長さを有するIllumina(登録商標)のショートリード配列決定ワークフローとは対照的に、ロングリードプラットフォームは、フローセル中で約20kbから莫大な2Mb/ポア(MinION, Oxford Nanopore Technologies)又はウェル(Sequel II, PacBio(登録商標))までの範囲の不確定なリード長さを有する。これらの大量のリード長は、バルク全ゲノム配列決定等の試みには最適であるが、全長転写物等の中間の長さの標的(500bp-10kb)には過剰である。

【0083】

個々のロングリード(図1A)からの複数のDNA標的の配列決定を可能にするキメラ

アレイ配列決定 (Caseq) は、中程度の長さの標的の拡張性のある捕捉のために、本明細書において、ロングリード配列決定プラットフォームをより良好に適合させるために開発された。本Caseq法では、DNA断片の多重化は、マルチフラグメントアレイへの所定数の断片のプログラムされたライゲーションの制御されたプロセスを介して行われる。本明細書に開示される線状核酸アレイは、ロングリードプラットフォーム上で配列決定され、配列決定されたDNA分子の全出力をアレイあたりの断片の数に多重化して、それらの個々の全長DNA断片に逆多重化され得る (図1B)。デオキシウラシル (dU) 消化を用いて断片の協調的アセンブリを駆動する技術によって、アレイへのDNAアンプリコンの制御された不偏ライゲーションが本明細書において達成される。簡潔には、DNAライブラリを、5'「相補配列 (complement sequence)」とそれに続くdUとを含有するプライマーを用いて増幅する。増幅後、dU含有アンプリコンをウラシルDNAグリコシラーゼ及びエンドヌクレアーゼVIIで消化すると、dUが除去され、DNAの残りの上流鎖が融解し、それによって一本鎖「相補配列」が露出する。次いで、これらのdU消化アンプリコンは、相補的な「相補配列」を含むアンプリコンとハイブリダイズして、標的化されたアセンブリを駆動することができる。アレイ長は、生成される「オーバーラップ配列」断片の数によって単純に変調される (図1C)。ひとたび組み立てられると、これらの多重化断片は、その後の配列決定のための標準的なNanopore又はPacBio (登録商標) ライブラリの準備ワークフローに入ることができる。非常に長い又は分子密度の高いアレイを生成するために、アレイを互いに連結してアレイのアレイを作製するようにプログラムすることもできる。特に、相補的配列の最小セットを有する非常に大きい又は高密度の多重化アレイを生成するために、アレイ自体をアレイに連結することができることが明確に企図される。実際には、これは、内部相補的配列の共通のコアセットを有する多数の一次アレイを最初に生成することによって達成することができる。したがって、これらの一次アレイの隣接断片は、一次アレイ間のプログラムされたライゲーションを駆動する固有の相補的配列を含むように設計することができる (一次アレイの初期形成と同様)。

10

20

30

40

【0084】

本明細書に開示されるCaseqプロセスは、限定するものではないが、(1) 10X Genomics (登録商標) のもの等の単一細胞遺伝子発現ワークフロー、例えば発現された核酸のバーコード化集団を構築し、任意選択でゲルビーズに分配することができるプロセス (例えば、PCT/US2018/16019を参照されたい)、(2) 10X Genomics (登録商標) Visium spatial genomics プロセス (特殊化された組織スライド上の捕捉領域内のスポットにグループ化された空間的にバーコード化されたmRNA結合オリゴヌクレオチドを使用するVisium Spatial Gene Expression; mRNAが処理された組織切片から放出されると、それは近傍の捕捉オリゴに結合し、次いで、これらの空間バーコードを組み込み、空間情報を保存するcDNAライブラリを、このmRNAから調製することができ、この遺伝子発現データは、その後、組織切片の高解像度顕微鏡画像上に重ねられ、どの遺伝子が発現されているか、及び組織試料全体のどこで発現されているかを視覚化することを可能にする) 及び例えば米国特許第2021/0123040号に開示されている「Slide-Seq」空間トランスクリプトームプロファイリング手法等の空間配列決定ワークフロー、(3) Caseqを使用する単一細胞遺伝子発現ワークフローから、例えば10X Genomics (登録商標) 試料からのミトコンドリア遺伝子の標的化増幅によって行うことができるミトコンドリア系統追跡、及び(4) とりわけ、B細胞受容体 (BCR) 及びT細胞受容体 (TCR) の高効率の自然対ロングリード配列決定と組み合わせることができるCaseqを含む、任意の数の当技術分野で認識されている技術と組み合わせ使用することもできることが明確に企図される。

【0085】

ある態様では、本開示のCaseq法は、配列又はライブラリの偏りなしに、DNA断片を制御可能かつ効率的に規定の断片番号のアレイに連結する能力を提供する。実施形態

50

では、本アプローチは、一方の鎖上に内部 d U を有する定義された配列（例えば、6 ~ 16 bp 長であるが、他の配列長、例えば、5 ~ 25 bp 又はそれを越える長さも実現性があると考えられる）で標的 DNA の末端を修飾する（例えば、5' - N6 - 16 __ d U __ target - DNA - 3'）。配列の末端は、ウラシル DNA グリコシラーゼ（UDG）及び DNA グリコシラーゼリアーゼエンドヌクラーゼ V I I I（NEB（登録商標）からの USER 酵素カクテル）を用いた d U の塩基切除によって一本鎖にされ、ハイブリダイゼーションのための定義された配列を明らかにする。これらの断片の複数のファミリーを作製及び処理して、ハイブリダイゼーション及びその後のライゲーションを指示することができる。次いで、長い配列断片をロングリードプラットフォームで配列決定することができ、それによって配列決定された分子のそれらの出力を増加させる。アレイ化された配列を調製するための現在の相補的配列媒介方法が本明細書において例示されているが、アレイを作製するための他の経路もまた、線状キメラアレイを作製するために使用されることが明確に企図され、例えば、ギブソンアセンブリ、重複伸長（例えば、遺伝子 SOE）等である。そのような用途のために、それぞれの反応に対する相補的末端配列を含有する増幅された断片がインキュベートされ、必要に応じて適切な条件でサイクルされ、それにより、キメラアレイが作製される。キメラ配列のアセンブリのために制限酵素を使用した、ロングリード核酸配列を作製するための 1 つの以前に開示されたアプローチも留意されるが、制限エンドヌクラーゼ媒介アプローチは、現在の C A s e q プロセスが克服する制限であるライブラリ多様性（Prabakaran et al., Genome Biology 20:134の「SMURF-seq」）の保持において有意な制限を示した。

【0086】

本開示の C A s e q プロセスは、配列決定の分野にわたって広範な適用性を有する。ゲノム配列決定のためには、リード長が長いほど配列再構成がより容易かつより正確になるので、リード長が重要である。ゲノムから 0.5 ~ 20 kb 断片を増幅し、次いで高効率ロングリード配列決定のためのアンプリコンアレイを生成する能力は、ゲノム再構成及びフェージングの精度及び忠実度を高める。C A s e q はまた、このアプローチが DNA のより長い領域からの SNP のフェージングを可能にするため、全エクソーム及び他の標的捕捉配列決定法にも有用である。更に、この C A s e q は、本明細書の他の箇所で更に詳細に説明されるように、アイソフォームの RNA 配列決定に適用可能である。ショートリードシーケンサは、従来の RNA seq ワークフローから RNA アイソフォームを捕捉するのにあまり適していない。ロングリードの最近の試みはスループットが低く、したがってパワー不足である。本開示の C A s e q プロセスは、ロングリード配列決定の出力を有意に増加させ、それによって C A s e q を試料中のアイソフォーム組成を理解するための実行可能なアプローチ、特にアイソフォーム scRNA seq にする。本開示の C A s e q プロセスはまた、TCR : TCR 及び V_H : V_L 対の天然にペアリングした（natively paired）配列決定に有用であり、抗原特異的タグの組込みに適していると考えられる。例えば、本開示の C A s e q プロセスは、全ゲノム及びエクソーム配列決定のための TCR 及び Ig レパートリー並びにライブラリアセンブリのハイスループットによる天然にペアリングした配列決定のための既存のプロセスに適用することができる。具体的には、本開示の C A s e q プロセスは、Tanno et al. (Science Advances, 6(17): eaay9093; DOI: 10.1126/sciadv.aay9093) に記載されているように、現在のワークフローに代わるロングリード配列決定として提供される。Tanno et al. は、TCR : TCR 対又は V_H : V_L 対に対して実施されるエマルジョン系の Overlap Extension RT-PCR によって天然にペアリングした配列決定が達成され、それによってそれらを 1 つの天然にペアリングした断片にスティッチングする方法を記載している。本明細書では、例えば、そのようなペアリングしたアンプリコンのプールを C A s e q ワークフローの入力配列として使用することができ、それによってそのようなペアの拡張性のあるロングリード配列決定を可能にすることが特に企図される。更に、そのようなキ

メラアレイの設計中に他の断片を重複伸長 R T - P C R に組み込むことができ、それにより、個々の細胞からのより多くの情報をそのような T C R : T C R 対及び / 又は V_H : V_L 対と対にし、そのようなアレイからの全ての配列情報の捕捉に必要なロングリード配列決定を与えることができると考えられる。

【 0 0 8 7 】

ある実施形態では、本開示の C A s e q プロセスは、アレイに組み立てられる D N A 分子を生成するための上流処理を最大化するように適合される。例としては、適切なアダプタを有するより大きなサイズの断片 (0 . 5 ~ 2 0 k b) を生成するための D N A の断片化及び増幅の様式的最適化、断片化された D N A からの特定の配列の誘引、及び / 又は標的化されたロングリード配列決定を可能にするための D N A 若しくは R N A からの標的化された増幅が挙げられる。標的化 D N A 又は R N A は、標的核酸のパネルを使用して配列決定の試みを指示することができるので、特に有利であると考えられる。例えば、標的化は、ゲノムの特定の領域のフェージングに特別な注意を払うため、ゲノムの複雑な / 反復的な特徴を解決するため、標的化アイソフォーム増幅のため、及び / 又は本明細書の他の箇所でも論じられるように、単一細胞遺伝子発現 / エピゲノム (A T A C) / ゲノム試料からの腫瘍ミトコンドリア系統追跡のために使用することができる。

10

【 0 0 8 8 】

本開示の特定の方法及び組成物の様々な明示的に企図される成分は、以下で更に詳細に考慮される。

【 0 0 8 9 】

20

核酸ライブラリ

本開示の C A s e q プロセスは、R N A、c D N A 及びゲノム D N A ライブラリを含む任意の核酸ライブラリに効果的に適用することができる。本 C A s e q 法を介して検出及び整列され得る R N A としては、m R N A、s n R N A、l n c R N A、s i R N A 及び g R N A が挙げられ、現在のアプローチでは、C A s e q プロセスを介した整列及び配列決定のために、そのような R N A の安定化された形態及び / 又は対応する D N A 配列を任意選択で使用 / 産生する。

【 0 0 9 0 】

プライマー / アダプタ

本 C A s e q プロセスの例示された態様では、アダプタ配列を入力核酸集団に付着させるためにテールプライマー (t a i l e d p r i m e r) が使用される。使用されるアダプタ配列は、最終的に、個々の入力核酸配列の一本鎖「粘着末端」のアニーリングを介してキメラアレイライゲーションを進行させることを可能にし、それぞれ末端に 1 つ又は 2 つのアダプタ配列が互いに結合している。任意選択で、アダプタ配列内の相補的一本鎖配列の設計は、各キメラアレイが正確な線形順序を有するように行うことができ、又はアダプタ配列の使用は、各キメラアレイ内の線形順序のより大きな柔軟性を可能にし得る。ある例示的な実施形態では、多重ライゲーションのために、15 塩基対 (b p) の相補的配列を増幅し、全長 c D N A ライブラリに付加するための d U 含有プライマーのファミリーが設計されている。アーチファクト配列の主要な供給源に対処するために、例示されたプロセスは、全長 c D N A アンプリコンの精製を可能にするためにビオチン化プライマーを使用した。効率的な多重化アセンブリを駆動し、不適切なライゲーション事象を軽減するために、本明細書に例示される 15 b p 相補的配列は、全ての配列が互いに少なくとも 11 ハミング距離単位離れていることを確実にすることによって、最小の類似性を有するように設計された。そのような品質を有するアダプタ配列の例示的な表を以下の表 1 に示す。

30

40

50

【表 1】

表 1 – 使用されるアダプタ配列の例示的なリスト

配列番号	相補的配列
1	AGCACCATTAATGTGT
2	CTTGTAAGCTGTCTA
3	CTCTGTCAGGTCCGA
4	CCTCCTCCTCCAGAA
5	TCGCTGGTATTCCAA
6	GCTTACTTGTGAAGA
7	TAACCGTATGGTTGA
8	GATGGCGCTATCTCA
9	CTACCAGTGAGGAAG
10	GAGTCCAATTCGCAG
11	ATCAAGGCTTAACGG
12	TGTTGAATCCTAGCG
13	GTGCGTTGCGAATTG
14	CGGTAATGTACCGGC
15	ATTGCGTAGTTGGCC
16	CACTTGGTCGCAATC
17	GTAAGCCTTCGTGTC
18	CCTAGATCAGAGCCT

10

20

【0091】

C A s e q プロセスにおける入力配列へのアダプタ配列の付加は、テール増幅プライマーを使用して本明細書で例示されているが、アダプタ配列を入力配列の集団に付加するための他の当技術分野で認識されている方法も使用することができることが明確に企図されている。例えば、特に断片の増幅（例えば、長さ、又は修正を維持することに起因する）を回避することが有利な場合、線状アレイの構築のための本明細書に開示される C A s e q プロセスの残りの部分の実施前に、入力配列（例えば、平滑末端入力配列）へのアダプタの直接ライゲーションを行うことができる。

30

【0092】

入力核酸の長さ（例えば c D N A）

入力核酸配列の長さは、本開示の具体的な用途に応じて、サイズの範囲が広くなり得る。入力核酸としての c D N A 集団の場合、長さは一般に 0 . 5 k b ~ 2 0 k b に分布する。しかしながら、本方法は、20ヌクレオチド以下という短い入力核酸配列長に、又は最大約メガベース以上の長さを有する入力核酸配列 / 断片に適用することができることが明確に企図される。実際、本開示の C A s e q 法は、例えば、C I T E s e q タグ又は他の生物学的に関連する情報等のライブラリからの捕捉のために、100bp未満の小さい断片に適用することができることが明確に企図される。上記のように、本開示の C A s e q プロセスは、約 350bp ~ 10kb の標準サイズ c D N A にも適用することができる。更に、ロングリード配列決定長が増加し続けると、C A s e q を適用して、多くの大きな (> 10kb) 核酸配列 / 断片の線状アレイを作製できることが明確に企図される。

40

【0093】

ウラシル D N A グリコシラーゼ

本開示のある態様は、ウラシル D N A グリコシラーゼを使用する。ウラシル - D N A グリコシラーゼ (U D G) は、D N A の突然変異を復帰させる酵素である。最も一般的な突然変異は、ウラシルへのシトシンの脱アミノ化である。U D G はこれらの突然変異を修復

50

する。UDGはDNA修復において非常に重要であるが、それがなければ、これらの突然変異はがんをもたらし得る (Pearl, L.H. *Mutat Res.* 460:165-81)。

【0094】

既知のウラシル-DNAグリコシラーゼ及び関連するDNAグリコシラーゼ(EC)としては、ウラシル-DNAグリコシラーゼ(Mol et al. *Cell.* 80:869-78)、好熱性ウラシル-DNAグリコシラーゼ(Sandigursky and Franklin. *Curr. Biol.* 9:531-4)、G:T/Uミスマッチ特異的DNAグリコシラーゼ(Mug)(Barrett et al. *Cell.* 92:117-29)、及び一本鎖選択的単機能性ウラシル-DNAグリコシラーゼ(single-strand selective monofunctional uracil-DNA glycosylase)(SMUG1; Buckley and Ehrenfeld. *J. Biol. Chem.* 262:13599-606)が挙げられる。

【0095】

ウラシルDNAグリコシラーゼは、シトシンの自発的な脱アミノ化によって又はDNA複製中のdAとは反対のdUの誤組込みによってのいずれかにより生じ得るDNAからウラシルを除去する。このファミリーの原型メンバーは、最初に発見されたグリコシラーゼの1つである大腸菌(*E. coli*)UDGである。UNG、SMUG1、TDG、及びMBD4を含む4つの異なるウラシル-DNAグリコシラーゼ活性が哺乳動物細胞で同定されている。それらは、基質特異性及び細胞内局在性が異なる。SMUG1は、基質として一本鎖DNAを好むが、二本鎖DNAからもUを除去する。非修飾ウラシルに加えて、SMUG1は、環C5に酸化基を有する5-ヒドロキシウラシル、5-ヒドロキシメチルウラシル及び5-ホルミルウラシルを切除することができる(Matsubara et al. *Nucleic Acids Res.* 32:5291-5302)。TDG及びMBD4は、二本鎖DNAに厳密に特異的である。TDGは、対向するグアニンが存在する場合、チミングリコール、並びに炭素5に修飾を有するUの誘導体を除去することができる。現在の証拠は、ヒト細胞では、TDG及びSMUG1が、自発的なシトシン脱アミノ化によって引き起こされるU:Gミスペアの修復に参与する主要な酵素であることを示唆しているが、dUの誤組込みによってDNAに生じるウラシルは主にUNGによって処理される。MBD4は、CpG部位における5-メチルシトシンのチミンへの脱アミノ化から生じるT:Gミスマッチを修正すると考えられている(Wu et al. *J. Biol. Chem.* 14:5285-5291)。MBD4変異マウスは正常に発達し、癌感受性の増加又は生存率の低下を示さない。しかし、それらは小腸の上皮細胞のCpG配列でより多くのC:T変異を獲得する(Wong et al. *PNAS.* 99:14937-14942)。制限酵素を使用して(相補的な末端配列を他の断片とアニーリングすることによって)キメラアレイを調製することができることが更に企図される。しかしながら、Caseqプロセスにおける制限酵素の使用は、特定の断片の消化を介してライブラリを偏らせる可能性が非常に高い。

