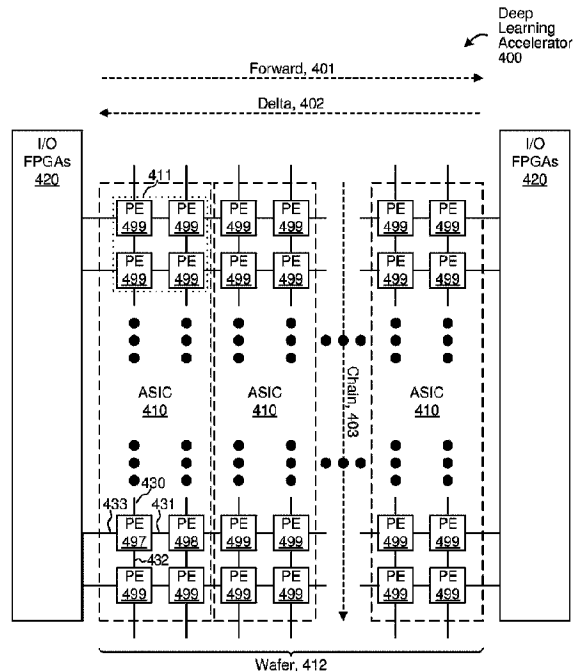




(22) **Date de dépôt/Filing Date:** 2018/02/23
 (41) **Mise à la disp. pub./Open to Public Insp.:** 2018/08/30
 (45) **Date de délivrance/Issue Date:** 2024/02/20
 (62) **Demande originale/Original Application:** 3 051 990
 (30) **Priorités/Priorities:** 2017/02/23 (US62/462,640);
 2017/04/17 (US62/486,372); 2017/06/11 (US62/517,949);
 2017/06/15 (US62/520,433); 2017/06/19 (US62/522,065);
 2017/06/19 (US62/522,081); 2017/08/08 (US62/542,645);
 2017/08/08 (US62/542,657); 2017/11/01 (US62/580,207);
 2018/02/09 (US62/628,784); 2018/02/09 (US62/628,773)

(51) **Cl.Int./Int.Cl.** *G06N 3/08* (2023.01),
G06N 3/045 (2023.01), *G06N 3/048* (2023.01),
G06N 3/063 (2023.01)
 (72) **Inventeurs/Inventors:**
 LIE, SEAN, US;
 MORRISON, MICHAEL, US;
 JAMES, MICHAEL EDWIN, US;
 LAUTERBACH, GARY R., US;
 AREKAPUDI, SRIKANTH, US
 (73) **Propriétaire/Owner:**
 CEREBRAS SYSTEMS INC., US
 (74) **Agent:** OYEN WIGGS GREEN & MUTALA LLP

(54) **Titre : APPRENTISSAGE PROFOND ACCELERER**
 (54) **Title: ACCELERATED DEEP LEARNING**



(57) **Abrégé/Abstract:**

Techniques in advanced deep learning provide improvements in one or more of accuracy, performance, and energy efficiency, such as accuracy of learning, accuracy of prediction, speed of learning, performance of learning, and energy efficiency of learning. An array of processing elements performs flow-based computations on wavelets of data. Each processing element has a respective compute element and a respective routing element. Each compute element has processing resources and memory resources. Each router enables communication via wavelets with at least nearest neighbors in a 2D mesh. Stochastic gradient descent, mini-batch gradient descent, and continuous propagation gradient descent are techniques usable to train weights of a neural network modeled by the processing elements. Reverse checkpoint is usable to reduce memory usage during the training.

Abstract

Techniques in advanced deep learning provide improvements in one or more of accuracy, performance, and energy efficiency, such as accuracy of learning, accuracy of prediction, speed of learning, performance of learning, and energy efficiency of learning. An array of processing elements performs flow-based computations on wavelets of data. Each processing element has a respective compute element and a respective routing element. Each compute element has processing resources and memory resources. Each router enables communication via wavelets with at least nearest neighbors in a 2D mesh. Stochastic gradient descent, mini-batch gradient descent, and continuous propagation gradient descent are techniques usable to train weights of a neural network modeled by the processing elements. Reverse checkpoint is usable to reduce memory usage during the training.

1 ACCELERATED DEEP LEARNING

2
3
4
5 [0001] [Intentionally blank]

6
7 BACKGROUND

8
9 [0002] Field: Advancements in accelerated deep learning are needed to provide
10 improvements in one or more of accuracy, performance, and energy efficiency.

11
12 [0003] Related Art: Unless expressly identified as being publicly or well known, mention
13 herein of techniques and concepts, including for context, definitions, or comparison purposes, should
14 not be construed as an admission that such techniques and concepts are previously publicly known or
15 otherwise part of the prior art.

16
17
18
19
20 SYNOPSIS

21
22 [0004] The invention may be implemented in numerous ways, e.g., as a process, an article of
23 manufacture, an apparatus, a system, a composition of matter, and a computer readable medium such
24 as a computer readable storage medium (e.g., media in an optical and/or magnetic mass storage device
25 such as a disk, an integrated circuit having non-volatile storage such as flash storage), or a computer
26 network wherein program instructions are sent over optical or electronic communication links. The
27 Detailed Description provides an exposition of one or more embodiments of the invention that enable
28 improvements in cost, profitability, performance, efficiency, and utility of use in the field identified
29 above. The Detailed Description includes an Introduction to facilitate understanding of the remainder
30 of the Detailed Description. The Introduction includes Example Embodiments of one or more of
31 systems, methods, articles of manufacture, and computer readable media in accordance with concepts
32 described herein. As is discussed in more detail in the Conclusions, the invention encompasses all
33 possible modifications and variations within the scope of the issued claims.

Brief Description of Drawings

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

[0005] Fig. 1 illustrates selected details of an embodiment of a system for neural network training and inference, using a deep learning accelerator.

[0006] Fig. 2 illustrates selected details of an embodiment of software elements associated with neural network training and inference, using a deep learning accelerator.

[0007] Fig. 3 illustrates selected details of an embodiment of processing associated with training a neural network and performing inference using the trained neural network, using a deep learning accelerator.

[0008] Fig. 4 illustrates selected details of an embodiment of a deep learning accelerator.

[0009] Fig. 5 illustrates selected details of an embodiment of a processing element of a deep learning accelerator.

[0010] Fig. 6 illustrates selected details of an embodiment of a router of a processing element.

[0011] Fig. 7 illustrates selected details of an embodiment of processing associated with a router of a processing element.

[0012] Fig. 8 illustrates selected details of an embodiment of a compute element of a processing element.

[0013] Fig. 9 illustrates selected details of an embodiment of processing a wavelet for task initiation.

[0014] Fig. 10 illustrates selected details of an embodiment of instruction processing associated with a compute element of a processing element.

[0015] Fig. 11 illustrates selected details of an embodiment of flow associated with dependency management via closeouts.

1
2 **[0016]** Fig. 12 illustrates selected details of an embodiment of flow associated with activation
3 accumulation and closeout, followed by partial sum computation and closeout.
4
5 **[0017]** Fig. 13A illustrates selected details of an embodiment of a sparse wavelet.
6
7 **[0018]** Fig. 13B illustrates selected details of an embodiment of a dense wavelet.
8
9 **[0019]** Fig. 14 illustrates selected details of an embodiment of creating and transmitting a
10 wavelet.
11
12 **[0020]** Fig. 15A illustrates selected details of an embodiment of receiving a wavelet.
13
14 **[0021]** Fig. 15B illustrates selected details of an embodiment of consuming a wavelet.
15
16 **[0022]** Fig. 16 illustrates selected details of an embodiment of block instruction and unblock
17 instruction execution.
18
19 **[0023]** Fig. 17 illustrates selected details of an embodiment of a neural network.
20
21 **[0024]** Fig. 18A illustrates selected details of a first embodiment of an allocation of
22 processing elements to neurons.
23
24 **[0025]** Fig. 18B illustrates selected details of a second embodiment of an allocation of
25 processing elements to neurons.
26
27 **[0026]** Fig. 19 illustrates selected details of an embodiment of smearing a neuron across a
28 plurality of processing elements.
29
30 **[0027]** Fig. 20 illustrates selected details of an embodiment of communication between
31 portions of split neurons.
32
33 **[0028]** Fig. 21A illustrates selected details of an embodiment of a Fabric Input Data Structure
34 Descriptor.

1
2 **[0029]** Fig. 21B illustrates selected details of an embodiment of a Fabric Output Data
3 Structure Descriptor.
4
5 **[0030]** Fig. 21C illustrates selected details of an embodiment of a 1D Memory Vector Data
6 Structure Descriptor.
7
8 **[0031]** Fig. 21D illustrates selected details of an embodiment of a 4D Memory Vector Data
9 Structure Descriptor.
10
11 **[0032]** Fig. 21E illustrates selected details of an embodiment of a Circular Memory Buffer
12 Data Structure Descriptor.
13
14 **[0033]** Fig. 22A illustrates selected details of an embodiment of a Circular Memory Buffer
15 Extended Data Structure Descriptor.
16
17 **[0034]** Fig. 22B illustrates selected details of an embodiment of a 4D Memory Vector
18 Extended Data Structure Descriptor.
19
20 **[0035]** Fig. 23 illustrates selected details of accessing operands in accordance with data
21 structure descriptors.
22
23 **[0036]** Fig. 24 illustrates selected details of an embodiment of decoding a data structure
24 descriptor.
25
26 **[0037]** Fig. 25A illustrates selected details of an embodiment of a multiple operand
27 instruction.
28
29 **[0038]** Fig. 25B illustrates selected details of an embodiment of a one source, no destination
30 operand instruction.
31
32 **[0039]** Fig. 25C illustrates selected details of an embodiment of an immediate instruction.
33

1 [0040] Fig. 26A illustrates an embodiment of a pipeline flow for Stochastic Gradient Descent
2 (SGD).
3
4 [0041] Fig. 26B illustrates an embodiment of a pipeline flow for Mini-Batch Gradient
5 Descent (MBGD).
6
7 [0042] Fig. 26C illustrates an embodiment of a pipeline flow for Continuous Propagation
8 Gradient Descent (CPGD).
9
10 [0043] Fig. 26D illustrates an embodiment of a pipeline flow for Continuous Propagation
11 Gradient Descent (CPGD) with Reverse CheckPoint (RCP).
12
13 [0044] Figs. 27A-27E illustrate various aspects of forward pass and backward pass
14 embodiments in accordance with SGD, MBGD, CPGD, and RCP processing.
15
16 [0045] Fig. 28A illustrates a generic operation of a matrix (m) multiplied by a vector (v).
17
18 [0046] Fig. 28B illustrates various representations of memory structures used in a forward
19 pass, a delta pass, and a chain pass.
20
21 [0047] Fig. 29 illustrates an embodiment of tasks as used in a forward pass state machine.
22