【0096】

エンドヌクレアーゼVII I

本開示の例示されたある態様は、エンドヌクレアーゼVII I酵素を使用する。大腸菌(*E. coli*)由来のエンドヌクレアーゼVII Iは、N-グリコシラーゼ及びAP-リアーゼの両方として作用する。N-グリコシラーゼ活性は、損傷したピリミジンを二本鎖DNAから放出し、アプリン(AP部位)を生成する。AP-リアーゼ活性は、AP部位に対し3'及び5'を切断して、5'リン酸及び3'リン酸を残す。エンドヌクレアーゼVII Iによって認識され除去される損傷を受けた塩基には、尿素、5,6-ジヒドロキシチミン、チミングリコール、5-ヒドロキシ-5-メチルヒダントイン、ウラシルグリコール、6-ヒドロキシ-5,6-ジヒドロチミン及びメチルタルトロニル尿素が含まれる。エンドヌクレアーゼVII IはエンドヌクレアーゼIII Iと類似しているが、エンドヌクレアーゼVII Iは 及び リアーゼ活性を有し、エンドヌクレアーゼIII Iは リア

ーゼ活性のみを有する。

【0097】

リガーゼ

ある態様において、アダプタのオーバーハング末端が Cas 9 プロセスにおいて互いにアニールすると、リガーゼが投与されて、キメラアレイ要素を固定し、要素を線状に取り付ける。リガーゼは、一般に、新しい化学結合を形成することによって、通常、大きな分子のうちの1つの上の小さなペンダント基 (pendant chemical group) の加水分解を伴って、2つの大きな分子の結合を触媒することができる酵素、又は2つの化合物と一緒に結合することを触媒する酵素、例えば、C-O、C-S、C-N等の結合を触媒する酵素を指す。一般に、リガーゼは以下の反応、 $A + B \rightarrow C + D$ 又は場合によっては $A + B + C \rightarrow D + E + F$ を触媒し、小文字はその従属する小さい基を示し得る。リガーゼは、核酸の2つの相補的断片を結合し、複製中に二本鎖 DNA に生じる一本鎖切断を修復することができる。一般的に使用されるリガーゼには、とりわけ、T4 DNA リガーゼ、T7 DNA リガーゼ、Taq DNA リガーゼ、及び大腸菌 (E. coli) DNA リガーゼが含まれるが、これらに限定されない。

10

【0098】

ロングリード配列決定プラットフォーム

本開示のある態様は、ロングリード配列決定を使用する核酸の調製を使用するか、又は含む。ロングリード配列決定 (LRS) は、現在活発に開発されている DNA 配列決定方法の一種である (Bleidorn, Christoph. Systematics and Biodiversity 14: 1-8)。ロングリード配列決定は、DNA の長い鎖を小さなセグメントに切断し、次いで増幅及び合成によってヌクレオチド配列を推測することを必要とする既存の方法とは対照的に、単一分子レベルでヌクレオチド配列を読み取ることによって機能する (「Illumina 配列決定技術」PDF)。

20

【0099】

上で定義した NGS は、その開発以来、DNA 配列決定分野で影響を及ぼしてきた。これにより、ゲノム全体にわたって非常に高いカバレッジで多数のリードをもたらすことができる超並列アプローチを可能にし、DNA 配列決定のコストを劇的に削減した (Treangen and Salzberg. Nature Reviews Genetics 13: 36-46)。

30

【0100】

NGS は、最初に DNA 分子を増幅し、次いで合成によって配列決定を行うことによって機能する。多数の増幅された同一の DNA 鎖の合成の結果得られる集合的な蛍光シグナルは、ヌクレオチド同一性の推論を可能にする。しかしながら、ランダムエラーのために、増幅された DNA 鎖間の DNA 合成は、次第に同期しなくなる。急速に、信号品質は、リード長が増大するにつれて劣化する。リード品質を維持するためには、長い DNA 分子を小さなセグメントに分割しなければならず、NGS 技術の重大な制限をもたらす (Treangen 及び Salzberg)。この課題を克服するための計算の試みは、正確なアセンブリをもたらさない可能性がある近似的なヒューリスティクスに依存することが多い。

40

【0101】

単一の DNA 分子の直接配列決定を可能にすることによって、ロングリード配列決定技術は、第2世代配列決定よりも実質的に長いリードを生成する能力を有する (Bleidorn)。このような利点は、ゲノム科学及び生物学全般の研究の両方に重大な意味を有する。しかしながら、ロングリード配列決定データは、以前の技術よりもはるかに高いエラー率を有し、下流ゲノム構築及び得られたデータの分析を複雑にする可能性がある (Gupta. Trends in Biotechnology 26: 602-611)。これらの技術は活発に開発されており、高いエラーレートの改善が期待されている。構造的変異体コーリング等のエラーレートに対してより耐性がある用途では、ロングリード配列決定が既存の方法よりも優れていることが見出されている。

50

【0102】

いくつかの企業、すなわち Pacific Biosciences、Oxford Nanopore Technology、Quantapore (CA - USA)、及び Stratos (WA - USA) が現在、ロングリード配列決定技術開発の中心にある。これらの企業は、単一の DNA 分子を配列決定するために根本的に異なるアプローチを取っている。

【0103】

PacBio (登録商標) は、ゼロモード導波路の特性に基づいて、単一分子リアルタイム配列決定 (SMRT) の配列決定プラットフォームを開発した。シグナルは、zLウェルの底部に結合した DNA ポリメラーゼによって組み込まれた各ヌクレオチドからの蛍光発光の形態である。本明細書で使用される PacBio (登録商標) ロングリード配列決定プラットフォームの現在の例は、ScISOr-seq である。

10

【0104】

Oxford Nanopore の技術は、DNA 分子をナノスケール細孔構造に通過させ、次いで細孔を取り囲む電場の変化を測定することを含み、一方、Quantapore は異なる独自のナノポアアプローチを有する。Stratos Genomics は、ナノポア ssDNA 読取りのノイズチャレンジに対するシグナルを回避するために、ポリマーインサート「Xpandomers」を用いて DNA 塩基を離間させる。R2C2 (Rolling Circle Amplification to Concatemeric Consensus) は、例示的なナノポアアイソフォーム配列決定方法として注目されている。

20

【0105】

ある実施形態では、ナノポア配列決定が使用される (例えば、参照により組み込まれる Astier et al., J. Am. Chem. Soc. 2006 Feb 8; 128 (5): 1705 - 10 を参照されたい)。ナノポア配列決定の背後にある理論は、ナノポアが導電性流体に浸漬され、それに電位 (電圧) が印加されたときに起こるものと関係がある。これらの条件下では、ナノポアを通るイオンの伝導によるわずかな電流を観察することができ、電流量はナノポアのサイズに非常に敏感である。核酸の各塩基がナノポアを通過すると (又はエキソヌクレアーゼベースの技術の場合には個々のヌクレオチドがナノポアを通過すると)、これにより、4つの塩基のそれぞれについて異なるナノポアを通る電流の大きさが変化し、それによって DNA 分子の配列を決定することが可能になる。

30

【0106】

本開示のある態様は、1つ又は複数の dU 残基で終結し、それぞれの配列要素の線状タンデムアレイを調製するために使用することができる別個の相補的配列を有するように設計された特殊なオリゴヌクレオチドプライマーを使用するが、追加の核酸プライマー/配列/アダプタも本開示の核酸ライブラリに付加することができることも企図される。そのような明示的に企図される更なるプライマー/配列/アダプタとしては、他の識別子及び/又はアダプタ配列の中でも、例えば、CITE-Seq プロセス (Stoeckius et al. Nature Methods, 14: 865 - 868)、REAP-Seq プロセス (Peterson et al. Nature Biotechnology, 35: 936 - 939) 又は他のプロセス、Smith et al. (Smith, A. M. Genome Research 19: 1836 - 1842) 及び他の場所で用いられているもの等の、固有の分子識別子 (UMI) において使用されるもの等の配列バーコードが挙げられるが、これらに限定されない。そのような配列は、CAsEQ プロセスのライゲーション工程の前の任意の時点でライブラリ配列に任意選択的に付加することができ、これにより、それぞれの線状キメラアレイ配列要素の順序が、ロングリード配列決定の実施に先立って固定される。

40

【0107】

バーコード配列及び他の識別配列は、任意の様々な長さであり得る。より長い配列、例

50

例えば、本 C A s e q プロセスによって調製されたものは、一般に、集団に対するより多くの数及び多様なバーコードを収めることができる。一般に、キメラアレイ内の複数の個々の要素は、(異なる配列を有するにもかかわらず)同じ長さのバーコードを有するが、単一のアレイの異なる要素に対して、又は異なる C A s e q ロングリード配列に対して異なる長さのバーコードを使用することも可能である。バーコード配列は、少なくとも 2、4、6、8、10、12、15、20 又はそれを超えるヌクレオチド長であり得る。代替的又は追加的に、バーコード配列の長さは、最大で 20、15、12、10、8、6、4 又はそれ未満のヌクレオチドであり得る。使用することができるバーコード配列の例は、例えば、それぞれ参照により本明細書に組み込まれる米国特許出願公開第 2014/0342921 号及び米国特許第 8,460,865 号に記載されている。

10

【0108】

本開示のあるオリゴヌクレオチドはまた、更なるリンカー(任意選択で切断可能なリンカー)、プライミング部位(当該技術分野で知られているように、例えば、国際公開第 2016/040476 号を参照されたい)ごとに異なる固有の分子識別子(U M I)、上記のバーコード配列、及び任意選択で、P C R 増幅を可能にするための共通配列(「P C R ハンドル」)を含むことができると考えられる。

【0109】

単一細胞配列決定/分子プロファイリング

単一細胞(S C)分子プロファイリング法は、そのような方法が最近主流に変遷し、F A C S のような既存の S C 感受性アプローチと共に変遷しているため、生物医学研究にす

20

【0110】

配列分析及びシステム

本開示は、本明細書で同定されるキメラアンプリコンアレイだけでなく、提供される方法を実施するためのコンピュータ及びシステムも包含する。

【0111】

試料を得るための一般的な方法、配列決定リードを生成するための一般的な方法、及び本開示を実施するために有用な様々なタイプの配列決定がここで記載される。これらの例示的な方法は限定的ではなく、当業者によって必要に応じて変更され得ることが理解されるべきである。

30

【0112】

複数の配列リードを得ることは、配列リードを生成するために試料から核酸を配列決定することを含み得る。複数の配列リードを得ることはまた、シーケンサから配列決定データを受け取ることを含み得る。試料中の核酸は、例えば、組織試料中のゲノム D N A、実験室試料中の特定の標的から増幅された c D N A、複数の生物由来の混合 D N A、合成核酸配列(例えば、バーコード及び固有の分子識別子(U M I))等を含む任意の核酸であり得る。一実施形態では、核酸鋳型分子(例えば、D N A 又は R N A)は、タンパク質、脂質、及び非鋳型核酸等の様々な他の成分を含有する生物学的試料から単離される。核酸鋳型分子は、動物、植物、細菌、真菌、又は任意の他の細胞生物から得られる任意の細胞材料から得ることができる。本開示で使用するための生物学的試料には、ウイルス粒子又は調製物も含まれる。核酸鋳型分子は、生物から直接得ることができ、又は生物から得られた生物学的試料、例えば血液、尿、脳脊髄液、精液、唾液、痰、糞便及び組織から得ることができる。任意の組織又は体液検体(例えば、体液標本のヒト組織)を、本開示において使用するための核酸の供給源として使用することができる。核酸鋳型分子は、初代細胞培養物又は細胞株等の培養細胞から単離することもできる。鋳型核酸が得られる細胞又は組織は、ウイルス又は他の細胞内病原体に感染し得る。試料はまた、生物学的検体、c

40

50

DNAライブラリ、ウイルス又はゲノムDNAから抽出された全RNAであり得る。試料はまた、非細胞起源から単離されたDNA、例えばフリーザからの増幅/単離されたDNAであり得る。

【0113】

一般に、核酸は、Green and Sambrook, *Molecular Cloning: A Laboratory Manual (Fourth Edition)*, Cold Spring Harbor Laboratory Press, Woodbury, N.Y., 2,028 pages (2012)に記載されているような、又は米国特許第7,957,913号、第7,776,616号、第5,234,809号、米国特許出願公開第2010/0285578号、及び米国特許出願公開第2002/0190663号に記載されているような様々な技術によって抽出、単離、増幅又は分析することができる。

10

【0114】

生物学的試料から得られた核酸を断片化して、分析に適した断片を生成することができる。鋳型核酸は、様々な機械的、化学的、及び/又は酵素的方法を使用して、所望の長さの断片化又は剪断され得る。DNAは、例えば、Covarisによって販売されている超音波処理機(Woburn, Mass.)、DNaseへの短時間の曝露、あるいは1つ又は複数の制限酵素の混合物、あるいはトランスポザゼ又はニッキング酵素を使用して、超音波処理によってランダムに剪断され得る。RNAは、RNase、熱+マグネシウムへの短時間の曝露によって、又は剪断によって断片化され得る。RNAをcDNAに変換することができる。断片化が使用される場合、RNAは、断片化の前又は後にcDNAに変換され得る。一実施形態では、核酸は超音波処理によって断片化される。別の実施形態では、核酸は、水素化剪断装置によって断片化される。一般に、個々の核酸鋳型分子は、約2kb塩基~約40kbであり得る。特定の実施形態では、核酸は約6kb~10kbの断片である。核酸分子は、一本鎖、二本鎖、又は一本鎖領域を有する二本鎖(例えば、ステム構造及びループ構造)であり得る。

20

【0115】

生物学的試料は、必要に応じて洗剤又は界面活性物質の存在下で溶解、ホモジナイズ又は分画され得る。適切な界面活性剤は、イオン性界面活性剤(例えば、ドデシル硫酸ナトリウム又はN-ラウロイルサルコシン)又は非イオン性界面活性剤(例えば、商標TWEENでUniqema Americas (Paterson, N.J.)により販売されているポリソルベート80又はTRITON X-100として知られているC₁₄H₂₂O(C₂H₄)_nを含み得る。核酸が試料から抽出又は単離されると、それは増幅され得る。

30

【0116】

増幅は、核酸配列の更なるコピーの産生を指し、一般にポリメラーゼ連鎖反応(PCR)又は当技術分野で公知の他の技術を使用して行われる。増幅反応は、PCR等の核酸分子を増幅する当技術分野で公知の任意の増幅反応であり得る。他の増幅反応には、ネステッドPCR、PCR-一本鎖コンフォメーション多型、リガーゼ連鎖反応、鎖置換増幅及び制限断片長多型、転写ベースの増幅システム、ローリングサークル増幅、及び超分岐ローリングサークル増幅、定量PCR、定量蛍光PCR(QF-PCR)、マルチプレックス蛍光PCR(MF-PCR)、リアルタイムPCR(RT-PCR)、制限断片長多型PCR(PCR-RFLP)、in situローリングサークル増幅(RCA)、ブリッジPCR、ピコチターPCR、エマルジョンPCR、転写増幅、自立配列複製、コンセンサス配列プライムPCR、任意プライムPCR、縮重オリゴヌクレオチド-プライムPCR、及び核酸ベースの配列増幅(NABSA)が含まれる。使用され得る増幅方法としては、米国特許第5,242,794、第5,494,810号、第4,988,617号、及び第6,582,938号に記載されているものが挙げられる。ある実施形態では、増幅反応は、例えば、米国特許第4,683,195及び第4,683,202号に記載されるようなPCRであり、参照により本明細書に組み込まれる。PCR、配列決定、及

40

50

び他の方法のためのプライマーは、クローニング、直接化学合成、及び当技術分野で公知の他の方法によって調製することができる。プライマーは、Eurofins MWG Operon (Huntsville, Ala.) 又はLife Technologies (Carlsbad, Calif.) 等の商業的供給元から入手することもできる。

【0117】

バーコード配列は、各配列が核酸の特定の部分に相関するように設計することができ、配列リードをそれらが由来する部分に相関させることができる。バーコード配列のセットを設計する方法は、例えば、米国特許第6,235,475号に示されており、その内容は参照によりその全体が本明細書に組み込まれる。ある実施形態では、バーコード配列は、約5ヌクレオチド～約15ヌクレオチドの範囲である。特定の実施形態では、バーコード配列は、約4ヌクレオチド～約7ヌクレオチドの範囲である。バーコード配列のセットを設計するための方法及びバーコード配列を取り付けるための他の方法は、米国特許第7,544,473号、第7,537,897号、第7,393,665号、第6,352,828号、第6,172,218号、第6,172,214号、第6,150,516号、第6,138,077号、第5,863,722号、第5,846,719号、第5,695,934、及び第5,604,097号に示されており、それぞれ参照により組み込まれる。

【0118】

配列決定は、当技術分野で公知の任意の方法によるものであり得る。DNA配列決定技術には、標識ターミネータ又はプライマー及びスラブ又はキャピラリーでのゲル分離を使用する古典的なジデオキシ配列決定反応(サンガー法)、可逆的に末端化された標識ヌクレオチドを使用する合成による配列決定、パイロシーケンシング、454配列決定、Illumina/Solexa配列決定、標識オリゴヌクレオチドプローブのライブラリへの対立遺伝子特異的ハイブリダイゼーション、標識クローンのライブラリへの対立遺伝子特異的ハイブリダイゼーションを使用した合成とそれに続くライゲーションによる配列決定、重合工程中の標識ヌクレオチドの取込みのリアルタイムモニタリング、ポロニー配列決定、及びSOLID配列決定が含まれる。分離された分子の配列決定は、最近になって、ポリメラーゼ又はリガーゼを用いた連続的又は単一の伸長反応、並びにプローブのライブラリとの単一又は連続的な差次的ハイブリダイゼーションによって実証された。

【0119】

使用され得る配列決定技術としては、例えば、Roche (Branford, Conn.) の454 Life Sciencesによって商標GS JUNIOR、GS FLX+及び454 SEQUENCINGとして販売されており、内容は、参照によりその全体が本明細書に組み込まれる、Margulies, M. et al., Genome sequencing in micro-fabricated high-density picotiter reactors, Nature, 437: 376-380 (2005)、米国特許第5,583,024号、第5,674,713号、及び第5,700,673号に記載されている合成による配列決定システムの使用が挙げられる。454配列決定は二段階を含む。これらのシステムの第1の工程では、DNAを約300～800塩基対の断片に切断し、断片を平滑末端化する。次いで、オリゴヌクレオチドアダプタを断片の末端にライゲーションする。アダプタは、断片の増幅及び配列決定のためのプライマーとして機能する。断片は、例えば5'-ビオチンタグを含むアダプタBを使用して、DNA捕捉ビーズ、例えばストレプトアビジン被覆ビーズに結合させることができる。ビーズに付着した断片は、油-水エマルジョンの液滴内でPCR増幅される。結果は、各ビーズ上のクローン増幅DNA断片の複数のコピーである。第2の工程では、ビーズをウェル(ピコリットルサイズ)に捕捉する。パイロシーケンシングは、各DNA断片に対して並行して行われる。1つ又は複数のヌクレオチドの付加は、配列決定機器においてCCDカメラによって記録される光信号を生成する。シグナル強度は、組み込まれるヌクレオチドの数に比例する。パイロシーケンシングは、ヌクレオチド付加時に放出されるピロホスファート(PPi)を利用する。PPiは、アデノシン5'ホスホスルフ