1
2
3

List of Reference Symbols in Drawings

[0048]

<u>Ref. Symbol</u>	<u>Element Name</u>
100	Neural Network System
110	Combined Server(s)
111	LAN
112	100Gb
113	Placements
114	Weights
115	Weights
120	Deep Learning Accelerator
121	FPGAs
122	PEs
123	Coupling
130	Autonomous Vehicle
131	CPUs
132	CRM
133	IEs
135	Camera
140	Cell Phone
141	CPUs
142	CRM
143	IEs
145	Camera
150	Placement Server(s)
151	CPUs
152	CRM
160	Connection Server(s)
161	CPUs
162	CRM
164	NICs
180	Internet
200	Neural Network Software
210	Placement Server(s) SW
212	Neuron to PE Mapping SW
220	Connection Server(s) SW
224	100Gb NIC Driver
225	Training Info Provider SW
226	Weight Receiver SW
230	Autonomous Vehicle SW
232	Video Camera SW
233	Inference Engine(s) SW
234	Navigating SW
240	Cell Phone SW
242	Still Camera SW
243	Inference Engine(s) SW
244	Posting SW

<u>Ref. Symbol</u>	<u>Element Name</u>
250	Misc SW on FPGAs
260	Task SW on PEs
300	Neural Network Training/Inference, Overall
310	Place Neurons
320	Initialize FPGAs
330	Initialize PEs
340	Training Data => PEs
350	Forward Pass, Delta Pass, Chain Pass, Update Weights
360	Training Complete?
370	Weights Out
380	Use Weights for Inference
400	Deep Learning Accelerator
401	Forward
402	Delta
403	Chain
410	ASIC
411	ASIC
412	Wafer
420	I/O FPGAs
430	North coupling
431	East coupling
432	South coupling
433	West coupling
497	Particular PE
498	Particular PE
499	PE
500	PE
510	Router
511	West
512	Skip West
513	North
514	Skip East
515	East
516	South
520	Compute Element
521	Off Ramp
522	On Ramp
600	Router
610	Data In
611	skipX+
612	skipX-
613	X+
614	X-
615	Y+
616	Y-
617	On Ramp
620	Data Out
621	skipX+

<u>Ref. Symbol</u>	<u>Element Name</u>
622	skipX-
623	X+
624	X-
625	Y+
626	Y-
627	Off Ramp
630	Stall Out
631	skipX+
632	skipX-
633	X+
634	X-
635	Y+
636	Y-
637	On Ramp
640	Sources
641	skipX+
642	skipX-
643	X+
644	X-
645	Y+
646	Y-
647	Off Ramp
650	Data Queues
651	Write Dec
652	Out
653	Sources
654	Router Sched
656	Gen Stall
657	Stall
660	Control Info
661	Dest
662	Sent
670	Src
710	Wavelet Ingress
711	Wait for Wavelet
712	Receive Wavelet
713	Wavelet=> Router Q
720	Stall Info
721	Router Q Full?
722	DeAssert Stall
723	Assert Stall
730	Wavelet Egress
731	Q Empty?
732	Choose?
733	Stalled?
734	Send Wavelet
800	CE
812	Terminate

<u>Ref. Symbol</u>	<u>Element Name</u>
820	Off Ramp
822	Hash
824	Qdistr
830	Picker
834	PC
836	I-Seq
840	Dec
842	RF
844	D-Seq
846	DSRs
848	D-Store
852	Data Path
854	Memory
860	On Ramp
890	Base
896	Scheduling Info
897	Qs
897.0	Q0
897.N	QN
898	Active Bits
898.0	Active Bit 0
898.N	Active Bit N
899	Block Bits
899.0	Block Bit 0
899.N	Block Bit N
900	Processing a Wavelet for Task Initiation, Overall
901	Start
905	Select Ready Wavelet for Task Initiation
908	Control/Data?
920	Add (Color * 4) to Base Register to Form Instruction Address
930	Add Lower Index Bits to Base Register to Form Instruction Address
950	Fetch Instructions From Memory at Instruction Address
960	Execute Fetched Instruction(s)
961	Not Terminate
962	Terminate
990	End
1000	Instruction Processing, Overall
1010	Check Control Inputs
1012	Branch Stall?
1014	Do Nothing
1016	Terminate => Scheduler
1020	EX Branch Resolution?
1022	D-Seq Stall?
1024	I-Seq Mode?
1026	Fetch Instr
1028	Terminate Instr?
1030	Branch Instr?

<u>Ref. Symbol</u>	<u>Element Name</u>
1032	Update PC Instr => Decode
1040	Process Next Task/Branch PC Task Addr => PC
1042	Stall Sequencer
1100	Dependency Management, Overall
1101	Activations From Prior Layer
1102	Receive and Accumulate Activations
1110	Closeout From Prior Layer
1111	Receive Activation Closeout
1112	Start Partial Sums
1113	Calculate Partial Sum
1114	Propagate Partial Sums
1120	Transmit Activations
1121	Activations to Next Layer
1122	Closeout to Next Layer
1123	Reschedule
1131	Flow Control Dependency
1132	Output Wavelet to Different PE
1133	Wake Wavelet to Self
1200	Activation Accumulation/Closeout and Partial Sum Computation/Closeout, Overall
1201	Start
1202	Receive Activation
1203	Accumulate Activations
1204	Receive Activation Closeout
1205	Start Partial Sum Ring
1206	Receive Partial Sum
1207	Compute Partial Sum
1208	Transmit Partial Sum
1209	Transmit Activations
1210	Transmit Closeout
1211	End
1301	Sparse Wavelet
1302	Sparse Wavelet Payload
1320	Control Bit
1321	Index
1321.1	Lower Index Bits
1321.2	Upper Index Bits
1322	Sparse Data
1324	Color
1331	Dense Wavelet
1332	Dense Wavelet Payload
1340	Control Bit
1343.1	Dense Data
1343.2	Dense Data
1344	Color
1400	Wavelet Creation Flow, Overall
1401	Start
1402	Initialize PEs

<u>Ref. Symbol</u>	<u>Element Name</u>
1403	Set Source
1404	Set Destination (Fabric) DSR
1404.5	Fetch/Decode Instruction with Destination DSR
1404.6	Read DSR(s)
1405	Read (Next) Source Data Element(s) from Queue/Memory
1406	Provide Data Element(s) as Wavelet to Router
1407	More Data Elements?
1408	Transmit Wavelet(s) to Fabric
1409	Receive Wavelet(s) from Fabric
1410	End
1420	CE of Transmitting PE
1430	Router of Transmitting PE
1440	Router of Receiving PE
1500	Wavelet Receive Flow, Overall
1501	Start
1502	Initialize PEs
1503	Receive Wavelet at Router
1504	To Other PE(s)?
1505	Transmit Wavelet to Output(s)
1506	For Local CE?
1507	Write Wavelet to Picker Queue
1510	End
1520	Router of Receiving PE
1530	CE of Receiving PE
1550	Wavelet Consumption Flow, Overall
1551	Start
1552	Picker Selects Wavelet for Processing
1553	Fetch, Execute Instructions
1554	End
1600	Block and Unblock Instruction Processing Flow, Overall
1601	Start
1602	Fetch, Decode Instruction
1603	Block Instruction?
1604	Block Color(s)
1610	Unblock Instruction?
1611	Unblock Color(s)
1620	Execute Instruction
1630	End
1700	Neural Network
1710	Input Layer
1711	N11
1712	N12
1713	N13
1720	Internal Layers
1721	N21
1721.1, 1721.2	1/2 N21 portions, respectively
1722	N22
1722.1, 1722.2	1/2 N22 portions, respectively

<u>Ref. Symbol</u>	<u>Element Name</u>
1723	N23
1723.1, 1723.2	1/2 N23 portions, respectively
1724	N24
1724.1, 1724.2	1/2 N24 portions, respectively
1731	N31
1731.1, 1731.2, 1731.3, 1731.4	1/4 N31 portions, respectively
1732	N32
1732.1, 1732.2, 1732.3, 1732.4	1/4 N32 portions, respectively
1733	N33
1740	Output Layer
1741	N41
1742	N42
1791	communication
1791.1	communication portion
1792	communication
1792.1	communication portion
1793	communication
1793.1	communication portion
1820	PE0
1821	PE1
1822	PE2
1823	PE3
1824	PE4
1825	PE5
1910	in0
1911	in1
1912	in2
1913	in3
1914	in4
1915	in5
1920	out0
1921	out1
1922	out2
1923	out3
1924	out4
1925	out5
1930.1	1/2 Local Compute
1930.2	1/2 Local Compute
1940.1	1/2 Local Storage
1940.2	1/2 Local Storage
1950.1	Additional Compute
1950.2	Additional Compute
1960.1	Additional Storage
1960.2	Additional Storage
1970	Additional Communication
2000	Wafer Portion
2040, 2041, 2043, 2044	coupling between adjacent PEs, respectively
2050, 2051, 2052, 2053, 2054,	portion of coupling between adjacent PEs, respectively

<u>Ref. Symbol</u>	<u>Element Name</u>
2055, 2056, 2057	
2060	communication
2100	Fabric Input Data Structure Descriptor
2101	Length
2102	UTID (Microthread Identifier)
2103	UE (Microthread Enable)
2104	SW (SIMD Width)
2105	AC (Activate Color)
2106	Term (Terminate Microthread on Control Wavelet)
2107	CX (Control Wavelet Transform Enable)
2108	US (Microthread Sparse Mode)
2109	Type
2110	SS (Single Step)
2111	SA (Save Address / Conditional Single Step Mode)
2112	SC (Color Specified, Normal Mode)
2113	SQ (Queue Specified, Normal Mode)
2114	CH (Color, High Bits)
2120	Fabric Output Data Structure Descriptor
2121	Length
2122	UTID (Microthread Identifier)
2123	UE (Microthread Enable)
2124	SW (SIMD Width)
2125	AC (Activate Color)
2126	Color
2127	C (Output Control Bit)
2128.1	Index Low
2128.2	Index High
2129	Type
2130	SS (Single Step)
2131	SA (Save Address / Conditional Single Step Mode)
2132	WLI (Wavelet Index Select)
2140	1D Memory Data Structure Descriptor
2141	Length
2142	Base Address
2149	Type
2150	SS (Single Step)
2151	SA (Save Address / Conditional Single Step Mode)
2152	WLI (Wavelet Index Select)
2153	Stride
2160	4D Memory Data Structure Descriptor
2161	Length
2161.1	Length Lower Bits
2161.2	Length Upper Bits
2162	Base Address
2169	Type
2170	SS (Single Step)
2171	SA (Save Address / Conditional Single Step Mode)
2172	WLI (Wavelet Index Select)