10

20

30

40

50

エートの存在下でATPスルフィラーゼによってATPに変換される。ルシフェラーゼはATPを使用してルシフェリンをオキシルシフェリンに変換し、この反応は検出及び分析される光を生成する。

【0120】

使用することができるDNA配列決定技術の別の例は、Life Technologies Corporation (Carlsbad, Calif.)のApplied BiosystemsによるSOLiD技術である。SOLiD配列決定では、ゲノムDNAを断片に切断し、断片の5'及び3'末端にアダプタを結合させて断片ライブラリを生成する。あるいは、アダプタを断片の5'末端及び3'末端にライゲーションし、断片を環状化し、環状化された断片を消化して内部アダプタを生成し、得られた断片の5'末端及び3'末端にアダプタを付着させて、メイト・ペア (mate-paired) ライブラリを生成することによって、内部アダプタを導入することができる。次に、クローンピース集団を、ビーズ、プライマー、鋳型及びPCR成分を含有するマイクロリアクタ中で調製する。PCRの後、鋳型を変性させ、ビーズを濃縮して、伸長した鋳型を有するビーズを分離する。選択されたビーズ上の鋳型は、スライドガラスへの結合を可能にする3'修飾に供される。配列は、部分的にランダムなオリゴヌクレオチドと、特定のフルオロフォアによって識別される中央の決定された塩基 (又は塩基対) との連続的なハイブリダイゼーション及びライゲーションによって決定することができる。色を記録した後、ライゲーションしたオリゴヌクレオチドを除去し、次いでプロセスを繰り返す。

【0121】

使用され得るDNA配列決定技術の別の例は、例えば、Life Technologies (South San Francisco, Calif.)によってIon TorrentによりION TORRENTの商標で販売されているシステムを使用するイオン半導体配列決定である。イオン半導体配列決定は、例えば、Rothberg, et al., An integrated semiconductor device enabling non-optical genome sequencing, Nature 475: 348-352 (2011)、米国特許公開第2010/0304982号、米国特許公開第2010/0301398号、米国特許公開第2010/0300895号、米国特許公開第2010/0300559号、及び米国特許公開第2009/0026082号に記載されており、これらの各々の内容は、参照によりその全体が

【0122】

使用され得る配列決定技術の別の例は、Illumina配列決定である。Illumina配列決定は、フォールドバック (fold-back) PCR及び固定プライマーを用いた固体表面上のDNAの増幅に基づく。ゲノムDNAを断片化し、断片の5'及び3'末端にアダプタを付加する。フローセルチャネルの表面に付着したDNA断片は伸長され、ブリッジ増幅される。断片は二本鎖になり、二本鎖分子は変性する。固相増幅とそれに続く変性の複数のサイクルは、フローセルの各チャネルに同じ鋳型の一本鎖DNA分子の約1,000コピーの数百万のクラスターを作製することができる。プライマー、DNAポリメラーゼ及び4つのフルオロフォア標識された可逆的に終結するヌクレオチドを使用して、順次配列決定 (sequential sequencing) を行う。ヌクレオチド取込み後、レーザーを使用してフルオロフォアを励起し、画像を取り込み、第1の塩基のアイデンティティを記録する。組み込まれた各塩基からの3'ターミネータ及びフルオロフォアを除去し、組込み、検出及び同定工程を繰り返す。この技術による配列決定は、米国特許第7,960,120号、第7,835,871号、第7,232,656号、第7,598,035号、第6,911,345号、第6,833,246号、第6,828,100号、第6,306,597号、第6,210,891号、米国特許公開第2011/0009278号、米国特許公開第2007/0114362号、米国特許公開第2006/0292611号、及び米国特許公開第2006/0024681号に記載されており、これらの各々は、参照によりその全体が組み込まれる。

【0123】

使用され得る配列決定技術の別の例としては、Pacific Biosciences (Menlo Park, Calif.)の単一分子リアルタイム (SMRT) 技術が挙げられる。SMRTでは、4つのDNA塩基のそれぞれが、4つの異なる蛍光色素のうちの一つに結合している。これらの色素は、リン連結されている。単一のDNAポリメラーゼは、ゼロモード導波路 (ZMW) の底部に鑄型一本鎖DNAの単一分子で固定化される。成長中の鎖にヌクレオチドを組み込むのに数ミリ秒かかる。この間、蛍光標識が励起され、蛍光シグナルが発生し、蛍光タグが切断される。色素の対応する蛍光の検出は、どの塩基が組み込まれたかを示す。このプロセスを繰り返す。

【0124】

使用され得る配列決定技術の別の例は、ナノポア配列決定である (Soni & Meller, 2007, Progress toward ultrafast DNA sequencing using solid-state nanopores, Clin Chem 53 (11): 1996 - 2001)。ナノポアは、直径1ナノメートル程度の小さな孔である。ナノポアを導電性流体に浸漬し、ナノポアの両端に電位を印加すると、ナノポアを通過するイオンの伝導に起因してわずかな電流が生じる。流れる電流の量は、ナノポアのサイズに対し感度を有する。DNA分子がナノポアを通過するとき、DNA分子上の各ヌクレオチドは、ナノポアを異なる程度で遮る。したがって、DNA分子がナノポアを通過するときナノポアを通過する電流の変化は、DNA配列の読取りを表す。

【0125】

使用され得る配列決定技術の別の例は、化学感応性電界効果トランジスタ (chemFET) アレイを使用してDNAを配列決定することを含む (例えば、米国特許出願公開第2009/0026082号)。この技術の一例では、DNA分子を反応チャンバ内に配置することができ、鑄型分子をポリメラーゼに結合した配列決定プライマーにハイブリダイズさせることができる。1つ又は複数の三リン酸の、配列決定プライマーの3'末端での新しい核酸鎖への組み込みは、chemFETによる電流の変化によって検出することができる。アレイは、複数のchemFETセンサを有することができる。別の例では、単一の核酸をビーズに付着させることができ、核酸をビーズ上で増幅させることができ、個々のビーズをchemFETアレイ上の個々の反応チャンバに移送することができ、各チャンバはchemFETセンサを有し、核酸を配列決定することができる。

【0126】

使用することができる配列決定技術の別の例は、例えばMoudrianakis, E. N. and Beer M., in Base sequence determination in nucleic acids with the electron microscope, III. Chemistry and microscopy of guanine-labeled DNA, PNAS 53: 564 - 71 (1965) によって記載されているような電子顕微鏡を使用することを含む。この技術の一例では、個々のDNA分子は、電子顕微鏡を用いて識別可能な金属標識を使用して標識される。次いで、これらの分子を平らな表面に伸ばし、電子顕微鏡を使用して画像化して配列を測定する。

【0127】

本開示の実施形態による配列決定は、複数のリードを生成する。本開示によるリードは、一般に、約150塩基長未満、又は約90塩基長未満のヌクレオチドデータの配列を含む。ある実施形態において、リードは、約80~約90塩基長、例えば、約85塩基長である。いくつかの実施形態において、本開示の方法は、非常に短いリード、すなわち、約50又は約30塩基長未満の長さに適用される。配列リードデータは、配列データ及びメタ情報を含み得る。配列リードデータは、当業者に知られているように、例えば、VCFファイル、FASTAファイル又はFASTQファイルを含む任意の適切なファイルフォーマットで保存することができる。

【0128】

10

20

30

40

50

FASTAは、元々、配列データベースを検索するためのコンピュータプログラムであり、FASTAという名称は、標準ファイルフォーマットも指すようになった。Pearson & Lipman, 1988, Improved tools for biological sequence comparison, PNAS 85:2444-2448を参照されたい。FASTAフォーマットの配列は、1行の記述で始まり、その後に配列データの行が続く。記述行は、1列目の大なり(「>」)記号によって配列データと区別される。「>」記号に続く単語は配列の識別子であり、行の残りは記述である(両方とも任意である)。「>」と識別子の最初の文字との間にスペースがあってはならない。テキストの全ての行が80文字未満であることが推奨される。「>」で始まる別の行が現れると配列は終了し、これは別の配列の開始を示す。

10

【0129】

FASTQフォーマットは、生物学的配列(通常はヌクレオチド配列)及びその対応する品質スコアの両方を保存するためのテキストベースのフォーマットである。これはFASTAフォーマットに類似しているが、配列データに続く品質スコアを有する。配列文字及び品質スコアの両方は、簡潔にするために単一のASCII文字で符号化される。FASTQフォーマットは、Illumina Genome Analyzer. Cock et al., 2009, The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants, Nucleic Acids Res 38(6):1767-1771.等のハイスループット配列決定装置の出力を保存するための事実上のスタンダードである。

20

【0130】

FASTA及びFASTQファイルの場合、メタ情報は記述行を含み、配列データの行を含まない。いくつかの実施形態では、FASTQファイルの場合、メタ情報は品質スコアを含む。FASTA及びFASTQファイルの場合、配列データは、記述行の後に始まり、典型的には、任意選択的に「-」を有するIUPAC多義性符号のいくつかのサブセットを使用して存在する。好ましい実施形態では、配列データは、必要に応じて「-」又はUを(例えば、間隙又はウラシルを表すために)任意選択的に含むA、T、C、G、及びN文字を使用する。

【0131】

30

上記及び他の場所で説明したように、NGS機器の出力量は増加している。例えば、Pinho & Pratas, 2013, MFCompress: a compression tool for FASTA and multi-FASTA data, Bioinformatics 30(1):117-8; Deorowicz & Grabowski, 2013, Data compression for sequencing data, Alg Mol Bio 8:25; Balzer et al., 2013, Filtering duplicate reads from 454 pyrosequencing data, Bioinformatics 29(7):830-836; Xu et al., 2012, FastUniq: A fast de novo duplicates removal tool for paired short reads, PLoS One 7(12):e52249; Bonfield and Mahoney, 2013, Compression of FASTQ and SAM format sequencing data, PLoS One 8(3):e59190; and Veeneman et al., 2012, Oculus: faster sequence alignment by streaming read compression, BMC Bioinformatics 13:297を参照されたい。NGS技術によって生成されるデータの量は、そのような配列決定情報を含むファイルを保存及び転送する際の困難を引き起こす。したがって、本開示の方法及びシステムは、核酸配列決定技術に由来するFASTA又はFASTQファイル(FASTA/Qファイル)に含まれる大量の配列データ等の情報を保存するために使

40

50

用することができる。

【0132】

いくつかの実施形態において、配列リードファイル及び/又は配列出力ファイルは、プレーンテキストファイル（例えば、ASCII、ISO/IEC 646、EBCDIC、UTF-8、又はUTF-16等の符号化を使用する）として保存される。本開示によって提供されるコンピュータシステムは、プレーンテキストファイルを開くことができるテキストエディタプログラムを含むことができる。テキストエディタプログラムは、コンピュータ画面上にテキストファイル（プレーンテキストファイル等）の内容を提示し、人がテキストを編集することを可能にすることができるコンピュータプログラムを指すことができる（例えば、モニタ、キーボード、及びマウスを使用する）。例示的なテキストエディタには、Microsoft Word、emacs、pico、vi、BBEdit、及びTextWranglerが含まれるが、これらに限定されない。好ましくは、テキストエディタプログラムは、コンピュータ画面上にプレーンテキストファイルを表示することができる、メタ情報及び配列リードを人が読める形式（例えば、バイナリエンコードされていない）で示すことができる。

10

【0133】

いくつかの実施形態では、本開示の工程のいずれか又は全ては自動化される。例えば、Perlスクリプト又はシェルスクリプトを記述して、上述の様々なプログラムのいずれか呼び出すことができる（例えば、Tisdall, Mastering Perl for Bioinformatics, O'Reilly & Associates, Inc., Sebastopol, CA 2003; Michael, R., Mastering Unix Shell Scripting, Wiley Publishing, Inc., Indianapolis, Ind. 2003を参照されたい）。あるいは、本開示の方法は、1つ又は複数の専用プログラムで全体的又は部分的に具体化されてもよく、例えば、それぞれ任意選択的にC++等のコンパイル型言語で記述され、次いでコンパイルされ、バイナリとして配布される。本開示の方法は、既存の配列分析プラットフォーム内のモジュールとして、又は既存の配列分析プラットフォーム内の機能呼び出すことによって、全体的又は部分的に実施され得る。ある実施形態では、本開示の方法は、単一の開始キュー（例えば、人の活動、別のコンピュータプログラム、又は機械から供給されるトリガーイベントの1つ又は組み合わせ）に反応して全て自動的に呼び出される多数の工程を含む。したがって、本開示は、任意の工程又は工程の任意の組み合わせがキューに反応して自動的に行われ得る方法を提供する。人の入力、影響、又は相互作用を介在させることのない自動的な一般的手段である（すなわち、元の又はプレキューの人の活動にのみ応答性である）。

20

30

【0134】

本開示はまた、対象核酸の正確かつ高感度な解釈を含む様々な形態の出力を包含する。出力は、コンピュータファイルの形式で提供することができる。ある実施形態では、出力は、FASTAファイル、FASTQファイル、又はVCFファイルである。出力を処理して、テキストファイル、又は参照ゲノムの配列にアライメントされた核酸の配列等の配列データを含むXMLファイルを生成することができる。他の実施形態では、処理は、参照ゲノムに対する対象核酸中の1つ又は複数の突然変異を記述する座標又は文字列を含む出力をもたらす。当技術分野で公知のアライメントストリングとしては、Simple UnGapped Alignment Report (SUGAR)、Verbose Useful Labeled Gapped Alignment Report (VULGAR)、及びCompact Idiosyncratic Gapped Alignment Report (CIGAR) (Ning, Z., et al., Genome Research 11(10): 1725-9 (2001))が挙げられる。これらの文字列は、例えば、European Bioinformatics Institute (Hinxton, UK)のExonerate配列アライメントソフトウェアに実装されている。

40

50

【0135】

いくつかの実施形態では、配列アライメントは、例えば配列アライメントマップ (SAM) 又はバイナリアライメントマップ (BAM) ファイル等、CIGAR文字列 (SAM形式は、例えば、Li et al., The Sequence Alignment / Map format and SAMtools, Bioinformatics, 2009, 25(16): 2078-9に記載されている) を含むものとして作成される。いくつかの実施形態では、CIGARは、ギャップのあるアライメントをラインごとに表示又は含む。CIGARは、CIGAR文字列として報告される圧縮されたペアワイズアライメントフォーマットである。CIGAR文字列は、長い (例えば、ゲノム) ペアワイズアライメントを表すのに有用である。CIGAR文字列は、参照ゲノム配列に対するリードのアライメントを表すためにSAM形式で使用される。

10

【0136】

CIGAR文字列は、確立されたモチーフに続く。各文字の前に数字が付けられ、イベントの塩基カウントが与えられる。使用される文字は、M、I、D、N、及びS (M = マッチ; I = 挿入; D = 欠失; N = ギャップ; S = 置換) を含むことができる。CIGAR文字列は、マッチ/ミスマッチ及び欠失 (又はギャップ) の配列を定義する。例えば、CIGAR文字列2MD3M2D2Mは、アライメントが2つのマッチ、1つの欠失 (いくらかのスペースを節約するために番号1は省略されている)、3つのマッチ、2つの欠失及び2つのマッチを含むことを意味する。

【0137】

本開示によって企図されるように、上述の機能は、ソフトウェア、ハードウェア、ファームウェア、ハード配線、又はこれらの任意の組み合わせを含む本開示のシステムを使用して実施することができる。機能を実装する特徴はまた、機能の一部が異なる物理的位置に実装されるように分散されることを含む、様々な位置に物理的に配置され得る。

20

【0138】

当業者であれば、本開示の方法の実行に必要な又は最も適していると認識するように、本開示のコンピュータシステム又はマシンは、バスを介して互いに通信する1つ又は複数のプロセッサ (例えば、中央プロセッシングユニット (CPU)、グラフィックスプロセッシングユニット (GPU)、又はその両方)、メインメモリ、及びスタティックメモリを含む。

30

【0139】

図12は、本開示の方法を実行するのに適したシステム701を示す。図12に示されるように、システム701は、サーバコンピュータ705、端末715、シーケンサ715、シーケンサコンピュータ721、コンピュータ749、又はそれらの任意の組み合わせのうちの1つ又は複数を含み得る。そのようなコンピュータデバイスの各々は、ネットワーク709を介して通信することができる。シーケンサ725は、任意選択的に、それ自体の、例えば専用のシーケンサコンピュータ721 (任意の入力/出力機構 (I/O)、プロセッサ、及び、例えばダイナミックランダムアクセスメモリDRAM又はDRAM729等のメモリを含む) を含むか、又はそれに動作可能に結合されてもよい。追加的又は代替的に、シーケンサ725は、ネットワーク709を介してサーバ705又はコンピュータ749 (例えば、ラップトップ、デスクトップ、又はタブレット) に動作可能に結合されてもよい。コンピュータ749は、1つ又は複数のプロセッサ、メモリ、及びI/Oを含む。本開示の方法がクライアント/サーバアーキテクチャを使用する場合、本開示の方法の任意の工程は、データ、命令等を取得するか、又はインターフェースモジュールを介して結果を提供するか、又はファイルとして結果を提供することができる、プロセッサ、メモリ、及びI/Oのうちの1つ又は複数を含むサーバ705を使用して実行され得る。サーバ705は、コンピュータ749又は端末715によりネットワーク709を介して係合されてもよく、又はサーバ705は、端末715に直接接続されてもよい。端末715は、好ましくはコンピュータデバイスである。本開示によるコンピュータは、好ましくは、I/O機構及びメモリに結合された1つ又は複数のプロセッサを含む。

40

50

【0140】

プロセッサは、例えば、シングルコア又はマルチコアプロセッサ（例えば、AMD Phenom II X2、Intel Core Duo、AMD Phenom II X4、Intel Core i5、Intel Core i&Extreme Edition 980X、又はIntel Xeon E7-2820）のうちの1つ又は複数を含む1つ又は複数のプロセッサによって提供され得る。

【0141】

I/O機構は、ビデオ表示ユニット（例えば、液晶ディスプレイ（LCD）又は陰極線管（CRT））、英数字入力デバイス（例えば、キーボード）、カーソル制御デバイス（例えば、マウス）、ディスク駆動ユニット、信号生成デバイス（例えば、スピーカ）、加 10
速度計、マイクロフォン、セルラー無線周波数アンテナ、及びネットワークインターフェースデバイス（例えば、ネットワークインターフェースカード（NIC）、Wi-Fiカード、セルラーモデム、データジャック、イーサネットポート、モデムジャック、HDMI（登録商標）ポート、ミニHDMI（登録商標）ポート、USBポート）、タッチスクリーン（例えば、CRT、LCD、LED、AMOLED、Super AMOLED）、ポインティングデバイス、トラックパッド、ライト（例えば、LED）、光/画像投影デバイス、又はそれらの組合わせを含むことができる。

【0142】

本開示によるメモリは、1つ又は複数の有形デバイスによって提供される非一時的メモリを指し、有形デバイスは、本明細書に記載の方法又は機能のいずれか1つ又は複数を具 20
現化する1つ又は複数の命令セット（例えば、ソフトウェア）が格納された1つ又は複数の機械可読媒体を含むことが好ましい。ソフトウェアはまた、システム501内のコンピュータによる実行中に、メインメモリ、プロセッサ、又はその両方内に完全に又は少なくとも部分的に存在してもよく、メインメモリ及びプロセッサはまた、機械可読媒体を構成する。ソフトウェアは、ネットワークインターフェース装置を介してネットワークにわたって更に送信又は受信することができる。

【0143】

機械可読媒体は、例示的な実施形態では単一の媒体であり得るが、「機械可読媒体」という用語は、1つ又は複数の命令セットを格納する単一の媒体又は複数の媒体（例えば、 30
集中型又は分散型データベース、並びに/あるいは関連するキャッシュ及びサーバ）を含むと解釈されるべきである。「機械可読媒体」という用語はまた、機械によって実行するための命令のセットを格納、エンコード、又は搬送することができ、機械に本開示の方法論のうちの任意の1つ又は複数を実行させる任意の媒体を含むと解釈されるべきである。メモリは、例えば、ハードディスクドライブ、ソリッドステートドライブ（SSD）、光ディスク、フラッシュメモリ、ジップディスク、テープドライブ、「クラウド」保存場所、又はそれらの組合わせのうちの1つ又は複数であってもよい。ある実施形態では、本開示の装置は、メモリ用の有形の非一時的コンピュータ可読媒体を含む。メモリとして使用するための例示的なデバイスには、半導体メモリデバイス（例えば、EPROM、EEPROM、ソリッドステートドライブ（SSD）、及びフラッシュメモリデバイス、例えば、SD、マイクロSD、SDXC、SDIO、SDHCカード）が含まれ、磁気ディスク 40
（例えば、内蔵ハードディスク又はリムーバブルディスク）、及び光ディスク（例えば、CD及びDVDディスク）が含まれる。

【0144】

コンティグを構築し、コンセンサス配列を生成する様々な方法を以下に論じる。

【0145】

コンティグは、一般に、核酸配列、例えばリードの複数のセグメント間、又はその中の 50
関係を指す。配列リードが重複する場合、コンティグは、重複リードの階層化画像として表すことができる。コンティグは、例えば、テキストファイル又はデータベース内の任意の特定の視覚的配置又は任意の特定の配置によって定義されず、それらに限定されない。コンティグは、一般に、配列決定された核酸の一部に対応するように編成された多数のR

ードからの配列データを含む。コンティグは、表示又は保存された、リードのセット又は互いに対する若しくは参照に対するそれらの位置に関する情報等の組み立て結果を含むことができる。コンティグは、行が個々の配列リードであり、列がその部位にアライメントすると推定される各リードの塩基を含むグリッドとして構成することができる。コンセンサス配列は、アセンブリの各カラム中の優勢な塩基を同定することによって作製することができる。本発明によるコンティグは、それらが互いに重なり合う（又は、重複せず、例えば、単に隣接する）ことを示すリードの視覚的表示を含むことができる。コンティグは、複数のリードに関連付けられ、互いに対するリードの位置を与える座標のセットを含み得る。コンティグは、リードの配列データを変換することによって得られたデータを含み得る。例えば、Burrows - Wheeler変換をリードに対して行うことができ、コンティグは、リードの非変換配列を必ずしも含まずに変換データを含むことができる。ヌクレオチド配列データのBurrows - Wheeler変換は、米国特許出願公開第2005/0032095号に記載され、その全体が参照により本明細書に組み込まれる。

10

【0146】

リードは、当技術分野で公知の任意の方法によってコンティグに組み立てることができる。複数の配列リードのデノボアセンブリのためのアルゴリズムは当技術分野において公知であるが、そのような公知のアルゴリズムは、本開示で記載されている構造化した配列リード入力のために本明細書において改良されている（個々の配列要素は、長い配列リードのより広い集団の各長配列リード内の反復シリーズ（キメラアレイ）として存在する、低複雑度のリンカー配列に隣接する、高複雑度のライブラリに由来していた）。

20

【0147】

配列リードを組み立てるための1つのアルゴリズムは、オーバーラップコンセンサスアセンブリとして知られている。オーバーラップコンセンサスアセンブリは、配列リード間のオーバーラップを使用してそれらの間のリンクを作成する。リードは、一般に、非ランダムな重複が想定されるのに十分に重複する領域によって連結される。このようにリードを一緒に連結すると、コンティグ又はオーバーラップグラフが生成され、各ノードはリードに対応し、エッジは2つのリード間のオーバーラップを表す。オーバーラップグラフによるアセンブリは、例えば、米国特許第6,714,874号に記載されている。

【0148】

いくつかの実施形態では、デノボアセンブリは、いわゆるグリーディアルゴリズムに従って進行する。グリーディアルゴリズムに従って組み立てるために、リードの一群のリードのうちの1つが選択され、それは、それがかなりの量の重複を示す別のリードと対にされ、一般に、それは他の全てのリードのうちの最も多くの重複を示すリードと対にされる。これらの2つのリードはマージされて新しいリード配列を形成し、次いでそのリード群に戻され、プロセスが繰り返される。グリーディアルゴリズムによるアセンブリは、例えば、Schatz, et al., Genome Res., 20:1165-1173 (2010)及び米国特許出願公開第2011/0257889号に記載され、これらの各々は、参照によりその全体が本明細書に組み込まれる。

30

【0149】

他の実施形態では、アセンブリは、ペアワイズアライメント、例えば網羅的又はヒューリスティック（例えば、網羅的ではない）ペアワイズアライメントによって進行する。アライメントについては、一般に、以下でより詳細に説明する。「力づく（brute force）」アプローチと呼ばれることもある網羅的なペアワイズアライメントは、セット内の配列の可能な全ての対の間の可能な全てのアライメントについてアライメントスコアを計算する。ヒューリスティック多重配列アライメントによるアセンブリは、数学的にありそうにない特定の組み合わせを無視し、計算的に高速であり得る。マルチプル配列アライメントによる組立ての1つのヒューリスティックな方法は、いわゆる「分割統治」ヒューリスティックであり、これは例えば、米国特許出願公開第2003/0224384号に記載される。マルチプル配列アライメントによる組立ての別のヒューリスティック方法

40

50

は、プログラム ClustalW によって実施されるプログレッシブアライメントである（例えば、Thompson, et al., Nucl. Acids. Res., 22: 4673-80 (1994) を参照されたい）。多重配列アライメントによるアセンブリは、一般に、Lecompte, O., et al., Gene 270: 17-30 (2001); Mullan, L. J., Brief Bioinform., 3: 303-5 (2002); Nicholas, H. B. Jr., et al., Biotechniques 32: 572-91 (2002); and Xiong, G., Essential Bioinformatics, 2006, Cambridge University Press, New York, N. Y. に記載されている。

【0150】

アライメントによる組立ては、リードを互いにアライメントすることによって、又はリードを参照にアライメントすることによって進行することができる。例えば、各リードを参照ゲノムに順にアライメントすることによって、全てのリードを互いに関連して配置してアセンブリを作製する。

【0151】

リードをコンティグに組み立てる1つの方法は、de Bruijn グラフを作成することを含む。de Bruijn グラフは、リードを k -mer と呼ばれる DNA のより小さな配列に分割することによって計算労力を削減し、パラメータ k はこれらの配列の塩基長を表す。de Bruijn グラフでは、全てのリードが k -mer (リード内の長さ k の全ての部分配列) に分割され、 k -mer 間の経路が計算される。この方法によるアセンブリでは、リードは、 k -mer を通る経路として表される。de Bruijn グラフは、これらの k -mer 間で長さ $k-1$ の重複を捕捉し、実際のリード間では捕捉しない。したがって、例えば、配列決定 CATGGA は、以下の 2-mer: CA, AT, TG, GG、及び GA を通る経路として表すことができる。de Bruijn グラフ手法は、冗長性をうまく扱い、複雑な経路の計算を扱いやすくする。データセット全体を k -mer 重複まで削減することにより、de Bruijn グラフは、ショートリードデータセットの高い冗長性を削減する。特定のアセンブリの最大効率的な k -mer サイズは、リード長並びにエラーレートによって決定される。パラメータ k の値は、アセンブリの品質に大きな影響を及ぼす。良好な値の推定は、組み立て前に行うことができ、又は、最適な値は、小さな範囲の値を試験することによって見つけることができる。de Bruijn グラフを使用したリードのアセンブリは、米国特許出願公開第 2011/0004413 号、米国特許出願公開第 2011/0015863 号、及び米国特許出願公開第 2010/0063742 号に記載され、これらの各々は、参照によりその全体が本明細書に組み込まれる。

【0152】

本発明による、リードをコンティグに組み立てる他の方法も可能である。例えば、リードは、配列決定中に鋳型核酸に挿入されたバーコード情報を含み得る。ある実施形態において、リードは、バーコード情報を参照することによってコンティグにアセンブルされる。例えば、バーコードを識別することができ、バーコードと一緒に配置することによってリードを組み立てることができる。

【0153】

リードのコンティグへのアセンブリは、Husemann, P. and Stoye, J., Phylogenetic Comparative Assembly, 2009, Algorithms in Bioinformatics: 9th International Workshop, pp. 145-156, Salzberg, S., 及び Warnow, T., Eds. Springer-Verlag, Berlin Heidelberg で更に論じられている。リードをコンティグに組み立てるためのいくつかの例示的な方法は、例えば、国特許出願公開第 6,223,128 号、国特許出願公開第 2009/0298064 号、米国特許出願公開 2010/0069263 号、及び米国特許出願公開第 2011/0257889 号に記載され、これらの各々は、参照によ

10

20

30

40

50

りその全体が本明細書に組み込まれる。

【0154】

リードを組み立てるためのコンピュータプログラムは、当技術分野において公知である。そのようなアセンブリプログラムは、単一の汎用コンピュータ上で、コンピュータのクラスター若しくはネットワーク上で、又は配列分析専用の専用コンピューティングデバイス上で実行することができる。

【0155】

アセンブリは、例えば、カナダのMichael Smith Genome Sciences Centre (Vancouver, B.C., CA)からのプログラム「The Short Sequence Assembly by k-mer search and 3' read Extension」(SSAKE)によって実施することができる(例えば、Warren, R., et al., Bioinformatics, 23:500-501(2007)を参照されたい)。SSAKEは、リードのテーブルを循環し、任意の2つの配列間の可能な限り長い重複についてプレフィックスツリーを検索する。SSAKEクラスターはコンティグに読み取る。

10

【0156】

別のリードアセンブリプログラムは、Darren Platt及びDirk Eversによって書かれたForge Genome Assemblerであり、Geeknet (Fairfax, Va.)によって管理されているSourceForgeウェブサイトを通じて入手可能である(例えば、DiGiustini, S., et al., Genome Biology, 10:R94(2009)を参照されたい)。Forgeは、利用可能であれば、その計算及びメモリ消費を複数のノードに分配し、したがって、大きなリードセットを組み立てる可能性を有する。Forgeは、並列MPIライブラリを使用してC++で書かれた。Forgeは、リードの混合物、例えば、Sanger、454及びIlluminaによるリードを扱うことができる。

20

【0157】

多重配列アライメントによるアセンブリは、例えば、University College Dublin (Dublin, Ireland)から入手可能なプログラムClustal Omega (Sievers F., et al., Mol Syst Biol 7(2011)), ClustalW, or ClustalX (Larkin M.A., et al., Bioinformatics, 23, 2947-2948(2007))によって行うことができる。

30

【0158】

当技術分野で知られている別の例示的なリードアセンブリプログラムは、European Bioinformatics Institute (Hinxton, UK)のウェブサイトを通じて入手可能なVelvetである(Zerbino D.R. et al., Genome Research 18(5):821-829(2008))。Velvetは、de Bruijnグラフに基づく手法を実装し、リード対からの情報を使用し、様々なエラー補正工程を実装する。

【0159】

リードアセンブリは、Beijing Genomics Institute (Beijing, CN)又はBGI Americas Corporation (Cambridge, Mass.)のウェブサイトを通じて入手可能なpackage SOAPからのプログラムを用いて実行することができる。例えば、SOAPdenovoプログラムは、de Bruijnグラフ手法を実装する。SOAPS/GPUは、短いリードを参照配列にアライメントする。

40

【0160】

別のリードアセンブリプログラムは、カナダのMichael Smith Genome Sciences Centre (Vancouver, B.C., CA) (Simpson, J.T., et al., Genome Res., 19(6):1117

50

- 23 (2009)) である。ABySSは、de Bruijnグラフ手法を使用し、並列環境で実行される。

【0161】

リードアセンブリは、Roche 454シーケンサ（例えば、Kumar, S. et al., Genomics 11: 571 (2010) 及び Margulies, et al., Nature 437: 376 - 380 (2005) に記載されている）からのリードをアセンブルするように設計された、gsAssembler又はNewbler (NEW assembler) として知られるRocheのGS De Novo Assemblerによって行うこともできる。Newblerは、454のFlex Standardリード及び454のTitaniumリード、並びに単一及びペアエンドリード、並びに任意選択でSangerのリードを受け付ける。Newblerは、32ビット又は64ビットのいずれかのバージョンでLinux（登録商標）上で実行される。Newblerは、コマンドライン又はJavaベースのGUIインターフェースを介してアクセスすることができる。

【0162】

オクスフォード大学のMario Caccamo及びZamin Iqbalによって作成されたCortexは、リードアセンブリを含むゲノム解析のためのソフトウェアフレームワークである。Cortexは、Spanu, P. D., et al., Science 330 (6010): 1543 - 46 (2010) に記載されているように使用される、コンセンサスゲノムアセンブリのためのcortex__conを含む。Cortexは、Iqbal, et al., De novo assembly and genotyping of variants using colored de Bruijn graphs, Nature Genetics (in press) に記載されており、Mills, R. E., et al., Nature 470: 59 - 65 (2010) に記載されているように使用される、変異及び集団アセンブリのためのcortex__varを含む。Cortexは、制作者のウェブサイトを通じて、及びGeeknet (Fairfax, Va.) が管理するSourceForgeウェブサイトから入手可能である。

【0163】

他のリードアセンブリプログラムには、Real Time Genomics, Inc. (San Francisco, Calif.) からのRTG Investigator; iAssembler (Zheng, et al., BMC Bioinformatics 12: 453 (2011)); TgiCL Assembler (Perte, et al., Bioinformatics 19 (5): 651 - 52 (2003)); Geeknet (Fairfax, Va.) が管理するSourceForgeのウェブサイトを通してダウンロード可能な、Heng LiによるMaq (Mapping and Assembly with Qualities); MIRA3 (Mimicking Intelligent Read Assembly)、Chevreux, B., et al., Genome Sequence Assembly Using Trace Signals and Additional Sequence Information, 1999, Computer Science and Biology: Proceedings of the German Conference on Bioinformatics (GCB) 99: 45 - 56; PG A4genomics (Zhao F., et al., Genomics .94 (4): 284 - 6 (2009) に記載; 及びPhrap (例えば、de la Bastide, M. and McCombie, W. R., Current Protocols in Bioinformatics, 17: 11.4.1 - 11.4.15 (2007) に記載) が含まれる。CLC cellは、CLC bio Germany (Muehlthal, Germany) から入手可能な、NGSリードのリードマッピング及びデノボアセンブリのためのde Bruijnグラフベースのコンピュータプログラムであ

る。

【0164】

リードのアセンブリは、1つ又は複数のコンティグを生成する。ホモ接合又は単一標的配列決定の場合、単一コンティグが生成される。ヘテロ接合性の二倍体標的、稀な体細胞変異又は混合試料の場合、例えば、2つ以上コンティグが生成され得る。各コンティグは、そのコンティグを構成するリードからの情報を含む。

【0165】

リードをコンティグに組み立てることは、各コンティグに対応するコンセンサス配列を生成するのに役立つ。ある実施形態において、コンセンサス配列は、アセンブルされたリードの中からの各位置における最も一般的な又は優勢なヌクレオチドを指す。コンセンサス配列は、そのコンティグによって表される核酸の配列の解釈を表すことができる。

10

【0166】

本明細書で使用されるアライメントは、一般に、1つの配列を別の配列に沿って配置すること、各配列に沿って反復的にギャップを導入すること、その2つの配列がどの程度よく一致するかをスコアリングすること、及び、好ましくは参照に沿った様々な位置について繰り返すこと、を含む。最良のスコアリング一致は、アライメントであると見なされ、配列間の歴史的関係性に関する推論を表す。アライメントにおいて、参照中の一致しない塩基と並んだリード中の塩基は、その時点で置換突然変異が起こったことを示す。同様に、一方の配列が他方の配列中の塩基と並んでギャップを含む場合、挿入又は欠失突然変異（「インデル」）が生じたと推測される。1つの配列が互いにアライメントされていることを明示することが望まれる場合、アライメントはペアワイズアライメントと呼ばれることがある。多重配列アライメントは、一般に、例えば、一連のペアワイズアライメントを含む、2つ以上の配列のアライメントを指す。