<u>Ref. Symbol</u>	<u>Element Name</u>
2180	Circular Memory Buffer Data Structure Descriptor
2181	Length
2182	Base Address
2184	SW (SIMD Width)
2188	FW (FIFO Wrap Bit)
2189	Type
2190	SS (Single Step)
2191	SA (Save Address / Conditional Single Step Mode)
2192	WLI (Wavelet Index Select)
2210	Circular Memory Buffer Extended Data Structure Descriptor
2211	Type
2212	Start Address
2213	End Address
2214	FIFO
2215	Push (Activate) Color
2216	Pop (Activate) Color
2240	4D Memory Vector Extended Data Structure Descriptor
2241	Type
2242	Dimensions
2243	DF (Dimension Format)
2244.1	Stride Select (for Dimension) 1
2244.2	Stride Select (for Dimension) 2
2244.3	Stride Select (for Dimension) 3
2244.4	Stride Select (for Dimension) 4
2245	Stride
2300	Data Structure Descriptor Flow, Overall
2301	Start
2302	Set DSR(s)
2303	Fetch/Decode Instruction with DSR(s)
2304	Read DSR(s)
2305	(optional) Set XDSR(s)
2306	(optional) Read XDSR(s)
2310	Read (Next) Source Data Element(s) from Queue/Memory
2311	Perform (Next) Operation(s) on Data Element(s)
2312	Write (Next) Destination Data Element(s) to Queue/Memory
2313	More Data Element(s)?
2316	End
2400	Data Structure Descriptor Decode Flow, Overall
2401	Start
2410	Fabric Vector
2411	Type = Fabric?
2412	Access via DSD
2420	Memory Vector
2421	Type = XDSR?
2422	Read XDSR Specified via DSD
2423	Type = 4D Vector?
2424	(optional) Read Stride Register(s)
2427	Access 1D via DSD

<u>Ref. Symbol</u>	<u>Element Name</u>
2428	Access 4D via XDSD
2429	Access Circular Buffer via XDSD
2499	End
2510	Multiple Operand Instruction
2511	Instruction Type
2512	Opcode
2513	Operand 0 Encoding
2513.1	Operand 0 Type
2513.2	Operand 0
2514	Operand 1 Encoding
2514.1	Operand 1 Type
2514.2	Operand 1
2515	Terminate
2520	One Source, No Destination Operand Instruction
2521	Instruction Type
2522	Opcode
2523	Operand 1 Encoding
2523.1	Operand 1 Type
2523.2	Operand 1
2524	Immediate
2525	Terminate
2530	Immediate Instruction
2531	Instruction Type
2532	Opcode
2533.2	Operand 0
2534.1	Immediate Low
2534.2	Immediate High
2534	Immediate
2611	First Forward Pass
2612	Second Forward Pass
2621	First Backward Pass
2622	Second Backward Pass
2631	Mini-Batch Size (N)
2632	Overhead
2633	Update Interval (U)
2651	Forward Pass
2661	Backward Pass
2665	Forward Pass
2666	Backward Pass
2667	Weight Update Use
2671	Forward Pass
2681	Backward Pass
2685	Activation Storage
2686	Recomputed Activation Storage
2701	Previous Layer
2702	Subsequent Layer
2703	Previous Layer
2704	Subsequent Layer

<u>Ref. Symbol</u>	<u>Element Name</u>
2710	Compute
2711	F
2712	B
2715	Storage
2716	A
2717	W
2718	W
2720	Compute
2721	F
2722	B
2725	Storage
2726	A
2727	W
2728	W
2729	A
2730	Compute
2735	Storage
2740	Compute
2745	Storage
2781	$A_{1,t}$
2782	$A_{2,t}$
2783	$A_{3,t}$
2784	$A'_{2,t}$
2791	$\Delta_{1,t}$
2792	$\Delta_{2,t}$
2793	$\Delta_{3,t}$
2794	$\Delta'_{1,t}$
2795	$\Delta'_{2,t}$
2796	$\Delta'_{3,t}$
2901	f rxact:acc
2902	f rxact:close
2903	f psum:prop
2904	f txact:tx
2911	Activations
2912	Closeouts
2913	Flow
2914	Wake
2915	Reschedule
2916	Start Psums
2921	Activations
2922	Closeouts
2930	Prop Psums
2931	Prop Psums

1

1
2
3 **DETAILED DESCRIPTION**

4 **[0049]** A detailed description of one or more embodiments of the invention is provided
5 below along with accompanying figures illustrating selected details of the invention. The invention is
6 described in connection with the embodiments. The embodiments herein are understood to be merely
7 exemplary, the invention is expressly not limited to or by any or all of the embodiments herein, and
8 the invention encompasses numerous alternatives, modifications, and equivalents. To avoid monotony
9 in the exposition, a variety of word labels (such as: first, last, certain, various, further, other,
10 particular, select, some, and notable) may be applied to separate sets of embodiments; as used herein
11 such labels are expressly not meant to convey quality, or any form of preference or prejudice, but
12 merely to conveniently distinguish among the separate sets. The order of some operations of disclosed
13 processes is alterable within the scope of the invention. Wherever multiple embodiments serve to
14 describe variations in process, system, and/or program instruction features, other embodiments are
15 contemplated that in accordance with a predetermined or a dynamically determined criterion perform
16 static and/or dynamic selection of one of a plurality of modes of operation corresponding respectively
17 to a plurality of the multiple embodiments. Numerous specific details are set forth in the following
18 description to provide a thorough understanding of the invention. The details are provided for the
19 purpose of example and the invention may be practiced according to the claims without some or all of
20 the details. For the purpose of clarity, technical material that is known in the technical fields related to
21 the invention has not been described in detail so that the invention is not unnecessarily obscured.

22
23 **INTRODUCTION**

24
25 **[0050]** This introduction is included only to facilitate the more rapid understanding of the
26 Detailed Description; the invention is not limited to the concepts presented in the introduction
27 (including explicit examples, if any), as the paragraphs of any introduction are necessarily an abridged
28 view of the entire subject and are not meant to be an exhaustive or restrictive description. For
29 example, the introduction that follows provides overview information limited by space and
30 organization to only certain embodiments. There are many other embodiments, including those to
31 which claims will ultimately be drawn, discussed throughout the balance of the specification.

32
33 **[0051]** In an aspect conceptually related to continuous propagation for accelerated deep
34 learning, techniques in advanced deep learning provide improvements in one or more of accuracy,

1 performance, and energy efficiency, such as accuracy of learning, accuracy of prediction, speed of
2 learning, performance of learning, and energy efficiency of learning. An array of processing elements
3 performs flow-based computations on wavelets of data. Each processing element has a respective
4 compute element and a respective routing element. Each compute element has processing resources
5 and memory resources. Each router enables communication via wavelets with at least nearest
6 neighbors in a 2D mesh. Stochastic gradient descent, mini-batch gradient descent, and continuous
7 propagation gradient descent are techniques usable to train weights of a neural network modeled by
8 the processing elements. Reverse checkpoint is usable to reduce memory usage during the training.
9

10 **[0052]** In an aspect conceptually related to fabric vectors for accelerated deep learning,
11 techniques in advanced deep learning provide improvements in one or more of accuracy, performance,
12 and energy efficiency. An array of processing elements performs flow-based computations on
13 wavelets of data. Each processing element has a respective compute element and a respective routing
14 element. Each compute element has memory. Each router enables communication via wavelets with
15 at least nearest neighbors in a 2D mesh. Routing is controlled by respective virtual channel specifiers
16 in each wavelet and routing configuration information in each router. Instructions executed by the
17 compute element include one or more operand specifiers, some of which specify a data structure
18 register storing a data structure descriptor. The data structure descriptor describes an operand as a
19 fabric vector or a memory vector. The data structure descriptor further describes the length of the
20 fabric vector, whether the fabric vector is eligible for microthreading, and a number of data elements
21 of the fabric vector to receive, transmit, and/or process in parallel. The data structure descriptor
22 further specifies virtual channel and task identification information relating to processing the fabric
23 vector, whether to terminate upon receiving a control wavelet, and whether to mark an outgoing
24 wavelet as a control wavelet.
25

26 **[0053]** In an aspect conceptually related to data structure descriptors for accelerated deep
27 learning, techniques in advanced deep learning provide improvements in one or more of accuracy,
28 performance, and energy efficiency. An array of processing elements performs flow-based
29 computations on wavelets of data. Each processing element has a respective compute element and a
30 respective routing element. Each compute element has memory. Each router enables communication
31 via wavelets with at least nearest neighbors in a 2D mesh. Routing is controlled by respective virtual
32 channel specifiers in each wavelet and routing configuration information in each router. Instructions
33 executed by the compute element include one or more operand specifiers, some of which specify a
34 data structure register storing a data structure descriptor. The data structure descriptor describes an

1 operand as a fabric vector or a memory vector. The data structure descriptor further describes the
2 memory vector as one of a one-dimensional vector, a four-dimensional vector, or a circular buffer
3 vector. Optionally, the data structure descriptor specifies an extended data structure register storing an
4 extended data structure descriptor. The extended data structure descriptor specifies parameters
5 relating to a four-dimensional vector or a circular buffer vector.
6

7 **[0054]** In an aspect conceptually related to neuron smearing for accelerated deep learning,
8 techniques in advanced deep learning provide improvements in one or more of accuracy, performance,
9 and energy efficiency. An array of processing elements performs flow-based computations on
10 wavelets of data. Each processing element has a respective compute element and a respective routing
11 element. Each compute element has memory. Each router enables communication via wavelets with
12 at least nearest neighbors in a 2D mesh. Routing is controlled by respective virtual channel specifiers
13 in each wavelet and routing configuration information in each router. At least a first single neuron is
14 implemented using resources of a plurality of the array of processing elements. At least a portion of a
15 second neuron is implemented using resources of one or more of the plurality of processing elements.
16 In some usage scenarios, the foregoing neuron implementation enables greater performance by
17 enabling a single neuron to use the computational resources of multiple processing elements and/or
18 computational load balancing across the processing elements while maintaining locality of incoming
19 activations for the processing elements.
20

21 **[0055]** In an aspect conceptually related to task synchronization for accelerated deep
22 learning, techniques in advanced deep learning provide improvements in one or more of accuracy,
23 performance, and energy efficiency. An array of processing elements performs flow-based
24 computations on wavelets of data. Each processing element has a respective compute element and a
25 respective routing element. Each compute element has memory. Each router enables communication
26 via wavelets with at least nearest neighbors in a 2D mesh. Routing is controlled by respective virtual
27 channel specifiers in each wavelet and routing configuration information in each router. A particular
28 one of the compute elements conditionally selects for task initiation a previously received wavelet
29 specifying a particular one of the virtual channels. The conditional selecting excludes the previously
30 received wavelet for selection until at least block/unblock state maintained for the particular virtual
31 channel is in an unblock state. The compute elements execute block/unblock instructions to modify
32 the block/unblock state.
33

1 **[0056]** In an aspect conceptually related to dataflow triggered tasks for accelerated deep
2 learning, techniques in advanced deep learning provide improvements in one or more of accuracy,
3 performance, and energy efficiency. An array of processing elements performs flow-based
4 computations on wavelets of data. Each processing element has a respective compute element and a
5 respective routing element. Each compute element has memory. Each router enables communication
6 via wavelets with at least nearest neighbors in a 2D mesh. Routing is controlled by respective virtual
7 channel specifiers in each wavelet and routing configuration information in each router. A particular
8 one of the compute elements receives a particular wavelet comprising a particular virtual channel
9 specifier and a particular data element. Instructions are read from the memory of the particular
10 compute element based at least in part on the particular virtual channel specifier. The particular data
11 element is used as an input operand to execute at least one of the instructions.

12

13 **[0057]** In an aspect conceptually related to control wavelet for accelerated deep learning,
14 techniques in advanced deep learning provide improvements in one or more of accuracy, performance,
15 and energy efficiency. An array of processing elements performs flow-based computations on
16 wavelets of data. Each processing element has a respective compute element and a respective routing
17 element. Each compute element has a memory. Each router enables communication via wavelets
18 with at least nearest neighbors in a 2D mesh. A particular one of the compute elements receives a
19 wavelet. If a control specifier of the wavelet is a first value, then instructions are read from the
20 memory of the particular compute element in accordance with an index specifier of the wavelet. If the
21 control specifier is a second value, then instructions are read from the memory of the particular
22 compute element in accordance with a virtual channel specifier of the wavelet. Then the particular
23 compute element initiates execution of the instructions.

24

25 **[0058]** In an aspect conceptually related to wavelet representation for accelerated deep
26 learning, techniques in advanced deep learning provide improvements in one or more of accuracy,
27 performance, and energy efficiency. An array of processing elements performs flow-based
28 computations on wavelets of data. Each processing element has a respective compute element and a
29 respective routing element. Each compute element has dedicated storage. Each router enables
30 communication with at least nearest neighbors in a 2D mesh. The communication is via wavelets in
31 accordance with a representation comprising an index specifier, a virtual channel specifier, an index
32 specifier, a data element specifier, and an optional control/data specifier. The virtual channel specifier
33 and the index specifier are associated with one or more instructions. The index specifier is associated

1 with at least a first instruction operand of the one or more instructions. The data element is associated
2 with at least a second instruction operand of the one or more instructions.

3

4 **[0059]** A first example of accelerated deep learning is using a deep learning accelerator to
5 train a neural network. A second example of accelerated deep learning is using a deep learning
6 accelerator to operate a trained neural network to perform inferences. A third example of accelerated
7 deep learning is using a deep learning accelerator to train a neural network and subsequently perform
8 inference with any one or more of the trained neural network, information from same, and a variant of
9 same.

10

11 **[0060]** Examples of neural networks include Fully Connected Neural Networks (FCNNs),
12 Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), Long Short-Term
13 Memory (LSTM) networks, autoencoders, deep belief networks, and generative adversarial networks.

14

15 **[0061]** An example of training a neural network is determining one or more weights
16 associated with the neural network, such as by hardware acceleration via a deep learning accelerator.
17 An example of making an inference is using a trained neural network to compute results by processing
18 input data based on weights associated with the trained neural network.

19

20 **[0062]** A neural network processes data according to a dataflow graph comprising layers of
21 neurons. Stimuli (e.g., input data) is received by an input layer of neurons and the computed results of
22 the data flow graph (e.g., output data) are provided by an output layer of neurons. Example layers of
23 neurons include input layers, output layers, rectified linear unit layers, fully connected layers,
24 recurrent layers, long short-term memory layers, convolutional layers, kernel layers, dropout layers,
25 and pooling layers. A neural network is conditionally and/or selectively trained, subject to hardware
26 acceleration. After being trained, a neural network is conditionally and/or selectively used for
27 inference, subject to hardware acceleration.

28

29 **[0063]** An example of a deep learning accelerator is one or more relatively specialized
30 hardware elements operating in conjunction with one or more software elements to train a neural
31 network and/or perform inference with a neural network relatively more efficiently than using
32 relatively less specialized hardware elements. Some implementations of the relatively specialized
33 hardware elements include one or more hardware logic circuitry elements such as transistors, resistors,
34 inductors, capacitors, wire interconnects, combinatorial logic (e.g., NAND, NOR) gates, latches,

1 register files, memory arrays, tags for memory arrays, content-addressable memories, flash, ROM,
2 DRAM, SRAM, Serializer/Deserializer (SerDes), I/O drivers, and the like, such as implemented via
3 custom logic, synthesized logic, ASICs, and/or FPGAs. Some of the relatively less specialized
4 hardware elements include conventional CPUs and conventional GPUs.

5

6 **[0064]** An example implementation of a deep learning accelerator is enabled to process
7 dataflow in accordance with computations performed for training of a neural network and/or inference
8 with a neural network. Some deep learning accelerators comprise processing elements coupled via a
9 fabric and enabled to communicate with each other via the fabric. Sometimes the processing elements
10 and the fabric are collectively referred to as a fabric of processing elements.

11

12 **[0065]** An example implementation of a processing element is enabled to communicate and
13 process wavelets. In various circumstances, the wavelets correspond to dataflow and/or instruction
14 flow in accordance with communication and/or processing enabling computations performed for
15 training of and/or inference using a neural network.

16

17 **[0066]** An example processing element comprises a router to communicate wavelets via the
18 fabric and a compute element to process the wavelets. An example router is coupled to a plurality of
19 elements: a fabric, an off ramp to the compute element, and an on ramp from the compute element.
20 An example coupling between the router and the fabric enables communication between the router
21 and, e.g., four logically and/or physically adjacent processing elements. The router variously receives
22 wavelets from the fabric and the on ramp. The router variously transmits wavelets to the fabric and
23 the off ramp.

24

25 **[0067]** An example implementation of a compute element is enabled to process wavelets by
26 initiating tasks and executing instructions associated with the wavelets, and accessing data associated
27 with the wavelets and/or the instructions. The instructions are in accordance with an instruction set
28 architecture comprising arithmetic instructions, control flow instructions, datatype conversion
29 instructions, configuration instructions, fabric management instructions, and load/store instructions.
30 The instructions operate on operands comprising various datatypes, e.g., integer datatypes and
31 floating-point datatypes of various widths. The operands variously comprise scalar operands and
32 vector operands. In various embodiments and/or usage scenarios, a vector variously represents
33 weights of a neural network, inputs or stimuli of a neural network, activations of a neural network,
34 and/or partial sums of a neural network. In some scenarios, a vector is a sparse vector (e.g., a vector

1 of neuron activations) and comprises sparse data elements (e.g., only non-zero elements). In some
2 other scenarios, a vector is a dense vector (e.g., pixel values) and comprises dense data elements (e.g.,
3 all elements of the vector, including zero elements).

4
5 **[0068]** An example compute element comprises hardware elements that collectively execute
6 the instructions associated with a wavelet by performing operations specified by the instructions (e.g.,
7 arithmetic operations, control flow operations, and load/store operations). Examples of the hardware
8 elements include picker queues, a picker, a task definition table, an instruction sequencer, an
9 instruction decoder, a data sequencer, a register file, a memory, a pseudo-random number generator,
10 and an ALU. Some implementations of the hardware elements are in accordance with hardware logic
11 circuitry elements as described elsewhere herein. Sometimes a compute element is referred to as a
12 compute engine. Sometimes the compute scheduler is referred to as a picker and the compute
13 scheduler queues are referred to as picker queues.

14
15 **[0069]** An example fabric is a collection of logical and/or physical couplings between
16 processing elements and/or within a single processing element. The fabric is usable to implement
17 logical and/or physical communication topologies such as a mesh, a 2D mesh, a 3D mesh, a
18 hypercube, a torus, a ring, a tree, or any combination thereof. An example of a physical coupling
19 between processing elements is a set of physical interconnects (comprising optional and/or selective
20 buffering) between physically-coupled processing elements. A first example of physically-coupled
21 processing elements is immediately physically adjacent processing elements, such as a first processing
22 element located directly beside (such as 'north', 'south', 'east', or 'west') of a second processing
23 element. A second example of physically-coupled processing elements is relatively physically nearby
24 processing elements, such as a first processing element located within a relatively small number of
25 intervening processing elements, e.g., one or two 'rows' and/or 'columns' away from a second
26 processing element. A third example of physically-coupled processing elements is relatively
27 physically far away processing elements, such as a first processing element located physical relatively
28 far away from a second processing element, such as a distance limited by signal propagation (with or
29 without optional and/or selective buffering) within a clock cycle and/or clock sub-cycle associated
30 with the processing elements. An example of physical coupling within a single processing element
31 (having, e.g., a compute element and a router) is an on ramp coupling output information from the
32 compute element to the router, and an off ramp coupling input information from the router to the
33 compute element. In some situations, the router routes information from the on ramp to the off ramp.

34

1 [0070] An example of a logical coupling between processing elements is a virtual channel as
2 implemented by routers within processing elements. A route between a first processing element and a
3 second processing element is implemented, e.g., by routers within processing elements along the route
4 forwarding in accordance with the virtual channel and routing configuration information. An example
5 of a logical coupling within a single particular processing element (having, e.g., a router) is a virtual
6 channel as implemented by the router, enabling the particular processing element to send information
7 via the virtual channel to the particular processing element. The router forwards “internally” with
8 respect to the particular processing element in accordance with the virtual channel and routing
9 configuration information.

10

11 [0071] An example wavelet is a bundle of information communicated between processing
12 elements via the fabric. An example wavelet comprises a wavelet payload and a color. A wavelet
13 payload comprises data and is associated with instructions. A first response to a wavelet received by a
14 compute element of a processing element comprises the compute element initiating a task, such as
15 corresponding to processing of instructions associated with the wavelet. A second response to a
16 wavelet received by a compute element of a processing element comprises the compute element
17 processing data of the wavelet. Example types of wavelets include dense wavelets and sparse
18 wavelets, as well as data wavelets and control wavelets.

19

20 [0072] Wavelets are used, for example, for communicating between processing elements. In
21 a first scenario, a first processing element transmits wavelets to a second processing element. In a
22 second scenario, an external device (e.g., an FPGA) transmits wavelets to a processing element. In a
23 third scenario, a processing element transmits wavelets to an external device (e.g., an FPGA).

24

25 [0073] An example virtual channel is one or more communication pathways specified by a
26 color and enabled, e.g., by a fabric and one or more routers. A wavelet comprising a particular color
27 is sometimes referred to as being associated with a particular virtual channel associated with the
28 particular color. A first example of a color is a fabric color specifying a virtual channel between two
29 different processing elements. In some embodiments, a fabric color is a 5-bit integer. A second
30 example of a color is a local color specifying a virtual channel from a processing element to the
31 processing element. In some embodiments, a color is a 6-bit integer and specifies one of a fabric color
32 and a local color.