20

【0167】

いくつかの実施形態では、アライメントのスコアリングは、置換及びインデルの確率の値を設定することを含む。個々の塩基がアライメントされる場合、マッチ又はミスマッチは、置換確率によってアライメントスコアに寄与し、これは、例えば、マッチについては1、ミスマッチについては0.33であり得る。インデルは、例えば、-1とすることができるギャップペナルティによってアライメントスコアから推定する。ギャップペナルティ及び置換確率は、配列がどのように変異するかについての経験的知識又は先験的仮定に基づくことができる。それらの値は、結果として生じるアライメントに影響を及ぼす。特に、ギャップペナルティと置換確率との間の関係は、得られるアライメントにおいて置換又はインデルが優先されるかどうかに影響を及ぼす。

30

【0168】

形式的に言えば、アライメントは、2つの配列 x と y との間の推測される関係性を表す。例えば、いくつかの実施形態では、配列 x 及び y のアライメント A は、 $(i) |x'| = |y'|$; $(ii) x'$ 及び y' からスペースを除去すると、それぞれ x 及び y に戻るはずであり、及び (iii) 任意の i について、 $x'[i]$ 及び $y'[i]$ は両方のスペースであることはできないようにスペースを含有し得る、別の2つの文字列 x' 及び y' にそれぞれ x 及び y をマップする。

40

【0169】

ギャップは、 x' 又は y' のいずれかにおける連続スペースの最大部分文字列である。アライメント A は、以下の3種類の領域を含む： (i) 一致した対（例えば、 $x'[i] = y'[i]$ ）； (ii) ミスマッチ対、（例えば、 $x'[i] \neq y'[i]$ であり、両方ともスペースではない）；又は (iii) ギャップ（例えば、 $x'[i..j]$ 又は $y'[i..j]$ のいずれかはギャップである）を含むことができる。ある実施形態では、一致した対のみが高い陽性スコア a を有する。いくつかの実施形態では、ミスマッチ対は一般に負のスコア b を有し、長さ r のギャップも負のスコア $g + rs$ を有し、ここで $g, s < 0$ である。DNA の場合、1つの一般的なスコアリングスキーム（例えば、BLAST によって使用される）は、スコア $a = 1$ 、スコア $b = -3$ 、 $g = -5$ 及び $s = -2$ とする。アラ

50

イメント A のスコアは、全ての一致した対、不一致の対及びギャップのスコアの合計である。x 及び y のアライメントスコアは、x 及び y の全ての可能なアライメントの中の最大スコアとして定義することができる。

【0170】

いくつかの実施形態では、任意の対は、置換確率の 4×4 マトリックス B によって定義されるスコア a を有する。例えば、 $B(i, i) = 1$ であり、 $0 < B(i, j) \text{ } i < j < 1$ が、1つの可能なスコアリングシステムである。例えば、転位 (transit ion) が転換 (transvers ion) よりも生物学的に可能性が高いと考えられる場合、マトリックス B は、 $B(C, T) = 0.7$ 及び $B(A, T) = 0.3$ 、又は当技術分野で公知の方法によって所望又は決定される任意の他の値のセットを含み得る。

10

【0171】

本発明のいくつかの実施形態によるアライメントは、ペアワイズアライメントを含む。ペアワイズアライメントは、一般に、m 個の文字及び n 個の文字の参照ゲノム T (標的) を有する配列 Q (クエリ) について、Q と T との間の可能な局所アライメントを発見及び評価することを含む。h i 及び k j である場合の、任意の $1 \leq i \leq n$ 及び $1 \leq j \leq m$ について、 $T[h \dots i]$ 及び $Q[k \dots j]$ の可能な最大アライメントスコアが計算される (すなわち、位置 i で終了する T の任意の部分文字列及び位置 j で終了する Q の任意の部分文字列の最良のアライメントスコア)。これは、c m 文字を有する全ての部分文字列を検査することを含むことができ、c は類似性モデルに応じた定数であり、各部分文字列を Q と別々に整列させる。各アライメントはスコア付けされ、好ましいスコアとのアライメントが、アライメントとして受け入れられる。いくつかの実施形態では、網羅的なペアワイズアライメントが実施され、これは一般に、Q と T との間の全ての可能な局所アライメント (任意選択的にいくつかの制限基準を受ける) がスコア付けされる上記のペアワイズアライメントを含む。

20

【0172】

いくつかの実施形態では、ペアワイズアライメントは、ドットマトリックス法、動的プログラミング法、又はワード法に従って進行する。動的プログラミング方法は、一般に、Smith-Waterman (SW) アルゴリズム又は Needleman-Wunsch (NW) アルゴリズムを実装する。NW アルゴリズムによるアライメントは、一般に、線形ギャップペナルティ d で類似度行列 $S(a, b)$ (例えば、前述のマトリックス B 等) に従ってアライメントされた文字をスコアリングする。行列 $S(a, b)$ は一般に置換確率を供給する。SW アルゴリズムは NW アルゴリズムと同様であるが、負のスコア行列セルはどれも 0 に設定される。SW 及び NW アルゴリズム、並びにそれらの実装形態は、米国特許出願公開第 5,701,256 号及び米国特許出願公開第 2009/0119313 号に更に詳細に記載され、両方ともその全体が参照により本明細書に組み込まれる。これらの方法を実施するための当技術分野で知られているコンピュータプログラムを以下により詳細に説明する。

30

【0173】

本発明によるアライメントは、当技術分野で公知の任意の適切なコンピュータプログラムを使用して実行することができる。

40

【0174】

BWT 手法を実装する 1 つの例示的なアライメントプログラムは、Geeknet (Fairfax, Va.) によって管理される SourceForge ウェブサイトから入手可能な Burrows-Wheeler Aligner (BWA) である。BWA は、リード、コンティグ又はコンセンサス配列を参照に対しアライメントすることができる。BWT は、ヌクレオチドあたり 2 ビットのメモリを占有し、典型的なデスクトップ又はラップトップコンピュータで 4 G 塩基対の長さのヌクレオチド配列をインデックスすることを可能にする。前処理は、BWT の構築 (すなわち、参照にインデックスを付ける) 及びサポート補助データ構造を含む。

【0175】

50

BWAは、両方ともBWTに基づいて、2つの異なるアルゴリズムを実装する。BWAによるアライメントは、約200bpまでの短いクエリに対して低いエラー率(<3%)で設計されたアルゴリズム**bwa-short**を使用して進行することができる(Li H. 及びDurbin R. *Bioinformatics*, 25:1754-60(2009))。第2のアルゴリズムであるBWA-SWは、より多くのエラーを伴うロングリード用に設計されている(Li H. 及びDurbin R. (2010) *Fast and accurate long-read alignment with Burrows-Wheeler Transform. Bioinformatics*, Epub.)。BWA-SWコンポーネントは、ヒューリスティックなSmith-Waterman様アライメントを実行して、高スコアの局所ヒットを見つける。当業者は、**bwa-sw**が「**bwa-long**」、「**bwa-long**アルゴリズム」等と呼ばれることがあることを認識するであろう。そのような使用は、一般にBWA-SWを指す。

10

【0176】

Smith-Watermanアルゴリズムのバージョンを実装するアライメントプログラムはMUMmerであり、Geeknet(Fairfax, Va.)が管理するSourceForgeウェブサイトから入手可能である。MUMmerは、完全形態又はドラフト形態にかかわらず、全ゲノムを迅速にアライメントするためのシステムである(Kurtz, S., et al., *Genome Biology*, 5:R12(2004); Delcher, A.L., et al., *Nucl. Acids Res.*, 27:11(1999))。例えば、MUMmer 3.0は、2.4GHzのLinux(登録商標)デスクトップコンピュータ上で、78MBのメモリを使用して、13.7秒で一对の5メガベースのゲノム間の20塩基対又はそれより長い完全一致を全て見つけることができる。MUMmerはまた、不完全なゲノムをアライメントすることができ、それは、ショットガン配列決定プロジェクトからの100s又は1000sのコンティグを容易に取り扱うことができ、システムに含まれるNUCmerプログラムを使用してそれらを別のセットのコンティグ又はゲノムにアライメントする。種が、類似性を検出するためのDNA配列アライメントにはあまりにも多様である場合、PROmerプログラムは、両方の入力配列の6フレーム翻訳に基づいてアライメントを生成することができる。

20

【0177】

本発明の実施形態による別の例示的なアライメントプログラムは、Kent Informatics(Santa Cruz, Calif.)からのBLATである(Kent, W.J., *Genome Research* 4:656-664(2002))。BLAT(BLASTではない)は、RAM等のメモリに参照ゲノムのインデックスを保持する。インデックスは、全ての重複しないk-mer(リピートに大きく関与するものを任意選択的に除く)を含み、デフォルトでk=11である。ゲノム自体はメモリに保持されない。インデックスは、可能性のある相同性の領域を見つけるために使用され、その領域は、その後、詳細なアライメントのためにメモリにロードされる。

30

【0178】

別のアライメントプログラムは、Beijing Genomics Institute(Beijing, CN)又はBGI Americas Corporation(Cambridge, Mass.)のSOAP2である。SOAP2は、双方向BWTを実装する(Liet al., *Bioinformatics* 25(15):1966-67(2009); Li, et al., *Bioinformatics* 24(5):713-14(2008))。

40

【0179】

配列を整列させるための別のプログラムは、Bowtieである(Langmead, et al., *Genome Biology*, 10:R25(2009))。Bowtieインデックスは、BWTを作製することによってゲノムを参照する。

【0180】

他の例示的なアライメントプログラムには、以下が挙げられる: Efficient

50

Large - Scale Alignment of Nucleotide Data bases (ELAND)又はConsensus Assessment of Sequence and Variation(CASAVA)ソフトウェアのELAND v2コンポーネント(Illumina, San Diego, Calif.); Real Time Genomics, Inc. (San Francisco, Calif.)からのRTG Investigator; Novocraft (Selangor, Malaysia)からのNovoalign; Exonerate、European Bioinformatics Institute (Hinxton, UK) (Slater, G., and Birney, E., BMC Bioinformatics 6:31(2005)); ユニバーシティ・カレッジ・ダブリン (Dublin, Ireland)からのClustal Omega (Sievers F., et al., Mol Syst Biol 7, article 539(2011)); ユニバーシティ・カレッジ・ダブリン (Dublin, Ireland)からのClustal W又はClustal X (Larkin M.A., et al., Bioinformatics, 23, 2947-2948(2007)); 及び、FASTA, European Bioinformatics Institute (Hinxton, UK) (Pearson W.R., et al., PNAS 85(8):2444-8(1988)); Lipman, D.J., Science 227(4693):1435-41(1985)。

【0181】

図13は、本開示の1つ又は複数の実施形態による最大状態経路を決定するための例示的な手順を示し、例示する。例えば、非汎用の、具体的に構成されたデバイス(例えば、システム701)は、格納された命令を実行することによって手順1200を実行することができる。手順1200は、工程1205で開始し、工程1210に進行し得、ここで、上で詳細に記載されるように、プロセスは、配列要素の線状アレイを有する個々の核酸配列リードを含む複数の核酸配列リードを得てもよい。実施形態において、高複雑度のライブラリから引き出された各核酸配列要素は、低複雑度の1つ若しくは複数の予想される核酸配列の、又は低複雑度の1つ若しくは複数の予想される核酸配列及び配列リード終端のいずれかに隣接し得る。

【0182】

工程1215において、プロセスは、高複雑度のライブラリ及び低複雑度のライブラリから引き出された個々の核酸配列要素の領域を予測するために、複数の核酸配列リードに1つ又は複数の統計的アノテーションモデルを適用し得る。実施形態において、1つ又は複数の統計的アノテーションモデルは、i)核酸配列リード全体に散在する1つ又は複数の予想される核酸配列を認識するための生成統計的アライメントモデル、又はii)既知ではない配列又は高複雑度の配列の辞書から引き出された配列を認識するためのランダム統計アライメントモデルを含み得る。実施形態では、予測された転位部位は、各モデルの末端に配置され、生成統計的アライメントモデルの内部位置内では許容されない。

【0183】

工程1220において、前の2つの工程を複数の核酸配列リードに対して繰り返すことができる。次いで、工程1225において、プロセスは、最大対数尤度値を有するモデルを識別することによって選択された最大事後状態経路の最終的リード当たりのモデル選択を決定することができる。このようにして、次いで、プロセスは、1つ又は複数の統計的モデルを複数の核酸配列リードの各核酸配列リードに順相補性配向及び逆相補性配向の両方で適用し、最大対数尤度値を有するモデルを識別することによって選択された最大事後状態経路の最終的リード当たりのモデル選択を決定することができる。

【0184】

次いで、工程1230において、プロセスは、複数の核酸配列リードの各核酸配列リードを、最大事後状態経路の最終的リード当たりのモデルによって識別される転位部位によって区画された個別の配列要素にセグメント化することができ、これにより、複数の核酸

配列リード内の個別の配列要素を識別することができる。

【0185】

次いで、工程1235において、プロセスは、複数の核酸配列リード内で同定された別の配列要素を、配列要素データファイルに保存し得る。簡略化された手順1700は、例示的に、新しいプロセスが開始されるまで、工程1240で終了することができる。

【0186】

キット

本開示はまた、本開示の方法で使用するための本開示の薬剤を含有するキットを提供する。本開示のキットは、本開示の薬剤及び/又は組成物を含む1つ又は複数の容器を含み得る。いくつかの実施形態において、キットは、本開示の方法に従って使用するための説明書を更に含む。

10

【0187】

本開示のキットで提供される説明書は、典型的には、ラベル又は添付文書(例えば、キットに含まれる紙のシート)に記載された説明書であるが、機械可読説明書(例えば、磁気又は光ストレージディスク上で実行される命令)も許容される。本明細書に記載の方法のいずれかを実施するための説明書を提供することができる。

【0188】

本開示のキットは適切な包装中にある。適切な包装としては、バイアル、ボトル、瓶、フレキシブル包装(例えば、密封されたマイラー又はビニール袋)等が挙げられるが、これらに限定されない。容器は、薬学的に活性な薬剤を更に含む得る。

20

【0189】

キットは、必要に応じて、バッファ及び説明的情報等の追加の構成要素を提供することができる。通常、キットは、容器と、容器上の又は容器に関連するラベル又は添付文書(複数可)とを含む。

【0190】

本開示の実施は、特に明記しない限り、当業者の技能の範囲内である化学、分子生物学、微生物学、組換えDNA、遺伝学、免疫学、細胞生物学、細胞培養及びトランスジェニック生物学の従来技術を使用する。例えば、Maniatis et al., 1982, Molecular Cloning (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y.); Sambrook et al., 1989, Molecular Cloning, 2nd Ed. (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y.); Sambrook and Russell, 2001, Molecular Cloning, 3rd Ed. (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y.); Ausubel et al., 1992, Current Protocols in Molecular Biology (John Wiley & Sons, including periodic updates); Glover, 1985, DNA Cloning (IRL Press, Oxford); Anand, 1992; Guthrie and Fink, 1991; Harlow and Lane, 1988, Antibodies, (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y.); Jakoby and Pastan, 1979; Nucleic Acid Hybridization (B.D. Hames & S.J. Higgins eds. 1984); Transcription And Translation (B.D. Hames & S.J. Higgins eds. 1984); Culture Of Animal Cells (R.I. Freshney, Alan R. Liss, Inc., 1987); Immobilized Cells And Enzymes (IRL Press, 1986); B. Perbal, A Practical Guide To Molecular Cloning (1984); th

30

40

50

e treatise, Methods In Enzymology (Academic Press, Inc., N.Y.); Gene Transfer Vectors For Mammalian Cells (J.H. Miller and M.P. Calos eds., 1987, Cold Spring Harbor Laboratory); Methods In Enzymology, Vols. 154 and 155 (Wu et al. eds.), Immunochemical Methods In Cell And Molecular Biology (Mayer and Walker, eds., Academic Press, London, 1987); Handbook Of Experimental Immunology, Volumes I - IV (D.M. Weir and C.C. Blackwell, eds., 1986); Riott, Essential Immunology, 6th Edition, Blackwell Scientific Publications, Oxford, 1988; Hogan et al., Manipulating the Mouse Embryo, (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y., 1986); Westerfield, M., The zebrafish book. A guide for the laboratory use of zebrafish (Danio rerio), (4th Ed., Univ. of Oregon Press, Eugene, 2000)を参照されたい。

【0191】

他に定義されない限り、本明細書で使用される全ての技術用語及び科学用語は、本開示が属する技術分野の当業者によって一般的に理解されるのと同じ意味を有する。本明細書に記載の方法及び材料と類似又は同等の方法及び材料を本開示の実施又は試験に使用することができるが、適切な方法及び材料を以下に記載する。本明細書で言及される全ての刊行物、特許出願、特許、及び他の参考文献は、その全体が参照により組み込まれる。矛盾する場合、定義を含む本明細書が優先する。更に、材料、方法、及び例は例示にすぎず、限定することを意図するものではない。

【0192】

ここで、本開示の例示的な実施形態を詳細に参照する。本開示は、例示的な実施形態に関連して説明されるが、本開示をそれらの実施形態に限定することを意図するものではないことが理解されよう。逆に、添付の特許請求の範囲によって定義される本開示の趣旨及び範囲内に含まれ得る代替、修正、及び均等物を網羅することが意図されている。当技術分野で周知の標準的な技術又は以下に具体的に記載される技術を利用した。

【実施例】

【0193】

実施例1：C A s e qプロセス

最近の試みは、単一細胞遺伝子発現試料からアイソフォーム配列決定を行うためにロングリード配列決定プラットフォームを活用してきたが、それらのワークフローは、これまで、不十分なスループット及び実質的な配列決定アーチファクトに悩まされており、リードの約35～50%しかフィルタを通過せず、フローセル当たり約300,000個の配列決定された転写物に相当する(約650～800ドル)。ある態様において、本開示は、例えばPacific Biosciences (PacBio (登録商標))からの最近更新されたSequel IIプラットフォーム上で、10x単一細胞遺伝子発現試料からのハイスループット完全転写配列決定を可能にする「C A s e q」プロセスを提供する。本開示のC A s e qプロセスの使用は、観察される配列決定アーチファクトの割合を10%未満に減少させることを可能にし、一方で、全長配列決定出力をフローセルあたり約25M全長転写物に増強することも可能にする。これを達成するために、多重ライゲーションのために、15塩基対(bp)の相補的配列を増幅し、全長cDNAライブラリに付加するためのdU含有プライマーのファミリーが設計されている。アーチファクト配列の主要な供給源に対処するために、例示されたプロセスは、全長cDNAアンプリコン