33

1 [0074] An example task comprises a collection of instructions executed in response to a
2 wavelet. An example instruction comprises an operation and optionally one or more operands
3 specifying locations of data elements to be processed in accordance with the operation. A first
4 example of an operand specifies data elements in memory. A second example of an operand specifies
5 data elements communicated (e.g., received or transmitted) via the fabric. An example of a data
6 sequencer determines the locations of data elements. An example of an instruction sequencer
7 determines an address in memory of instructions associated with a wavelet.

8
9 [0075] An example picker queue is enabled to hold wavelets received via an off ramp of the
10 fabric for processing in the compute element. An example of a picker selects a wavelet from the
11 picker queue for processing.

12
13 [0076] An example of an Integrated Circuit (IC) is a collection of circuitry implemented on a
14 single portion of semiconductor material. An example of an Application-Specific Integrated Circuit
15 (ASIC) is an IC designed for a particular use. An example of wafer-scale integration is implementing
16 a system using all or a significant portion of a wafer as an element of the system, e.g., by leaving the
17 wafer whole or substantially whole.

18
19 [0077] In some embodiments and/or usage scenarios, wafer-scale integration enables
20 connecting multiple elements in a system via wafer interconnect formed using silicon fabrication
21 processes instead of via inter-chip interconnect, and thus improves any one or more of improved
22 performance, cost, reliability, and energy efficiency. As a specific example, a system implemented
23 using wafer-scale integration technology enables implementation of three million PEs on a single
24 wafer, each of the PEs having bandwidth to nearest physical neighbors that is greater than a
25 comparable system using other-than wafer-scale integration technology. The greater bandwidth
26 enables the system implemented using wafer-scale integration technology to relatively efficiently train
27 and/or perform inferences for larger neural networks than the system implemented using other-than
28 wafer-scale integration technology.

29
30
31 Acronyms

32
33 [0078] At least some of the various shorthand abbreviations (e.g., acronyms) defined here
34 refer to certain elements used herein.

<u>Acronym</u>	<u>Description</u>
ASIC	Application Specific Integrated Circuit
CE	Compute Element
CNN	Convolutional Neural Network
CPGD	Continuous Propagation Gradient Descent
CPU	Central Processing Unit
CRM	Computer Readable Media
DSD	Data Structure Descriptor
DSP	Digital Signal Processor
DSR	Data Structure Register
FCNN	Fully Connected Neural Network
FPGA	Field-Programmable Gate Array
GPU	Graphics Processing Unit
HPC	High-Performance Computing
HW	HardWare
IC	Integrated Circuit
IE	Inference Engine
LFSR	Linear Feedback Shift Register
LSB	Least Significant Bit
LSTM	Long Short-Term Memory
MBGD	Mini-Batch Gradient Descent
ML	Machine Learning
MSB	Most Significant Bit
PE	Processing Element
PRNG	Pseudo Random Number Generator
RNN	Recurrent Neural Network
RCP	Reverse CheckPoint
SGD	Stochastic Gradient Descent
SW	SoftWare
XDSD	eXtended Data Structure Descriptor
XDSR	eXtended Data Structure Register

1

2

1 EXAMPLE EMBODIMENTS

2
3 **[0079]** In concluding the introduction to the detailed description, what follows is a collection
4 of example embodiments, including at least some explicitly enumerated as “ECs” (Example
5 Combinations), providing additional description of a variety of embodiment types in accordance with
6 the concepts described herein; these examples are not meant to be mutually exclusive, exhaustive, or
7 restrictive; and the invention is not limited to these example embodiments but rather encompasses all
8 possible modifications and variations within the scope of the issued claims and their equivalents.
9

10 **[0080]** EC100) A system comprising:
11 a fabric of processor elements, each processor element comprising a fabric router and a
12 compute engine enabled to perform dataflow-based and instruction-based processing;
13 wherein each processor element selectively communicates fabric packets with others of the
14 processor elements; and
15 wherein each compute engine selectively performs the processing in accordance with a virtual
16 channel specifier and a task specifier of each fabric packet the compute engine
17 receives.
18

19 **[0081]** EC100b) A system comprising:
20 a fabric of processor elements, each processor element comprising a fabric router and a
21 compute engine;
22 wherein each processor element selectively communicates fabric packets with others of the
23 processor elements; and
24 wherein each compute engine selectively performs dataflow processing and instruction
25 processing respectively in accordance with a dataflow field and an instruction field of
26 each fabric packet the compute engine receives.
27

28 **[0082]** EC100c) The system of EC100, wherein the processing is in accordance with a data-
29 flow graph.
30

31 **[0083]** EC100d) The system of EC100, wherein a workload is executed comprising
32 predominantly dataflow-based processing with minimal instruction-based processing.
33

1 [0084] EC100e) The system of EC100d, wherein the system implements a Long Short Term
2 Memory (LSTM) neural network model.

3

4 [0085] EC100f) The system of EC100, wherein a workload is executed comprising
5 predominantly instruction-based processing with minimal dataflow-based processing.

6

7 [0086] EC100g) The system of EC100, wherein the system is implemented at least in part
8 using wafer-scale integration.

9

10 [0087] EC100h) The system of EC100, wherein the fabric of processor elements is
11 implemented at least in part using VLSI fabrication.

12

13 [0088] EC101) The system of EC100, wherein the virtual channel specifier selects
14 independent respective routing paths in the fabric.

15

16 [0089] EC101b) The system of EC100, wherein the virtual channel specifier selects routing
17 paths in the fabric to perform multicast.

18

19 [0090] EC101c) The system of EC100, wherein the virtual channel specifier selects routing
20 paths in the fabric to perform load splitting.

21

22 [0091] EC102) The system of EC100, wherein the task specifier selects one or more
23 operations to perform.

24

25 [0092] EC103) The system of EC100, wherein the fabric comprises a 2D array of the
26 processor elements.

27

28 [0093] EC103b) The system of EC100, wherein the fabric comprises a processor element
29 interconnection topology selected from the group consisting of fully connected, star, ring, array, mesh,
30 hypercube, torus, and tree.

31

- 1 **[0094]** EC103c) The system of EC100, wherein the fabric comprises a processor element
2 interconnection topology dimension selected from the group consisting of 1D, 2D, 3D, and a
3 dimension greater than 3D.
4
- 5 **[0095]** EC104) The system of EC100, wherein the system is enabled to execute machine
6 learning workloads.
7
- 8 **[0096]** EC105) The system of EC100, wherein the system is trained to perform an inference
9 application.
10
- 11 **[0097]** EC105b) The system of EC100, wherein the system performs an inference
12 application.
13
- 14 **[0098]** EC106) The system of EC100, wherein the system implements a deep neural network
15 trained to perform object classification and/or detection.
16
- 17 **[0099]** EC107) The system of EC100, wherein the system implements a deep neural network
18 trained to perform an inference application selected from the group consisting of text translation,
19 optical character recognition, image classification, facial recognition, scene recognition for a self-
20 driving car, speech recognition, data analysis for high energy physics, and drug discovery.
21
- 22 **[0100]** EC108) The system of EC100, wherein the fabric is organized as a plurality of
23 periphery processor elements and a plurality of interior processor elements, and each of the interior
24 processor elements is coupled in at least four logical directions respectively to at least four others of
25 the plurality of processor elements.
26
- 27 **[0101]** EC109) The system of EC100, wherein each compute engine comprises a memory, a
28 data path, and a hybrid dataflow and instruction execution controller.
29
- 30 **[0102]** EC110) The system of EC109, wherein each compute engine operates in accordance
31 with a multi-stage compute engine pipeline having a plurality of compute engine pipeline stages.
32

1 [0103] EC111) The system of EC109, wherein the instruction execution controller comprises
2 an instruction sequencer implemented using one or more of microcode, PLAs, one or more counters,
3 and a gate-level state machine.

4
5 [0104] EC112) The system of EC109, wherein each compute engine further comprises a
6 register file, an instruction decoder, an instruction cache, and a data cache.

7
8 [0105] EC112b) The system of EC109, wherein each compute engine further comprises a
9 register file, an instruction decoder, an instruction buffer, and a data buffer.

10
11 [0106] EC113) The system of EC100, wherein:
12 each compute engine is configured to perform a predefined set of basic operations in response
13 to receiving a corresponding basic instruction selected from a predefined native
14 instruction set of codes; and further comprising
15 a training workload comprising
16 a first set of machine codes selected from the native instruction set for performing a
17 mapping of at least a part of a neuron onto the compute engine of the
18 processor element, the mapping comprising managing at least one partial-
19 neuron weight,
20 a second set of machine codes selected from the native instruction set for performing
21 a forward pass to propagate activations in a forward logical direction based at
22 least in part on the at least one partial-neuron weight, the forward pass
23 initiated responsive to an input sample,
24 a third set of machine codes selected from the native instruction set for performing a
25 delta pass in a backward logical direction to generate deltas, the delta pass
26 initiated responsive to completion of the forward pass,
27 a fourth set of machine codes selected from the native instruction set for performing a
28 chain pass to calculate gradients based on the deltas, and
29 a fifth set of machine codes selected from the native instruction set for performing a
30 selective update of the at least one partial-neuron weight in accordance with a
31 predetermined learning rule and based at least in part on the deltas; and
32 wherein each compute engine comprises storage for the at least one partial-neuron weight.

33

1 [0107] EC113a) The system of EC113, wherein each basic instruction is performed in
2 accordance with the task specifier of a respective fabric packet of the fabric packets.
3

4 [0108] EC113b) The system of EC113, wherein the fabric comprises a 2D array of the
5 processor elements comprising a first, second, third, and fourth physical directions, the first and
6 second physical directions being collinear and opposite, the third and fourth physical directions being
7 collinear and opposite, the first and third physical directions being orthogonal, and the forward logical
8 direction is in the first physical direction and the backward logical direction is in the second physical
9 direction.
10

11 [0109] EC113c) The system of EC113, wherein the training workload further comprises a
12 sixth set of machine codes selected from the native instruction set for performing a nonlinear
13 activation function.
14

15 [0110] EC113d) The system of EC113c, wherein the nonlinear activation function is
16 selected from the group consisting of sigmoid, tanh, and ReLU.
17

18 [0111] EC114) The system of EC113, wherein the mapping is in accordance with initializing
19 the fabric to implement a partitioning of a neural network into a plurality of layers, the neuron is a first
20 neuron of a plurality of neurons of the neural network, the first neuron is comprised in a first layer of
21 the plurality of layers, and each of the plurality of neurons is mapped in a distributed manner across a
22 plurality of the processor elements of the fabric.
23

24 [0112] EC115) The system of EC114, wherein the mapping is in accordance with each input
25 sample of a training set completing all of the passes for each layer in the same amount of time.
26

27 [0113] EC115b) The system of EC114, wherein the mapping is in accordance with each
28 input sample of a training set completing all of the passes for each layer within a same predetermined
29 amount of time.
30

31 [0114] EC115c) The system of EC114, wherein the mapping is in accordance with each
32 input sample of a training set completing all of the passes for each layer within a same time period
33 determined in real time.
34

1 **[0115]** EC116) The system of EC114, wherein the plurality of layers operates as a logical
2 fabric pipeline comprising logical fabric pipeline stages, each logical fabric pipeline stage comprising
3 completion of all of the passes for each layer, the completion for each layer taking a time step
4 comprising the same amount of time.

5

6 **[0116]** EC116b) The system of EC114, wherein each of the plurality of layers operates as a
7 logical fabric pipeline stage of a respective logical fabric pipeline of each of the passes, the
8 completion for each layer taking a time step comprising the same amount of time.

9

10 **[0117]** EC117) The system of EC114, wherein as each input sample of a training set streams
11 through at least a first plurality of the processor elements across the plurality of layers, the neuron
12 weights are selectively updated in the first plurality of the processor elements across the plurality of
13 layers.

14

15 **[0118]** EC117b) The system of EC118, wherein as each input sample of a training set
16 streams through at least a first plurality of the processor elements across the plurality of layers, the
17 neuron weights are selectively updated in the first plurality of the processor elements across the
18 plurality of layers, and the streaming and updating is ongoing for each time step over a plurality of
19 time steps.

20

21 **[0119]** EC119) The system of EC120, further comprising a digital clock, and wherein the
22 time step is an integral multiple of a clock-cycle of the digital clock.

23

24 **[0120]** EC118b) The system of EC120, further comprising a digital clock, and wherein the
25 time step is a variable amount of time.

26

27 **[0121]** EC118c) The system of EC121 or EC118b, wherein the time step is determined in
28 real-time.

29

1 **[0122]** EC122) The system of EC114, further comprising:
2 wherein each compute engine operates in accordance with a multi-stage compute engine
3 pipeline having a plurality of compute engine pipeline stages, a compute engine
4 machine cycle comprising the time to complete each compute engine pipeline stage, a
5 compute engine pipeline cycle comprising the time to complete the plurality of
6 compute engine pipeline stages;
7 wherein the compute engine machine cycle comprises a first multiple of a clock-cycle of a
8 digital clock;
9 wherein the plurality of layers operates as a logical fabric pipeline comprising logical fabric
10 pipeline stages, each logical fabric pipeline stage comprising completion of all of the
11 passes for each layer, a time step comprising the time to complete each logical fabric
12 pipeline stage; and
13 wherein the time step comprises a second multiple of the compute engine pipeline cycle.
14

15 **[0123]** EC123) The system of EC122, wherein the first multiple is one.
16

17 **[0124]** EC124) The system of EC122, wherein the second multiple is in the hundreds to
18 thousands.
19

20 **[0125]** EC125) The system of EC120, wherein for each time step over a plurality of time
21 steps while forward propagation of activations are ongoing, the at least one partial-neuron weight is
22 selectively updated within a first plurality of the processor elements in response to changes in
23 backward propagating data within the first plurality of the processor elements.
24

25 **[0126]** EC126) The system of EC120, wherein the at least one partial-neuron weight is
26 selectively updated each time step over a plurality of time steps.
27

28 **[0127]** EC123b) The system of EC117, EC117b, EC122, or EC123, wherein the selective
29 updating is in accordance with a continuous propagation gradient descent process.
30

31 **[0128]** EC127) The system of EC114, wherein the neural network comprises over a thousand
32 layers.
33

- 1 **[0129]** EC128) The system of EC114, wherein the plurality of neurons comprises billions of
2 neurons.
3
- 4 **[0130]** EC125b) The system of EC114, wherein the plurality of neurons comprises millions
5 of neurons.
6
- 7 **[0131]** EC125c) The system of EC114, wherein the neural network comprises at least 10
8 weights per neuron for at least some of the plurality of neurons.
9
- 10 **[0132]** EC125d) The system of EC114, wherein the neural network comprises at least 1000
11 weights per neuron for at least some of the plurality of neurons.
12
- 13 **[0133]** EC129) The system of EC114, wherein the neural network comprises billions of
14 weights per layer.
15
- 16 **[0134]** EC126b) The system of EC114, wherein the neural network comprises millions of
17 weights per layer.
18
- 19 **[0135]** EC130) The system of EC114, wherein for each layer of the neural network,
20 incoming activations are weighted to create partial sums that are accumulated to generate output
21 activations for the layer, and the accumulated weighted partial sums represent the neurons and
22 associated synapses of the neural network.
23
- 24 **[0136]** EC127b) The system of EC127, wherein each weight corresponds to a synapse, each
25 partial sum corresponds to a stimulus, the accumulated weighted partial sums correspond to a total
26 stimulus, and each output activation for the layer corresponds to a neuron output.
27
- 28 **[0137]** EC131) The system of EC113, wherein an iteration of the training workload is
29 performed for each of a plurality of input samples collectively comprising a training set.
30
- 31 **[0138]** EC132) The system of EC131, wherein the predetermined learning rule specifies that
32 the at least one partial-neuron weight is updated after the completion of all the passes for the entire
33 training set.
34

- 1 **[0139]** EC129b) The system of EC129, wherein the predetermined learning rule is in
2 accordance with a stochastic gradient descent process.
3
- 4 **[0140]** EC129c) The system of EC129, wherein the predetermined learning rule is in
5 accordance with a mini-batch gradient descent process.
6
- 7 **[0141]** EC129d) The system of EC129, wherein the predetermined learning rule is in
8 accordance with a continuous propagation gradient descent process.
9
- 10 **[0142]** EC133) The system of EC131, wherein the training set is partitioned into a plurality
11 of so-called mini-batches and the predetermined learning rule specifies that the at least one partial-
12 neuron weight is updated after the completion of all the passes for the input samples comprised in
13 each of the mini-batches.
14
- 15 **[0143]** EC134) The system of EC131, wherein the training set is partitioned into a plurality
16 of so-called mini-batches and the predetermined learning rule specifies that the at least one partial-
17 neuron weight is updated after the completion of all the passes for each input sample of each of the
18 mini-batches.
19
- 20 **[0144]** EC131b) The system of EC131, wherein the predetermined learning rule is in
21 accordance with a continuous propagation gradient descent process.
22
- 23 **[0145]** EC135) The system of EC134, wherein the forward pass incorporates weight updates
24 within a first plurality of the processor elements while the mini-batch learning is ongoing within the
25 first plurality of the processor elements.
26
- 27 **[0146]** EC136) The system of EC113, wherein the storage is comprised in a memory local to
28 the compute engine.
29
- 30 **[0147]** EC133b) The system of EC113, wherein the storage is comprised in the compute
31 engine.
32

- 1 **[0148]** EC133b) The system of EC113, wherein the storage is a respective memory attached
2 to each compute engine.
3
- 4 **[0149]** EC137) The system of EC113, wherein the storage is enabled to store a 2D matrix
5 data structure.
6
- 7 **[0150]** EC134b) The system of EC113, wherein the storage is enabled to store a
8 multidimensional data structure.
9
- 10 **[0151]** EC134c) The system of EC113, wherein the storage is enabled to store a tensor data
11 structure comprising a dimension selected from the group consisting of 2D, 3D, 4D, 5D, and 6D.
12
- 13 **[0152]** EC138) The system of EC113, wherein each compute engine further comprises
14 storage for gradient accumulation, forward partial sums, delta partial sums, and forward pass
15 activations.
16
- 17 **[0153]** EC139) The system of EC114, wherein data propagates to a logical end of the neural
18 network during the forward pass and circulates back in a reverse logical direction during the delta and
19 chain passes.
20
- 21 **[0154]** EC140) The system of EC113, wherein the forward pass saves the activations for use
22 by the delta and chain passes.
23
- 24 **[0155]** EC141) The system of EC113, wherein each processor element is time shared across
25 the forward, delta and chain passes.
26
- 27 **[0156]** EC142) The system of EC131, wherein for each input sample, the system is enabled
28 to selectively update the at least one partial-neuron weight in accordance with the predetermined
29 learning rule responsive to completion of the forward pass, the delta pass, and the chain pass
30 corresponding to the input sample.
31
- 32 **[0157]** EC139b) The system of EC139, wherein the predetermined learning rule is in
33 accordance with a continuous propagation gradient descent process.
34

1 **[0158]** EC143) The system of EC142, wherein the system is enabled for each forward pass to
2 use weight information provided by the most recent selective update of the at least one partial-neuron
3 weight.

4
5 **[0159]** EC144) The system of EC143, wherein the system is enabled to initiate a forward
6 pass of a particular iteration of the training workload independent of whether the selective update of
7 the at least one partial-neuron weight corresponding to a prior iteration of the training workload has
8 occurred.

9
10 **[0160]** EC145) The system of EC143, wherein the system is enabled to initiate a forward
11 pass of a particular iteration of the training workload independent of whether the delta pass of a prior
12 iteration of the training workload has begun.

13
14 **[0161]** EC146) The system of EC143, wherein at least one compute engine is enabled to
15 perform at least a portion of a forward pass for a subsequent iteration of the training workload after
16 performing at least a portion of a forward pass for a prior iteration of the training workload and before
17 performing a portion of the selective update of the at least one partial-neuron weight corresponding to
18 the prior iteration of the training workload.

19
20 **[0162]** EC147) The system of EC143, wherein the system is enabled to perform the delta
21 pass and the chain pass for each input sample based at least in part on activations that are recomputed
22 based at least in part on a first partial-neuron weight.

23
24 **[0163]** EC148) The system of EC147, wherein the first partial-neuron weight is the partial-
25 neuron weight produced by the most recent selective update.

26
27 **[0164]** EC145b) The system of EC145, wherein the recomputed activations need not be
28 stored between computations, thereby decreasing the total memory required for a given system
29 training configuration.

30
31 **[0165]** EC145c) The system of EC139, EC140, EC141, or EC142, wherein concurrent layer
32 training enables achieving a predetermined accuracy goal at a faster convergence rate, thereby
33 decreasing total training time required for a given system training configuration.

34

1 **[0166]** EC145d) The system of EC139, EC140, EC141, or EC142, wherein concurrent layer
2 training enables increased accuracy for a given total training time and system training configuration.

3

4 **[0167]** EC149) The system of EC143, wherein each compute element is enabled to perform
5 portions of a delta pass and portions of a chain pass for an input sample based at least in part on
6 activations that are recomputed based at least in part on a first partial-neuron weight.

7

8 **[0168]** EC150) The system of EC149, wherein the first partial-neuron weight is the partial-
9 neuron weight produced by the most recent selective update.

10

11 **[0169]** EC200) A method comprising:

12 in each of a fabric of processor elements, selectively communicating fabric packets with
13 others of the processor elements, each processor element comprising a fabric router
14 and a compute engine enabled to perform dataflow-based and instruction-based
15 processing; and

16 in each compute engine, selectively performing the processing in accordance with a virtual
17 channel specifier and a task specifier of each fabric packet the compute engine
18 receives.