の精製を可能にするためにビオチン化プライマーを使用する。効率的な多重化アセンブリを駆動し、不適切なライゲーション事象を軽減するために、本明細書に例示される15bp相補的配列は、全ての配列が互いに少なくとも11ハミング距離単位離れていることを確実にすることによって、最小の類似性を有するように設計された(Buschmann, T. Bioconductor version: Release (3.11). DOI: 10.18129/B9.bioc.DNABarcodes)。更なる設計上の考慮事項は、15~20kbの多重化アレイの生成、すなわち、Sequel IIの出力及び塩基呼出し精度のバランスをとるための現在の最適な長さを保証することであった。適切なサイズのライブラリは、cDNAのサイズ分布に基づいて、集められた断片の数をプログラムすることによって構築される。多重ロングリード及び単一細胞遺伝子発現データを処理及び統合するために、分析パイプラインも調製する。

10

【0194】

実施例2：予備実験で効率的に生成された線状キメラアレイのCaseq

予備的なCaseqランにおいて、1.2kbの平均断片サイズを有するcDNAライブラリからの8断片多重化アセンブリを行い、ライゲーション時に約10kbの多重化断片を得た(図2A)。多重化ライブラリをSequel IIで配列決定し、これにより、合計約2.5Mのリードが得られ、逆多重化後に約23Mの転写物が得られ、これはスループットの約9倍の増加を表した(図2B)。逆多重化されたリードの分析により、元のcDNAライブラリと同様のサイズ分布が確認された(図2A)。

【0195】

例示されたcDNAライブラリサイズ分布は、効果的な線状キメラアレイを形成することを可能にしたが、サイズ選択はまた、ある特定の状況下では、キメラアレイからの効果的な配列収率を増加させることが予想される入力核酸ライブラリ(例えば、キメラアレイライゲーションプロセスの実施前に、電気泳動又は入力核酸ライブラリの他の分離を介して)に対して行うこともでき、特に個々のリード長がメガベースである場合、配列された別個の配列の総数が多く、及び/又は核酸サイズ範囲の元の分布が分散していることが更に企図される。

20

【0196】

実施例3：改善されたデータアノテーション、逆多重化及びセグメント化方法によるCaseqリード収率の向上

30

本開示のキメラアンプリコンアレイの最初の処理は、既存のゲノムリードアライメントソフトウェアに基づく反復アダプタ発見戦略を用いた既存のサーキュラーコンセンサスシーケンシング(circular consensus sequencing)(CCS)の正確な高忠実度ロングリード(Hifiリード)プロセスを用いた。このプロセスは、本キメラアンプリコンアレイのロングリードからの配列データの抽出に最適ではないと確認され、Caseqリードの分析のための改良された方法の開発が開始された。それによって、以下の実施による、キメラアンプリコンアレイ配列決定リードの統計的アノテーション、逆多重化、及びセグメント化を含む「Longbow」と呼ばれる改良されたCaseqリード解析プロセスが設計された：

(1)アンプリコンアレイ配列及びそれらの間の転位を識別するための1つ又は複数の統計的アノテーションモデル(例えば、複数のリンクされたサブモデルを有するプロファイル隠れマルコフモデル)を使用したキメラアンプリコンアレイ配列決定データのアノテーション。その1つ又は複数の統計的アノテーションモデルは、(a)キメラアンプリコンアレイ配列決定リード全体に散在する先験的に予想される核酸配列(すなわち、アダプタ配列)を認識するための生成統計的アライメントモデル；(b)先験的に知られていない配列(例えばcDNA転写物配列)を認識するための、又は後の処理工程(例えば、単一細胞バーコード配列、固有の分子識別子)で異なる考慮事項に値するほど大きい配列の辞書からの、ランダム統計アライメントモデルを含み、転位が各モデルの末端に配置され、そのアダプタ配列モデル内の内部位置内では許容されない；

40

(2)最大対数尤度値を有するモデルを評価し、それによってキメラアンプリコンアレイ

50

イ配列決定リードを逆多重化することによって決定された、最大事後状態経路の最終的リード当たりのモデル選択の決定による順相補配向及び逆相補配向の両方における各長いリードへの上記工程(1)の統計的アノテーションモデルの反復適用；及び

(3)上記の工程(1)及び(2)の実行によって同定された部位におけるキメラアンブリコンアレイ配列決定リードのセグメント化。

【0197】

上に開示された「Longbow」プロセスは、少なくとも、(1)サーキュラーコンセンサスシーケンシング(CCS)ソフトウェアによって最初に同定されたリードの集団から実際には低品質の配列リードを、名目上高品質であると識別し、除去すること、(2)サーキュラーコンセンサスシーケンシング(CCS)ソフトウェアによって最初に使用不可能な品質であると主張されるリードの集団から高品質の配列リードを救済すること、(3)「Longbow」プロセスから新たに同定された高品質リードの品質を概算すること、への適用を照らして、本開示のキメラアンブリコンアレイからの品質管理及び配列データ収率の増強に有用であると更に確認された。そのような各用途は、以下で更に詳細に検討される。

【0198】

本開示のキメラアンブリコンアレイ配列決定から潜在的に低品質のデータを識別するために、方法は、(a)シーケンサによって高品質であると確認されたキメラアンブリコンアレイ配列決定リードに(上記のような)Longbowモデルを適用すること(それによって、これらのリードのそれぞれにおける各ヌクレオチドを、それが由来するライブラリアダプタ配列で標識すること)；(b)等しい隣接するLongbowのヌクレオチド標識を、その標識された部分全体を含む領域にマージすること；及び(c)全ての標識されたリードにわたって反復し、その順序で生じない標識された部分を有する任意のリードを、そのライブラリ調製により予想される順序で同定すること、を含む。最初の予想されるセグメントの後に始まるが、その残りのセクションが順番通りであるリード、及び最後の予想されるセグメントの前に終わるが、その前のセクションが全て順番通りであるリード、並びにこれらの場合の組み合わせはこれから除外される。予想されるライブラリに適合しないリードは、低品質と見なされる。

【0199】

シーケンサによって低品質で使用不可能であると報告されたサブセットから高品質の配列決定データを同定するために、方法は以下の工程を含む。(a)シーケンサが使用不可能な品質として報告したデータ(すなわち、リード)を同定すること。そのような使用不可能な品質データは、データに非常に低いリード品質スコア(0未満の値、0~0.5の値、及び0.5~1.0の値を含むが、これらに限定されない)を割り当てるサーキュラーコンセンサスシーケンシングのソフトウェアによって、又はリードを「ZMWパスフィルタ」以外の任意のカテゴリに割り当てるサーキュラーコンセンサスシーケンシングのソフトウェアのいずれかによって決定される。(b)使用不可能な品質のこれらのリードに(上記のような)Longbowモデルを適用し、それにより、これらのリードのそれぞれにおける各ヌクレオチドを、それが由来するライブラリアダプタ配列で標識すること。(c)等しい隣接するLongbowのヌクレオチド標識を、標識された部分全体を含む領域にマージすること。そして、(d)全ての標識されたリードを反復し、最初の予想されるセグメントの後に始まるが、その残りのセクションが順番であるリード、及び最後の予想されるセグメントの前に終わるが、その前のセクションが順番であるリード、並びにこれらの場合の任意の組み合わせを含めて、そのライブラリ調製により出現すると予想される順序で標識された部分を有する任意のリードを同定すること。そのようなリードは、そのリードが更なる分析のために十分に高品質であることを示す予想されるライブラリ調製に適合する。前述のプロセスは、サーキュラーコンセンサスシーケンシングのソフトウェアによって、例えば、0.99未満のリード品質が割り当てられた、又は「ZMWパスフィルタ」以外の任意のカテゴリが割り当てられた、使用不可能なデータに適用するための例示であるが、このプロセスはまた、任意の主張される品質の任意のリード又はリードの

10

20

30

40

50

集団にも適用できることが明示的に記載される。

【0200】

Longbowプロセスの新たに同定された高品質リードの品質を近似するために、方法は以下の工程を含む。(a)新たに同定された高品質リードの各々の標識された部分について、標識された部分のヌクレオチドとその部分に対する予想される配列との間のアライメントスコアを計算すること。このアライメントスコアは、Smith-Waterman又はNeedleman-Wunschアルゴリズム等の動的プログラミングアルゴリズムを使用して直接計算することができ、又は、標識された部分と予想される配列との間のレーベンシュタイン距離を計算して、その距離をその予想される配列の長さから減算することによって直接計算することができる。(b)このアライメントスコアを最良のアライメントスコア(予想される配列とそれ自体との間のアライメントスコアを計算することによって得ることができる)で除算して、各セクションの品質を得ること。そして、(c)(a)で計算された全てのアライメントスコアを合計して、全体のアライメントスコアを得ること。(b)で計算された全ての最良のアライメントスコアを合計して、全体の最良のアライメントスコアを得る。全体のアライメントスコアと全体の最良のアライメントスコアとの比は、リードの推定品質である。

10

【0201】

実施例4:COVID-19患者試料の評価のための拡張性のある単一細胞アイソフォーム配列決定ワークフローにおけるCaseqの実装

単一細胞遺伝子発現研究からの遺伝子アイソフォーム組成の解明は、以前は不可能であった。選択的スプライシングは、転写物の成熟中に差次的なエクソンスプライシングによって内在性タンパク質の構造及び機能を調節するコア調節プロセスである。選択的スプライシングから得られる遺伝子アイソフォームは、細胞のシグナル伝達及び機能の媒介において中心的役割を果たすことが示されている(Baralle and Giudice, Nat Rev Mol Cell Biol 18:437-451)。細胞発生及び恒常性維持を超えて、遺伝子アイソフォームは、複数の病状又は腫瘍の進行及び耐性を駆動する異常なスプライシングに関連する顕著なアイソフォームを有する複数の病状に関連している(Kim et al, Pflugers Arch-Eur J Physiol 470:995-1016; Scotti and Swanson, Nat Rev Genet 17:19-32)。単一細胞解像度でアイソフォーム組成物を効果的に捕捉することができないことは、不均一な生物学的系を効果的に特徴付けるための上述の方法の能力に重大な欠陥があることを強調する。

20

30

【0202】

本実施例では、本開示のCaseqプロセスを用いて、単一細胞遺伝子発現試料に対してハイスループットアイソフォーム配列決定を実施する。アイソフォーム及び単一細胞遺伝子発現データを処理及び統合するためのパイプラインは、当技術分野で認識されている分析ツールを使用して開発される。標的化アイソフォーム配列決定のために、遺伝子パネルも開発されている。免疫応答及び感染組織の両方を特徴付けるために、COVID-19患者を評価する。

【0203】

COVID-19症状は、部分的には、SARS-CoV-2感染に対する過活動免疫応答に起因して生じる。本開示の実施例では、CaseqをCOVID-19試料(300名のCOVID-19患者の血液及び約10名の剖検由来の組織からの免疫区画の進行中の単一細胞ゲノム研究に由来する)に対し使用し、疾患の重症度に関連する免疫細胞クラスターにおいて差次的に発現されるアイソフォームを発見することを目的とする。

40

【0204】

(非Caseq)予備データの初期セットは、健康な患者と軽度及び重度のCOVID-19患者との間の単球区画における著しい転写の違いを確認している(図10A~10D)。アイソフォーム分析は、炎症及び単球活性化経路に関連する遺伝子に焦点を当てているが、これらに限定されない(doi.org/10.1093/nar/gky40

50

1 及び doi.org/10.1038/s41467-019-11076-1 を参照されたい)。アイソフォーム分析の出力を高めるために、Leiden クラスタを一緒にグループ化して、クラスタ間の差別的なアイソフォーム組成のより堅牢な統計的比較を可能にする。SARS-CoV2 感染試料を健康な対照患者と比較して、遺伝子発現の違い及び選択的スプライシングの役割を特徴付けた。SARS-CoV2 は、そのゲノムからの転写の複雑な不連続プロセスを利用することが示され、ショートリード配列決定は特にウイルス遺伝子発現の解明に適さないので、SARS-CoV2 トランスクリプトームの再構築は洞察力が期待される。感染過程にわたる潜在的な転写動態に光を当てるために、感染細胞でのウイルス転写物の組成及び量との潜在的な関連がそれによって調査される。

10

【0205】

実施例 5：単一細胞遺伝子発現試料からのミトコンドリア系統追跡

腫瘍内不均一性及びクローン進化は、腫瘍進行及び治療耐性を可能にする推進力である。クローン動態を追跡する能力は、治療に直面して腫瘍がどのように進化しているかを理解するために重要である。最近のアプローチは、ミトコンドリア変異がクローン同一性を推論するためのマーカーとして役立ち得ることを実証している (Ludwig et al., Cell 176:1325-1339)。そのようなアプローチは、ミトコンドリアゲノムが核ゲノムと比較してはるかに高い割合 (10~100 倍) で突然変異を起こし、配列決定データに高度に表されているという事実の部分的に依存している。当技術分野で認識されているショートリード単一細胞遺伝子発現ワークフローからのカバレッジは制限されるために、研究者らはこれまで、クローン推論に必要なミトコンドリアゲノムの均一かつ十分なカバレッジを提供するために、単一細胞 ATAC (トランスポザーゼアクセシブルクロマチンのアッセイ (Assay for Transposase Accessible Chromatin)) 配列決定に依存してきた。本実施例では、本開示の CasEq アプローチを適用して、単一細胞遺伝子発現試料からの完全ミトコンドリア転写物の標的化ロングリード配列決定を実施し、それによって遺伝子発現試料とのクローン同一性の統合を可能にする。現在のミトコンドリア系統追跡バイオインフォマチックパイプラインを適用し、当技術分野で認識されている方法に対してベンチマークを実施して、全長転写物データで動作するように適合させる。次いで、患者の腫瘍試料を、本 CasEq プロセスを使用して評価して、治療の過程にわたるクローン動態を明らかにする。全ミトコンドリア転写物の CasEq 対応標的化ロングリード配列決定によってクローン情報を抽出する能力は、クローン性と同じ試料からの遺伝子発現との連結を提供する。クローン性及び遺伝子発現のこのような協調的評価は、進行及び治療耐性の過程にわたって腫瘍におけるクローン進化の研究を劇的に向上させる。

20

30

【0206】

実施例 6：単一細胞遺伝子発現試料からのミトコンドリア転写物捕捉及び多重ライゲーションの最適化

これまで、単一細胞遺伝子発現ワークフローは、対立遺伝子情報を、個々の細胞からのクローン関係の堅牢な再構築を可能にする程度まで捕捉するには不十分であった。広く使用されている単一細胞遺伝子発現データから得られたクローン関係を明らかにする能力は、深い洞察を促進し、遺伝子発現状態、クローン性及び細胞運命間の連結を同定することを可能にするため、このことは計り知れない機会の損失を表している。これまで単一細胞遺伝子発現試料からのクローン再構築を妨げてきた低いカバレッジに対処するために、本明細書に開示される CasEq はまた、全長ミトコンドリア転写物配列情報を得ることを標的とする。ミトコンドリア転写物の高効率配列決定は、本明細書の他の箇所に記載されている多重化プライマーを使用してミトコンドリアから発現される 13 個の遺伝子の標的化増幅を行うことによって達成される。配列決定出力及び忠実度のバランスをとりながら、15~20 kb の最適な多重化アレイ長を確保するために、ミトコンドリア cDNA プールの長さ分布を考慮して、組み立てられた断片の数が増える。配列決定されると、全長転写物は、マッピング及び塩基品質のために逆多重化及びフィルタリングされる。リ

40

50

ード通過フィルタは、ミトコンドリアゲノムのカバレッジを定量化するために使用される。既存のミトコンドリア系統追跡パイプラインはまた、クローン関係の再構築のために全長ミトコンドリア転写物を使用するように適合されている。

【0207】

実施例7：全長ミトコンドリア転写物系統追跡のベンチマーキング

完全長ミトコンドリア転写物系統追跡を検証するために、安定に組み込まれたDNAバーコードを有するHeLa細胞株集団からクローン関係を再構築する能力を定量化し、これは、クローン同一性のためのグラウンドトゥルースを確立するのに役立ち得る。具体的には、ClonMapper発現バーコードシステム(単一細胞RNA配列決定によるクローン同定を可能にする以前に開発されたシステム)でタグ付けされた細胞を使用する。更に、Ludwig et al. (Cell 176:1325-1339)に記載の方法を、バーコード化集団の並行試料に対して実施し、特異性及びリコールに関連する測定値をクローン同一性の割り当てのために計算し、比較する。

10

【0208】

したがって、本明細書に開示されるCaseqプロセスは、既存のプラットフォームによってこれまで達成できなかった配列決定スループット及びリード長を可能にするため、配列決定の分野における重要な進歩を提供する。更に、本Caseqプロセスは、高度に適合可能であり、目的の遺伝的特徴を捕捉するために容易に特殊化することができる。本開示に記載されたCaseqの実装形態は、発見のための新しいプラットフォームとして提供され、多くの科学分野に広く適用可能である。本Caseqアプローチは、ロングリードプラットフォームと共進化する能力を有し、それらのリード長が増加し続けるにつれてそれらの分子出力を更に高めるのに役立つ。

20

【0209】

参考文献

1. I. Gupta et al., Single-cell isoform RNA sequencing characterizes isoforms in thousands of cerebellar cells. Nat Biotechnol. 36:1197-1202 (2018).

2. R. Volden et al., Improving nanopore read accuracy with the R2C2 method enables the sequencing of highly multiplexed full-length single-cell cDNA. Proc Natl Acad Sci U S A 115:9726-9731 (2018).

30

3. M. Singh et al., High-throughput targeted long-read single cell sequencing reveals the clonal and transcriptional landscape of lymphocytes. Nat Commun. 10:3120 (2019).