19

20 **[0170]** EC200b) A method comprising:

21 in each of a fabric of processor elements, selectively communicating fabric packets with
22 others of the processor elements, each processor element comprising a fabric router
23 and a compute engine; and

24 in each compute engine, selectively performing dataflow processing and instruction
25 processing respectively in accordance with a dataflow field and an instruction field of
26 each fabric packet the compute engine receives.

27

28 **[0171]** EC200c) The method of EC200, wherein the processing is in accordance with a data-
29 flow graph.

30

31 **[0172]** EC200d) The method of EC200, further comprising executing a workload comprising
32 predominantly dataflow-based processing with minimal instruction-based processing.

33

- 1 **[0173]** EC200e) The method of EC200d, wherein performing the method implements a
2 Long Short Term Memory (LSTM) neural network model.
3
- 4 **[0174]** EC200f) The method of EC200, further comprising executing a workload comprising
5 predominantly instruction-based processing with minimal dataflow-based processing.
6
- 7 **[0175]** EC200g) The method of EC200, wherein the fabric of processor elements is
8 implemented at least in part using wafer-scale integration.
9
- 10 **[0176]** EC200h) The method of EC200, wherein the fabric of processor elements is
11 implemented at least in part using VLSI fabrication.
12
- 13 **[0177]** EC201) The method of EC200, wherein the virtual channel specifier selects
14 independent respective routing paths in the fabric.
15
- 16 **[0178]** EC201b) The method of EC200, wherein the virtual channel specifier selects routing
17 paths in the fabric to perform multicast.
18
- 19 **[0179]** EC201c) The method of EC200, wherein the virtual channel specifier selects routing
20 paths in the fabric to perform load splitting.
21
- 22 **[0180]** EC202) The method of EC200, wherein the task specifier selects one or more
23 operations to perform.
24
- 25 **[0181]** EC203) The method of EC200, wherein the fabric comprises a 2D array of the
26 processor elements.
27
- 28 **[0182]** EC203b) The method of EC200, wherein the fabric comprises a processor element
29 interconnection topology selected from the group consisting of fully connected, star, ring, array, mesh,
30 hypercube, torus, and tree.
31

1 **[0183]** EC203c) The method of EC200, wherein the fabric comprises a processor element
2 interconnection topology dimension selected from the group consisting of 1D, 2D, 3D, and a
3 dimension greater than 3D.
4

5 **[0184]** EC204) The method of EC200, wherein performing the method enables executing
6 machine learning workloads.
7

8 **[0185]** EC205) The method of EC200, wherein performing the method enables training an
9 inference application.
10

11 **[0186]** EC205b) The method of EC200, wherein performing the method performs an
12 inference application.
13

14 **[0187]** EC206) The method of EC200, wherein performing the method implements a deep
15 neural network trained to perform object classification and/or detection.
16

17 **[0188]** EC207) The method of EC200, wherein performing the method implements a deep
18 neural network trained to perform an inference application selected from the group consisting of text
19 translation, optical character recognition, image classification, facial recognition, scene recognition for
20 a self-driving car, speech recognition, data analysis for high energy physics, and drug discovery.
21

22 **[0189]** EC208) The method of EC200, wherein the fabric is organized as a plurality of
23 periphery processor elements and a plurality of interior processor elements, and each of the interior
24 processor elements is coupled in at least four logical directions respectively to at least four others of
25 the plurality of processor elements.
26

27 **[0190]** EC209) The method of EC200, wherein each compute engine comprises a memory, a
28 data path, and a hybrid dataflow and instruction execution controller.
29

30 **[0191]** EC210) The method of EC209, wherein each compute engine operates in accordance
31 with a multi-stage compute engine pipeline having a plurality of compute engine pipeline stages.
32

1 [0192] EC211) The method of EC209, wherein the instruction execution controller comprises
2 an instruction sequencer implemented using one or more of microcode, PLAs, one or more counters,
3 and a gate-level state machine.

4
5 [0193] EC212) The method of EC209, wherein each compute engine further comprises a
6 register file, an instruction decoder, an instruction cache, and a data cache.

7
8 [0194] EC212b) The method of EC209, wherein each compute engine further comprises a
9 register file, an instruction decoder, an instruction buffer, and a data buffer.

10
11 [0195] EC213) The method of EC200, wherein:
12 each compute engine is configured to perform a predefined set of basic operations in response
13 to receiving a corresponding basic instruction selected from a predefined native
14 instruction set of codes; and further comprising
15 processing a training workload comprising
16 a first set of machine codes selected from the native instruction set for performing a
17 mapping of at least a part of a neuron onto the compute engine of the
18 processor element, the mapping comprising managing at least one partial-
19 neuron weight,
20 a second set of machine codes selected from the native instruction set for performing
21 a forward pass to propagate activations in a forward logical direction based at
22 least in part on the at least one partial-neuron weight, the forward pass
23 initiated responsive to an input sample,
24 a third set of machine codes selected from the native instruction set for performing a
25 delta pass in a backward logical direction to generate deltas, the delta pass
26 initiated responsive to completion of the forward pass,
27 a fourth set of machine codes selected from the native instruction set for performing a
28 chain pass to calculate gradients based on the deltas, and
29 a fifth set of machine codes selected from the native instruction set for performing a
30 selective update of the at least one partial-neuron weight in accordance with a
31 predetermined learning rule and based at least in part on the deltas; and
32 wherein each compute engine comprises storage for the at least one partial-neuron weight.

33

1 **[0196]** EC213a) The method of EC213, wherein each basic instruction is performed in
2 accordance with the task specifier of a respective fabric packet of the fabric packets.
3

4 **[0197]** EC213b) The method of EC213, wherein the fabric comprises a 2D array of the
5 processor elements comprising a first, second, third, and fourth physical directions, the first and
6 second physical directions being collinear and opposite, the third and fourth physical directions being
7 collinear and opposite, the first and third physical directions being orthogonal, and the forward logical
8 direction is in the first physical direction and the backward logical direction is in the second physical
9 direction.
10

11 **[0198]** EC213c) The method of EC213, wherein the training workload further comprises a
12 sixth set of machine codes selected from the native instruction set for performing a nonlinear
13 activation function.
14

15 **[0199]** EC213d) The method of EC213c, wherein the nonlinear activation function is
16 selected from the group consisting of sigmoid, tanh, and ReLU.
17

18 **[0200]** EC214) The method of EC213, wherein the mapping is in accordance with initializing
19 the fabric to implement a partitioning of a neural network into a plurality of layers, the neuron is a first
20 neuron of a plurality of neurons of the neural network, the first neuron is comprised in a first layer of
21 the plurality of layers, and each of the plurality of neurons is mapped in a distributed manner across a
22 plurality of the processor elements of the fabric.
23

24 **[0201]** EC215) The method of EC214, wherein the mapping is in accordance with each input
25 sample of a training set completing all of the passes for each layer in the same amount of time.
26

27 **[0202]** EC215b) The method of EC214, wherein the mapping is in accordance with each
28 input sample of a training set completing all of the passes for each layer within a same predetermined
29 amount of time.
30

31 **[0203]** EC215c) The method of EC214, wherein the mapping is in accordance with each
32 input sample of a training set completing all of the passes for each layer within a same time period
33 determined in real time.
34

1 **[0204]** EC216) The method of EC214, wherein the plurality of layers operates as a logical
2 fabric pipeline comprising logical fabric pipeline stages, each logical fabric pipeline stage comprising
3 completion of all of the passes for each layer, the completion for each layer taking a time step
4 comprising the same amount of time.

5

6 **[0205]** EC216b) The method of EC214, wherein each of the plurality of layers operates as a
7 logical fabric pipeline stage of a respective logical fabric pipeline of each of the passes, the
8 completion for each layer taking a time step comprising the same amount of time.

9

10 **[0206]** EC217) The method of EC214, wherein as each input sample of a training set streams
11 through at least a first plurality of the processor elements across the plurality of layers, the neuron
12 weights are selectively updated in the first plurality of the processor elements across the plurality of
13 layers.

14

15 **[0207]** EC217b) The method of EC216, wherein as each input sample of a training set
16 streams through at least a first plurality of the processor elements across the plurality of layers, the
17 neuron weights are selectively updated in the first plurality of the processor elements across the
18 plurality of layers, and the streaming and updating is ongoing for each time step over a plurality of
19 time steps.

20

21 **[0208]** EC218) The method of EC216, wherein at least one of the processor elements
22 comprises a digital clock, and the time step is an integral multiple of a clock-cycle of the digital clock.

23

24 **[0209]** EC218b) The method of EC216, wherein at least one of the processor elements
25 comprises a digital clock, and wherein the time step is a variable amount of time.

26

27 **[0210]** EC218c) The method of EC218 or EC218b, wherein the time step is determined in
28 real-time.

29

- 1 **[0211]** EC219) The method of EC214, further comprising:
2 operating each compute engine in accordance with a multi-stage compute engine pipeline
3 having a plurality of compute engine pipeline stages, a compute engine machine cycle
4 comprising the time to complete each compute engine pipeline stage, a compute
5 engine pipeline cycle comprising the time to complete the plurality of compute engine
6 pipeline stages;
7 wherein the compute engine machine cycle comprises a first multiple of a clock-cycle of a
8 digital clock;
9 wherein the plurality of layers operates as a logical fabric pipeline comprising logical fabric
10 pipeline stages, each logical fabric pipeline stage comprising completion of all of the
11 passes for each layer, a time step comprising the time to complete each logical fabric
12 pipeline stage; and
13 wherein the time step comprises a second multiple of the compute engine pipeline cycle.
14
- 15 **[0212]** EC220) The method of EC219, wherein the first multiple is one.
16
- 17 **[0213]** EC221) The method of EC219, wherein the second multiple is in the hundreds to
18 thousands.
19
- 20 **[0214]** EC222) The method of EC216, further comprising, for each time step over a plurality
21 of time steps while forward propagation of activations are ongoing, selectively updating the at least
22 one partial-neuron weight within a first plurality of the processor elements in response to changes in
23 backward propagating data within the first plurality of the processor elements.
24
- 25 **[0215]** EC223) The method of EC216, further comprising selectively updating the at least
26 one partial-neuron weight each time step over a plurality of time steps.
27
- 28 **[0216]** EC223b) The method of EC217, EC217b, EC222, or EC223, wherein the selectively
29 updating is in accordance with a continuous propagation gradient descent process.
30
- 31 **[0217]** EC224) The method of EC214, wherein the neural network comprises over a
32 thousand layers.
33