【0210】

本明細書で言及される全ての特許及び刊行物は、本開示が関係する当業者の技術レベルを示す。本開示において引用された全ての参考文献は、あたかも各参考文献が個別にその全体が参照により組み込まれたのと同程度に、参照により組み込まれる。

40

【0211】

当業者は、本開示が目的を実行し、言及された目的及び利点、並びにそれらに固有の目的及び利点を得るようによく適合されていることを容易に理解するであろう。本開示の好ましい実施形態の代表例として本明細書に記載される方法及び組成物は例示的なものであり、本開示の範囲に対する限定として意図されるものではない。その中の変更及び他の使用が当業者には思い浮かぶであろうが、それらは本開示の趣旨の範囲内に包含され、特許請求の範囲によって定義される。

【0212】

50

更に、本開示の特徴又は態様がマーカッシュ群又は代替物の他のグループ化に関して記載されている場合、当業者は、本開示がそれによってマーカッシュ群又は他の群の任意の個々のメンバー又はメンバーのサブグループに関する記載されることを認識するであろう。

【0213】

本開示を説明する文脈において（特に以下の特許請求の範囲の文脈において）「a」及び「an」及び「the」という用語並びに同様の指示対象の使用は、本明細書で特に指示されない限り、又は文脈と明らかに矛盾しない限り、単数及び複数の両方を包含すると解釈されるべきである。「含む（comprising）」、「有する（having）」、「を含む（including）」、及び「含有する（containing）」という用語は、特に明記しない限り、オープンエンド用語（すなわち、「を含むが、限定されない」を意味する）と解釈されるべきである。本明細書における値の範囲の列挙は、本明細書に別段の指示がない限り、範囲内に含まれる各別個の値を個別に参照する簡略方法として作用することを意図しているにすぎず、各別個の値は、本明細書に個別に列挙されているかのように本明細書に組み込まれる。

10

【0214】

本明細書に記載の全ての方法は、本明細書に別段の指示がない限り、又は文脈と明らかに矛盾しない限り、任意の適切な順序で実行することができる。本明細書で提供されるありとあらゆる例又は例示的な言語（例えば、「等」）の使用は、単に本開示をよりよく明らかにすることを意図しており、別段の請求がない限り、本開示の範囲を限定するものではない。本明細書におけるいかなる言語も、特許請求されていない要素を本開示の実施に必須であると示すと解釈されるべきではない。

20

【0215】

開示された発明を実施するための本発明者らに知られている最良の形態を含む、本開示の実施形態を本明細書で説明する。これらの実施形態の変形は、前述の説明を読めば当業者には明らかとなり得る。

【0216】

本明細書に例示的に記載された開示は、本明細書に具体的に開示されていない任意の1つ又は複数の要素、1つ又は複数の制限がない状態で適切に実施することができる。したがって、例えば、本明細書の各例では、「を含む（comprising）」、「から本質的になる（consisting essentially of）」、及び「からなる（consisting of）」という用語のいずれかは、他の2つの用語のいずれかと置き換えることができる。使用された用語及び表現は、限定ではなく説明の用語として使用され、そのような用語及び表現の使用において、示され説明された特徴又はその一部の均等物を除外することは意図されていないが、特許請求される発明の範囲内で様々な修正が可能であることが認識される。したがって、本開示は好ましい実施形態を提供するが、本明細書に開示された概念の任意選択の特徴、修正及び変形は当業者によって使用されてもよく、そのような修正及び変形は、説明及び添付の特許請求の範囲によって定義される本開示の範囲内にあると見なされることを理解されたい。

30

【0217】

本発明の範囲及び趣旨から逸脱することなく、本明細書に開示された発明に対して様々な置換及び修正を行うことができることは、当業者には容易に明らかであろう。したがって、そのような追加の実施形態は、本開示及び以下の特許請求の範囲の範囲内である。本開示は、改善されたコントラスト、診断及び/又はイメージング活性を有するコンジュゲートを生成するために、本明細書に記載の化学修飾の様々な組み合わせ及び/又は置換を試験することを当業者に教示する。したがって、本明細書に記載されるある実施形態は限定的ではなく、当業者は、本明細書に記載される修飾のある組み合わせが、改善されたコントラスト、診断及び/又はイメージング活性を有するコンジュゲートを同定するための過度の実験なしに試験され得ることを容易に理解することができる。

40

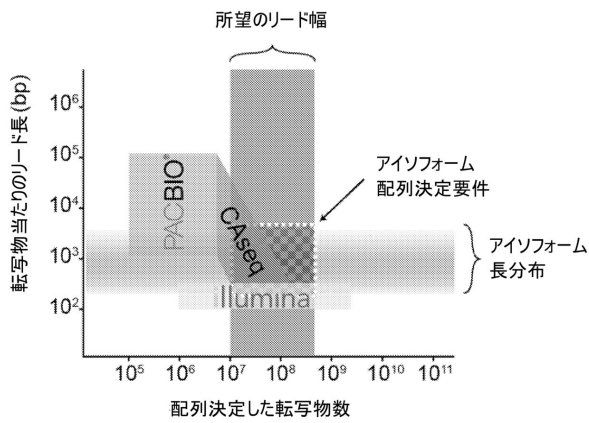
【0218】

50

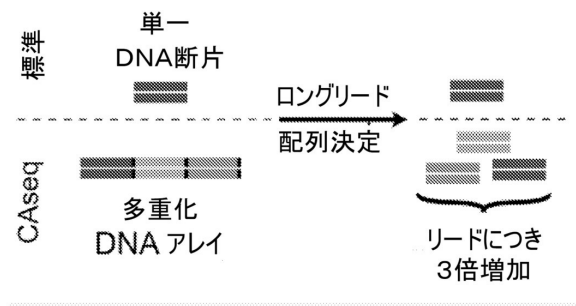
本発明者らは、当業者がそのような変形形態を適切に使用することを期待しており、本発明者らは、本開示が本明細書に具体的に記載されている以外の方法で実施されることを意図している。したがって、本開示は、適用法によって許容されるように、添付の特許請求の範囲に列挙された主題の全ての修正及び均等物を含む。更に、本明細書に別段の指示がない限り、又は文脈と明らかに矛盾しない限り、その全ての可能な変形における上述の要素の任意の組み合わせが本開示に含まれる。当業者は、本明細書に記載の開示のある実施形態に対する多くの均等物を認識するか、又は日常的な実験のみを使用して確認することができるであろう。そのような均等物は、以下の特許請求の範囲に含まれることが意図されている。

【 図 面 】

【 図 1 A 】



【 図 1 B 】



10

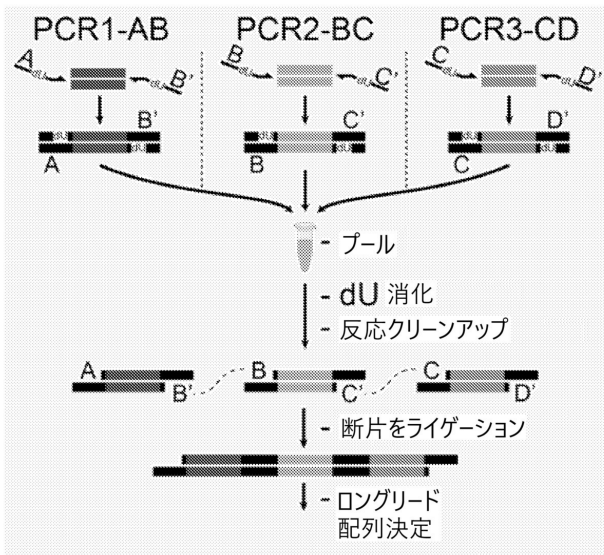
20

30

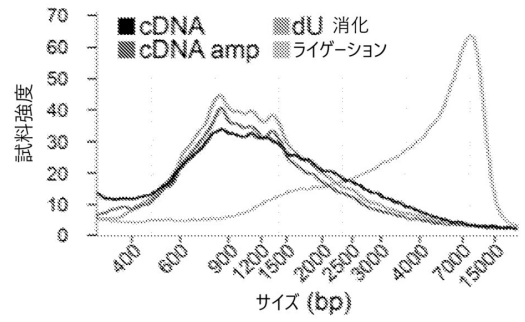
40

50

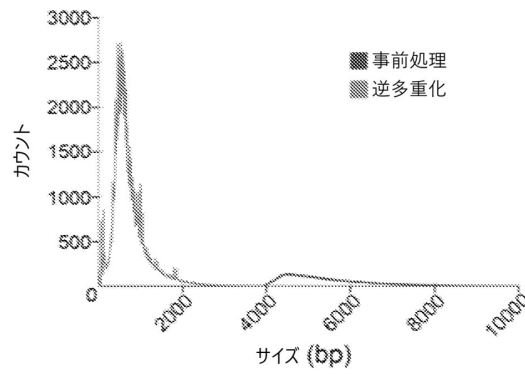
【 図 1 C 】



【 図 2 A 】

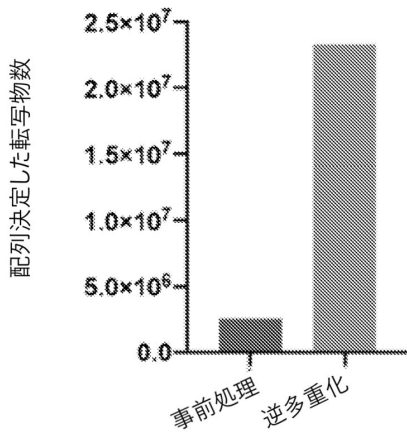


10

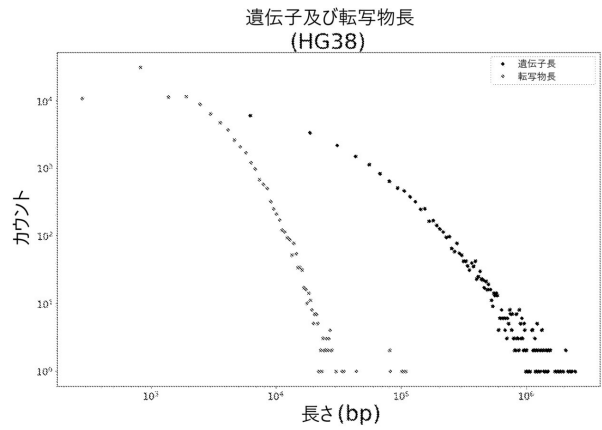


20

【 図 2 B 】



【 図 3 A 】



30

40

50

【 図 7 】

	ライゲーションプロファイル	長さ	カウント	全カウントのうちの百分率
0	ABCDEFGHIJKLMNPO	16	484869	29.88%
1	A'B'C'D'E'F'G'H'I'J'K'L'M'N'O'P'	16	476205	29.34%
2	C'D'E'F'G'H'I'J'K'L'M'N'O'P'	14	54879	3.38%
3	CDEFGHIJKLMNPO	14	54471	3.36%
4	ABCDEFGHIJKLMN	14	49074	3.02%
5	ABCDEFGHIJKLM	13	46094	2.84%
6	A'B'C'D'E'F'G'H'I'J'K'L'M'N'	14	40890	2.52%
7	A'B'C'D'E'F'G'H'I'J'K'L'M'	13	39838	2.45%
8	ABCDEFGHIJKLM	12	31918	1.97%
9	D'E'F'G'H'I'J'K'L'M'N'O'P'	13	30846	1.90%
10	DEFGHIJKLMNPO	13	29839	1.84%
11	A'B'C'D'E'F'G'H'I'J'K'L'	12	28700	1.77%
12	E'F'G'H'I'J'K'L'M'N'O'P'	12	28590	1.76%
13	EFGHIJKLMNPO	12	28128	1.73%
14	ABCDEFGHIHIJK	11	14138	0.87%
15	ABCDEFGHIJKLMNO	15	13311	0.82%
16	B'C'D'E'F'G'H'I'J'K'L'M'N'O'P'	15	12297	0.76%
17	A'B'C'D'E'F'G'H'I'J'K'	11	11336	0.70%
18	A'B'C'D'E'F'G'H'I'J'K'L'M'N'O'	15	10711	0.66%
19	BCDEFGHIJKLMNPO	15	10495	0.65%

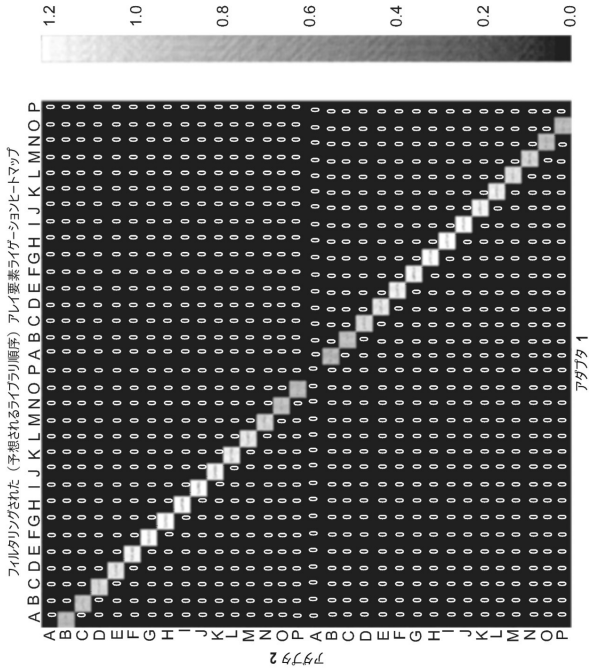
【 図 8 】

HiFiリードの数	キメラアンプリコンアレイ配列決定リードの数	現行のリード分析法によってリクレームされたキメラアンプリコンアレイ配列決定リードの数	キメラアンプリコンアレイ配列決定リードの総数	総キメラアンプリコンアレイ配列決定率増加
2477400	29789780	10194059	39983839	16.14x
1622783	22662840	12794497	35457337	21.84x

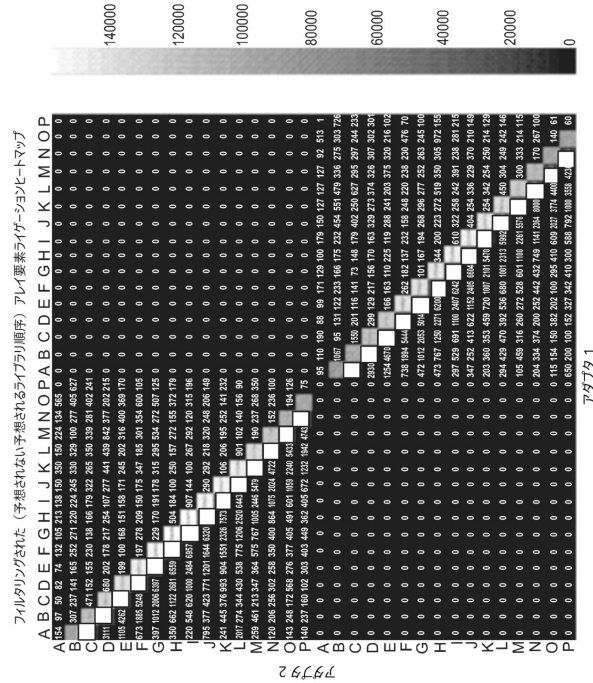
10

20

【 図 9 A 】



【 図 9 B 】

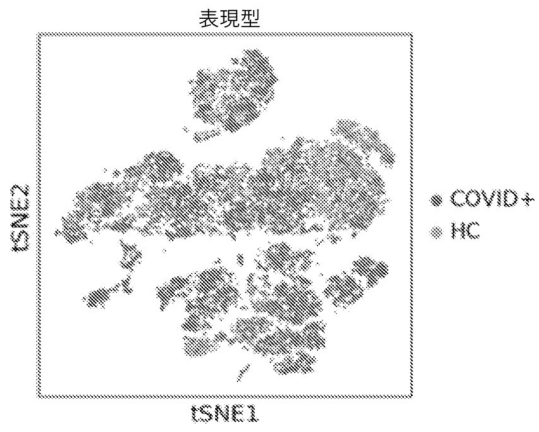


30

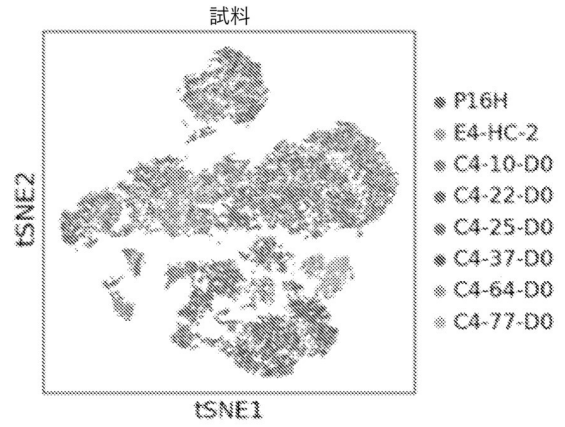
40

50

【 図 1 0 A 】

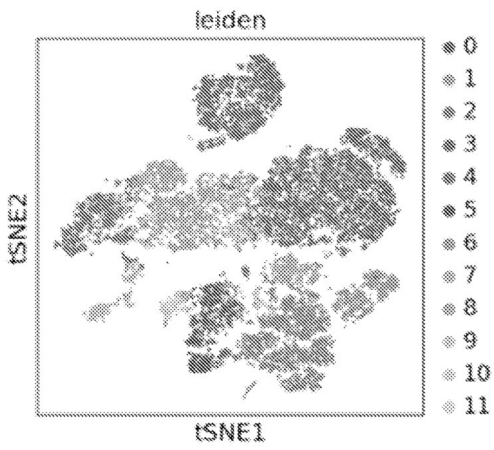


【 図 1 0 B 】

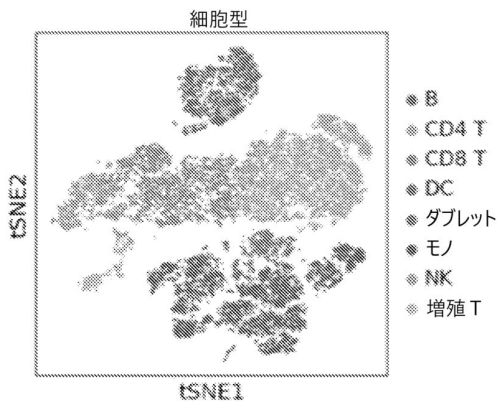


10

【 図 1 0 C 】



【 図 1 0 D 】



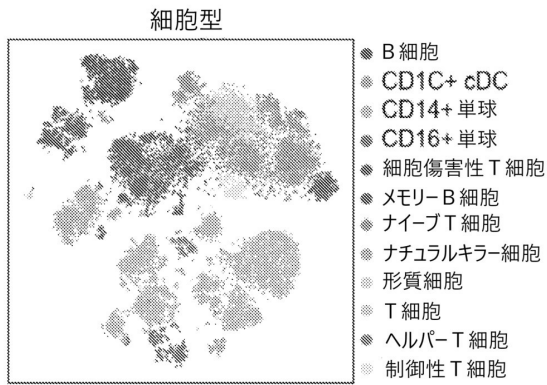
20

30

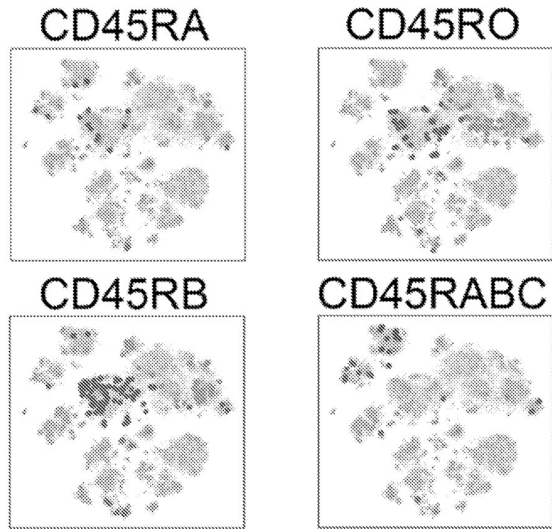
40

50

【 図 1 1 A 】



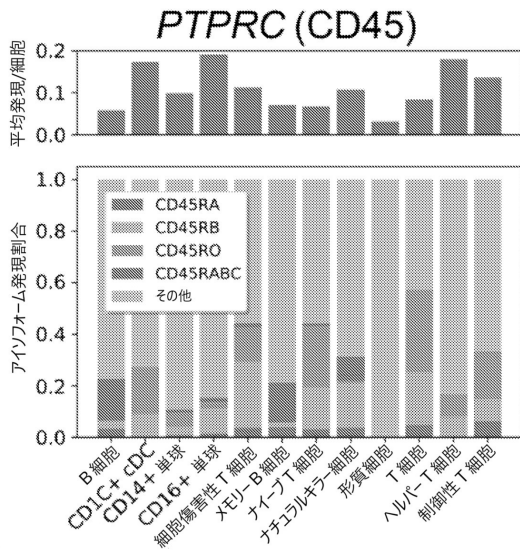
【 図 1 1 B 】



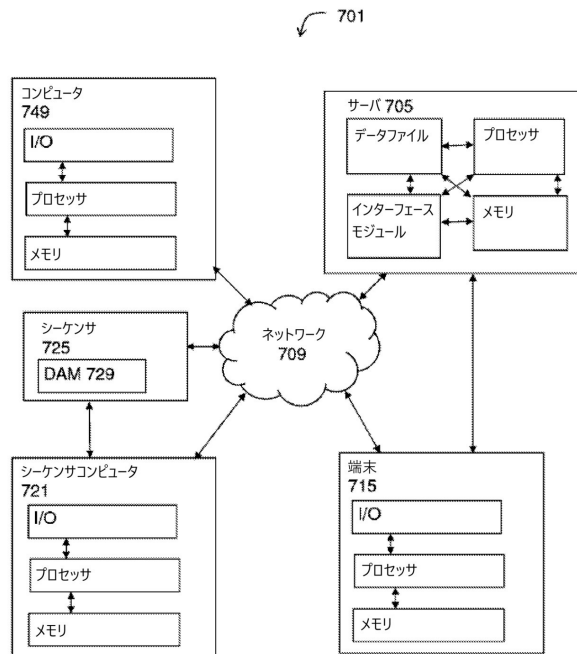
10

20

【 図 1 1 C 】



【 図 1 2 】

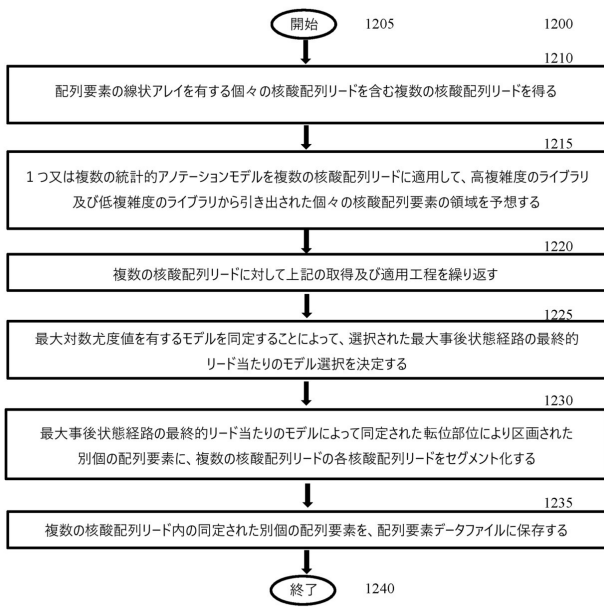


30

40

50

【 図 1 3 】



10

20

【 配列表 】

202353488200001.app

30

40

50

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US 21/37226

A. CLASSIFICATION OF SUBJECT MATTER
 IPC - C12Q 1/68, C12P 19/34 (2021.01)
 CPC - C12Q 1/6837, C12Q 1/6844, C12Q 1/6858, C12Q 1/6874, C12Q 1/6806

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)
 See Search History document

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched
 See Search History document

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)
 See Search History document

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X - Y - A	US 2016/0264958 A1 (TWIST BIOSCIENCE CORPORATION) 15 September 2016 (15.09.2016); especially abstract; para [0004]-[0006], [0009]-[0010], [0030]-[0031], [0034]-[0035], [0044], [0072], [0075], [0084], [0127], [0162]	1-4, 6-9, 11-13, 15, 20-24, 27-29 - 5, 14, 16-19 -
Y	US 2010/0105052 A1 (DRMANAC et al.) 29 April 2010 (29.04.2010); especially para [0067], [0145], [0258], [0366], [0369], [0511]	5, 14, 16, 18
Y	WO 2018/227025 A1 (ARC BIO, LLC) 13 December 2018 (13.12.2018); especially para [00448]-[00447], [00456]-[00457]	17
Y	GIBSON et al. Enzymatic assembly of DNA molecules up to several hundred kilobases. Nature Methods. May 2009, Vol. 6, No. 5, pg 343-345; especially abstract; Figure 1	19
A	US 2009/0099034 A1 (AHLQUIST et al.) 16 April 2009 (16.04.2009); especially para [0018]; SEQ ID NO: 40	10
A	WO 2008/133643 A2 (MONSANTO TECHNOLOGY, LLC) 6 November 2008 (06.11.2008); especially SEQ ID NO: 6165	10

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
 "D" document cited by the applicant in the international application
 "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
10 November 2021

Date of mailing of the international search report
DEC 08 2021

Name and mailing address of the ISA/US
Mail Stop PCT, Attn: ISA/US, Commissioner for Patents
P.O. Box 1450, Alexandria, Virginia 22313-1450
Facsimile No. 571-273-8300

Authorized officer
Kari Rodriguez
Telephone No. PCT Helpdesk: 571-272-4300

10

20

30

40

50

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US 21/37226

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.:
because they relate to subject matter not required to be searched by this Authority, namely:

- 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).

10

20

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:
This application contains the following inventions or groups of inventions which are not so linked as to form a single general inventive concept under PCT Rule 13.1. In order for all inventions to be searched, the appropriate additional search fees must be paid.

Continued on Supplemental Page

- 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 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.:
1-24, 27-29 limited to SEQ ID NO: 1

30

40

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.

50

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US 21/37226

Continued from Box No. III Observations where unity of invention is lacking

Groups I+: Claims 1-24, 27-29, drawn to a method for preparing an array nucleic acid sequence, the method comprising: obtaining a plurality of input nucleic acid sequence within each input nucleic acid is of approximately 30 kilobases in length or shorter; attaching one or more adapter sequences to the plurality of input nucleic acid sequences. The method will be searched to the extent that the adapter sequence encompasses SEQ ID NO: 1. It is believed that claims 1-24, 27-29 encompass this first named invention, and thus these claims will be searched without fee to the extent that they encompass SEQ ID NO: 1. Additional adapter sequence(s) will be searched upon the payment of additional fees. Applicants must specify the claims that encompass any additionally elected adapter sequence(s). Applicants must further indicate, if applicable, the claims which encompass the first named invention, if different than what was indicated above for this group. Failure to clearly identify how any paid additional invention fees are to be applied to the "+" group(s) will result in only the first claimed invention to be searched. An exemplary election would be an adapter sequence comprising SEQ ID NO: 2 (Claims 1-24, 27-29).

10

Group II: Claim 25, drawn to a method for obtaining isoform sequencing information.

Group III: Claim 26, drawn to a method for performing mitochondrial lineage tracing.

Groups IV+: Claims 30-31, drawn to a composition comprising a plurality of adapter sequences. Group IV+ will be searched upon payment of additional fees. The composition may be searched, for example, to the extent that the adapter sequence encompasses SEQ ID NO: 1 for an additional fee and election as such. It is believed that claims 30-31 read on this exemplary invention. Additional adapter sequence(s) will be searched upon the payment of additional fees. Applicants must specify the claims that encompass any additionally elected adapter sequence(s). Failure to clearly identify how any paid additional invention fees are to be applied to the "+" group(s) will result in only the first claimed invention to be searched. Another exemplary election would be an adapter sequence comprising SEQ ID NO: 2 (Claims 30-31).

20

Group V: Claims 32-52, drawn to a method/system for identifying discrete sequence elements within individual nucleic acid sequence reads of a population of nucleic acid sequence reads, wherein each of the linear array of sequence elements comprises two or more nucleic acid sequence elements drawn from a library of high complexity flanked nucleic acid sequences drawn from a library of low complexity.

The inventions listed as Groups I+, II, III, IV+, V do not relate to a single general inventive concept under PCT Rule 13.1 because, under PCT Rule 13.2, they lack the same or corresponding special technical features for the following reasons:

Special Technical Features

Groups I+, II and III include the special technical feature of a method which differs from the special technical feature of a composition, as disclosed by Groups IV+.

Groups I+ include the special technical feature of a method of preparing an array nucleic acid sequence wherein each input nucleic acid sequence within the plurality of input nucleic acid sequences is of approximately 30 kilobases in length or shorter, not required by Groups II, III and V.

Group II includes the special technical feature of a method of analyzing the sequence information obtained from the linear array nucleic acid sequence to obtain isoform sequencing information, not required by Groups I+, III and V.

30

Group III includes the special technical feature of a method of analyzing the sequence information obtained from the array nucleic acid sequence to trace mitochondrial lineage, not required by Groups I+, II and V.

Group V includes the special technical feature of a method for identifying discrete sequence comprising: applying one or more statistical annotation models to sequence data of the population of nucleic acid sequence reads, to predict within the population of nucleic acid sequence reads regions of individual nucleic acid sequence elements drawn from a library of high complexity and regions of nucleic acid sequences drawn from a library of low complexity, not required by Groups I+, II, III and IV+.

No technical features are shared between the polynucleotide sequences of Groups I+ and IV+, accordingly, these groups lack unity a priori.

Common Technical Features

The inventions of Groups I+, II, III, IV+ and V share the technical feature of a linear array of sequence elements.

The inventions of Groups I+, II, III and IV+ share the technical feature of a plurality of nucleic acid sequences, wherein at least two of the plurality of nucleic acid sequences comprise an adapter sequences.

The inventions of Groups I+, II and III share the technical feature of a method for preparing an array nucleic acid sequence.

40

Continued on Next Supplemental Page

INTERNATIONAL SEARCH REPORT

International application No.

PCT/US 21/37226

Continued from Previous Supplemental Page

However, these shared technical features do not represent a contribution over prior art in view of US 2018/0264958 A1 to Twist Bioscience Corporation (hereinafter 'Twist').

Twist discloses (instant claim 27) a method for preparing an array nucleic acid sequence (abstract - "Methods and compositions are provided for assembly of large nucleic acids..."), the method comprising:

i) obtaining a plurality of input nucleic acid sequences (para [0005] - "Methods are provided herein for nucleic acid assembly, comprising: providing a predetermined nucleic acid sequence..."; para [0072] - "In some cases, a precursor nucleic acid sequence is assembled with another precursor nucleic acid sequence via annealing and ligation of complementary sticky ends, followed by additional rounds of sticky end generation and assembly with other precursor fragment(s) to generate a long target nucleic acid sequence"), wherein each input nucleic acid sequence within the plurality of input sequences is of approximately 300 kilobases in length or shorter (para [0005] - "Methods are further provided wherein the predetermined nucleic acid sequence is 1 kb to 100 kb in length.");

ii) contacting the plurality of input nucleic acid sequences with paired amplification primers (para [0005] - "...each precursor double-stranded nucleic acid fragment having two strands, wherein each of the two strands comprises a sticky end sequence... providing primers comprising a nicking endonuclease recognition site and a sequence... corresponding to each of the different sticky end sequences..."), wherein at least one primer within the paired amplification primers comprises an adapter sequence comprising an internal dU on one strand (para [0005] - "...providing primers comprising a nicking endonuclease recognition site and a sequence comprising (i) 5'-A (Nx) M-3' (SEQ ID NO.: 82) corresponding to each of the different sticky end sequences... wherein M is a non-canonical base... Methods are further provided wherein the uracil is in a deoxyuridine-deoxyadenosine base pair. ... Methods are further provided wherein the precursor double-stranded nucleic acid fragments comprise an adaptor sequence comprising the nicking endonuclease recognition site."), and performing at least one round of amplification, thereby generating a population of adapted nucleic acid sequences (para [0005] - "...performing a polynucleotide extension reaction to form double-stranded nucleic acid fragments...");

iii) contacting the population of adapted nucleic acid sequences with Uracil DNA glycosylase and Endonuclease VIII (para [0006] - "Methods are further provided wherein the first nicking enzyme comprises uracil-DNA glycosylase (UDG). ... Methods are further provided wherein the first nicking enzyme comprises endonuclease VIII."; para [0084] - "For example a glycosylase, such as UDG, alone or in combination with an AP endonuclease, such as endonuclease VIII, is used to excise uracil and create a gap."), thereby forming a population of adapted nucleic acid sequences having single-stranded ends (para [0006] - "Methods are further provided wherein the first nicking enzyme comprises endonuclease VIII."); and

iv) contacting the population of adapted nucleic acid sequences having single-stranded ends with a ligase, thereby forming an array nucleic acid sequence (para [0008] - "...and annealing the double-stranded nucleic acid fragments to form a nucleic acid encoding for the predetermined nucleic acid sequence that does not include the nicking endonuclease recognition site. ... Methods are further provided wherein the method further comprises ligating the annealed double-stranded nucleic acid fragments."; para [0072] - "In some cases, a plurality of precursor nucleic acid fragments are prepared with sticky ends, and the sticky ends are annealed and ligated to generate the predetermined target nucleic acid sequence"; para [0127] - "Nucleic Acid Assembly. Provided herein are methods where two or more of the cleavage, annealing and ligation reactions are performed concurrently within the same mixture and the mixture comprises a ligase.").

10

20

Twist discloses (instant claim 30) a composition comprising a plurality of nucleic acid sequences, wherein at least two of the plurality of nucleic acid sequences comprise an adapter sequence (para [0005] - "... synthesizing a plurality of precursor double-stranded nucleic acid fragments... Methods are further provided wherein the precursor double-stranded nucleic acid fragments comprise an adaptor sequence comprising the nicking endonuclease recognition site.").

30

As said technical features were known in the art at the time of the invention, these cannot be considered special technical features that would otherwise unify the groups.

Groups I+, II, III, IV+, V therefore lack unity under PCT Rule 13 because they do not share a same or corresponding special technical feature.

*Note: Claims 31 and 32 run together as there is no whitespace between the end of claim 31 and the claim numbered 32.

40

フロントページの続き

(51)国際特許分類 F I テーマコード (参考)
 C 1 2 N 15/09 (2006.01) C 1 2 N 15/09 Z

MK,MT,NL,NO,PL,PT,RO,RS,SE,SI,SK,SM,TR),OA(BF,BJ,CF,CG,CI,CM,GA,GN,GQ,GW,KM,ML,MR,N
 E,SN,TD,TG),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,IT,JO,JP,K
 E,KG,KH,KN,KP,KR,KW,KZ,LA,LC,LK,LR,LS,LU,LY,MA,MD,ME,MG,MK,MN,MW,MX,MY,MZ,NA,N
 G,NI,NO,NZ,OM,PA,PE,PG,PH,PL,PT,QA,RO,RS,RU,RW,SA,SC,SD,SE,SG,SK,SL,ST,SV,SY,TH,TJ,TM,
 TN,TR,TT,TZ,UA,UG,US,UZ,VC,VN,WS,ZA,ZM,ZW

(特許庁注：以下のものは登録商標)

1 . T R I T O N

2 . J A V A

アメリカ合衆国 マサチューセッツ州 ボストン フルーツ ストリート 5 5

(74)代理人 100189131

弁理士 佐伯 拓郎

(74)代理人 100182486

弁理士 中村 正展

(74)代理人 100147289

弁理士 佐伯 裕子

(72)発明者 ハコーヘン, ニア

アメリカ合衆国 マサチューセッツ州 0 2 1 1 4 ボストン, フルーツ ストリート 5 5, ザ
 ジェネラル ホスピタル コーポレーション内

(72)発明者 アルカファジ, アジズ

アメリカ合衆国 マサチューセッツ州 0 2 1 4 2 ケンブリッジ, メイン ストリート 4 1 5,
 ザ ブロード インスティテュート インコーポレイテッド内

(72)発明者 ブレイニー, ポール

アメリカ合衆国 マサチューセッツ州 0 2 1 4 2 ケンブリッジ, メイン ストリート 4 1 5,
 ザ ブロード インスティテュート インコーポレイテッド内

(72)発明者 ババディ, メルタシュ

アメリカ合衆国 マサチューセッツ州 0 2 1 4 2 ケンブリッジ, メイン ストリート 4 1 5,
 ザ ブロード インスティテュート インコーポレイテッド内

(72)発明者 ガリメラ, キラン ヴィ

アメリカ合衆国 マサチューセッツ州 0 2 1 4 2 ケンブリッジ, メイン ストリート 4 1 5,
 ザ ブロード インスティテュート インコーポレイテッド内

(72)発明者 スミス, ジョナサン セオドル

アメリカ合衆国 マサチューセッツ州 0 2 1 4 2 ケンブリッジ, メイン ストリート 4 1 5,
 ザ ブロード インスティテュート インコーポレイテッド内

F ターム (参考) 4B029 AA07 AA27 BB20 FA15

4B063 QA13 QQ02 QQ04 QQ42 QR10 QR20 QR32 QR55 QS34 QS36