- 1 **[0218]** EC225) The method of EC214, wherein the plurality of neurons comprises billions of
2 neurons.
3
- 4 **[0219]** EC225b) The method of EC214, wherein the plurality of neurons comprises millions
5 of neurons.
6
- 7 **[0220]** EC225c) The method of EC214, wherein the neural network comprises at least 10
8 weights per neuron for at least some of the plurality of neurons.
9
- 10 **[0221]** EC225d) The method of EC214, wherein the neural network comprises at least 1000
11 weights per neuron for at least some of the plurality of neurons.
12
- 13 **[0222]** EC226) The method of EC214, wherein the neural network comprises billions of
14 weights per layer.
15
- 16 **[0223]** EC226b) The method of EC214, wherein the neural network comprises millions of
17 weights per layer.
18
- 19 **[0224]** EC227) The method of EC214, further comprising, for each layer of the neural
20 network, weighting incoming activations to create partial sums that are accumulated to generate output
21 activations for the layer, and wherein the accumulated weighted partial sums represent the neurons
22 and associated synapses of the neural network.
23
- 24 **[0225]** EC227b) The method of EC227, wherein each weight corresponds to a synapse, each
25 partial sum corresponds to a stimulus, the accumulated weighted partial sums correspond to a total
26 stimulus, and each output activation for the layer corresponds to a neuron output.
27
- 28 **[0226]** EC228) The method of EC213, further comprising performing an iteration of the
29 training workload for each of a plurality of input samples collectively comprising a training set.
30
- 31 **[0227]** EC229) The method of EC228, wherein the predetermined learning rule specifies that
32 the at least one partial-neuron weight is updated after the completion of all the passes for the entire
33 training set.
34

- 1 **[0228]** EC229b) The method of EC229, wherein the predetermined learning rule is in
2 accordance with a stochastic gradient descent process.
3
- 4 **[0229]** EC229c) The method of EC229, wherein the predetermined learning rule is in
5 accordance with a mini-batch gradient descent process.
6
- 7 **[0230]** EC229d) The method of EC229, wherein the predetermined learning rule is in
8 accordance with a continuous propagation gradient descent process.
9
- 10 **[0231]** EC230) The method of EC228, further comprising partitioning the training set into a
11 plurality of so-called mini-batches and the predetermined learning rule specifies that the at least one
12 partial-neuron weight is updated after the completion of all the passes for the input samples comprised
13 in each of the mini-batches.
14
- 15 **[0232]** EC231) The method of EC228, further comprising partitioning the training set into a
16 plurality of so-called mini-batches and the predetermined learning rule specifies that the at least one
17 partial-neuron weight is updated after the completion of all the passes for each input sample of each of
18 the mini-batches.
19
- 20 **[0233]** EC231b) The method of EC231, wherein the predetermined learning rule is in
21 accordance with a continuous propagation gradient descent process.
22
- 23 **[0234]** EC232) The method of EC231, wherein the forward pass incorporates weight updates
24 within a first plurality of the processor elements while the mini-batch learning is ongoing within the
25 first plurality of the processor elements.
26
- 27 **[0235]** EC233) The method of EC213, wherein the storage is comprised in a memory local to
28 the compute engine.
29
- 30 **[0236]** EC233b) The method of C213, wherein the storage is comprised in the compute
31 engine.
32

- 1 **[0237]** EC233b) The method of C213, wherein the storage is a respective memory attached to
2 each compute engine.
3
- 4 **[0238]** EC234) The method of EC213, wherein the storage is enabled to store a 2D matrix
5 data structure.
6
- 7 **[0239]** EC234b) The method of C213, wherein the storage is enabled to store a
8 multidimensional data structure.
9
- 10 **[0240]** EC234c) The method of C213, wherein the storage is enabled to store a tensor data
11 structure comprising a dimension selected from the group consisting of 2D, 3D, 4D, 5D, and 6D.
12
- 13 **[0241]** EC235) The method of EC213, wherein each compute engine further comprises
14 storage for gradient accumulation, forward partial sums, delta partial sums, and forward pass
15 activations.
16
- 17 **[0242]** EC236) The method of EC214, wherein data propagates to a logical end of the neural
18 network during the forward pass and circulates back in a reverse logical direction during the delta and
19 chain passes.
20
- 21 **[0243]** EC237) The method of EC213, wherein the forward pass saves the activations for use
22 by the delta and chain passes.
23
- 24 **[0244]** EC238) The method of EC213, further comprising time sharing each processor
25 element across the forward, delta and chain passes.
26
- 27 **[0245]** EC239) The method of EC228, further comprising, for each input sample, selectively
28 updating the at least one partial-neuron weight in accordance with the predetermined learning rule
29 responsive to completion of the forward pass, the delta pass, and the chain pass corresponding to the
30 input sample.
31
- 32 **[0246]** EC239b) The method of EC239, wherein the predetermined learning rule is in
33 accordance with a continuous propagation gradient descent process.
34

1 [0247] EC240) The method of EC239, further comprising, for each forward pass, selectively
2 using weight information provided by the most recent selective update of the at least one partial-
3 neuron weight.

4
5 [0248] EC241) The method of EC240, further comprising initiating a forward pass of a
6 particular iteration of the training workload independent of whether the selective update of the at least
7 one partial-neuron weight corresponding to a prior iteration of the training workload has occurred.

8
9 [0249] EC242) The method of EC240, further comprising selectively initiating a forward
10 pass of a particular iteration of the training workload independent of whether the delta pass of a prior
11 iteration of the training workload has begun.

12
13 [0250] EC243) The method of EC240, further comprising, in at least one of the compute
14 engines, performing at least a portion of a forward pass for a subsequent iteration of the training
15 workload after performing at least a portion of a forward pass for a prior iteration of the training
16 workload and before performing a portion of the selective update of the at least one partial-neuron
17 weight corresponding to the prior iteration of the training workload.

18
19 [0251] EC244) The method of EC240, further comprising selectively performing the delta
20 pass and the chain pass for each input sample based at least in part on activations that are recomputed
21 based at least in part on a first partial-neuron weight.

22
23 [0252] EC245) The method of EC244, wherein the first partial-neuron weight is the partial-
24 neuron weight produced by the most recent selective update.

25
26 [0253] EC245b) The method of EC245, wherein the recomputed activations need not be
27 stored between computations, thereby decreasing the total memory required for a given system
28 training configuration.

29
30 [0254] EC245c) The method of EC239, EC240, EC241, or EC242, wherein concurrent layer
31 training enables achieving a predetermined accuracy goal at a faster convergence rate, thereby
32 decreasing total training time required for a given system training configuration.

33

1 **[0255]** EC245d) The method of EC239, EC240, EC241, or EC242, wherein concurrent layer
2 training enables increased accuracy for a given total training time and system training configuration.

3

4 **[0256]** EC246) The method of EC240, further comprising, in each compute element,
5 selectively performing portions of a delta pass and portions of a chain pass for an input sample based
6 at least in part on activations that are recomputed based at least in part on a first partial-neuron weight.

7

8 **[0257]** EC247) The method of EC246, wherein the first partial-neuron weight is the partial-
9 neuron weight produced by the most recent selective update.

10

11 **[0258]** EC300) A system comprising:
12 in each of a fabric of processor elements, means for selectively communicating fabric packets
13 with others of the processor elements, each processor element comprising a fabric
14 router and a compute engine enabled to perform dataflow-based and instruction-based
15 processing; and

16 in each compute engine, means for selectively performing the processing in accordance with a
17 virtual channel specifier and a task specifier of each fabric packet the compute engine
18 receives.

19

20 **[0259]** EC300b) A system comprising:
21 in each of a fabric of processor elements, means for selectively communicating fabric packets
22 with others of the processor elements, each processor element comprising a fabric
23 router and a compute engine; and

24 in each compute engine, means for selectively performing dataflow processing and instruction
25 processing respectively in accordance with a dataflow field and an instruction field of
26 each fabric packet the compute engine receives.

27

28 **[0260]** EC300c) The system of EC300, wherein the processing is in accordance with a data-
29 flow graph.

30

31 **[0261]** EC300d) The system of EC300, further comprising means for executing a workload
32 comprising predominantly dataflow-based processing with minimal instruction-based processing.

33

1 **[0262]** EC300e) The system of EC300d, wherein the system implements a Long Short Term
2 Memory (LSTM) neural network model.

3

4 **[0263]** EC300f) The system of EC300, further comprising means for executing a workload
5 comprising predominantly instruction-based processing with minimal dataflow-based processing.

6

7 **[0264]** EC300g) The system of EC300, wherein the system is implemented at least in part
8 using wafer-scale integration.

9

10 **[0265]** EC300h) The system of EC300, wherein the fabric of processor elements is
11 implemented at least in part using VLSI fabrication.

12

13 **[0266]** EC301) The system of EC300, wherein the virtual channel specifier selects
14 independent respective routing paths in the fabric.

15

16 **[0267]** EC301b) The system of EC300, wherein the virtual channel specifier selects routing
17 paths in the fabric to perform multicast.

18

19 **[0268]** EC301c) The system of EC300, wherein the virtual channel specifier selects routing
20 paths in the fabric to perform load splitting.

21

22 **[0269]** EC302) The system of EC300, wherein the task specifier selects one or more
23 operations to perform.

24

25 **[0270]** EC303) The system of EC300, wherein the fabric comprises a 2D array of the
26 processor elements.

27

28 **[0271]** EC303b) The system of EC300, wherein the fabric comprises a processor element
29 interconnection topology selected from the group consisting of fully connected, star, ring, array, mesh,
30 hypercube, torus, and tree.

31

- 1 **[0272]** EC303c) The system of EC300, wherein the fabric comprises a processor element
2 interconnection topology dimension selected from the group consisting of 1D, 2D, 3D, and a
3 dimension greater than 3D.
4
- 5 **[0273]** EC304) The system of EC300, wherein the system is enabled to execute machine
6 learning workloads.
7
- 8 **[0274]** EC305) The system of EC300, wherein the system is trained to perform an inference
9 application.
10
- 11 **[0275]** EC305b) The system of EC300, wherein the system performs an inference
12 application.
13
- 14 **[0276]** EC306) The system of EC300, wherein the system implements a deep neural network
15 trained to perform object classification and/or detection.
16
- 17 **[0277]** EC307) The system of EC300, wherein the system implements a deep neural network
18 trained to perform an inference application selected from the group consisting of text translation,
19 optical character recognition, image classification, facial recognition, scene recognition for a self-
20 driving car, speech recognition, data analysis for high energy physics, and drug discovery.
21
- 22 **[0278]** EC308) The system of EC300, wherein the fabric is organized as a plurality of
23 periphery processor elements and a plurality of interior processor elements, and each of the interior
24 processor elements is coupled in at least four logical directions respectively to at least four others of
25 the plurality of processor elements.
26
- 27 **[0279]** EC309) The system of EC300, wherein each compute engine comprises a memory, a
28 data path, and a hybrid dataflow and instruction execution controller.
29
- 30 **[0280]** EC310) The system of EC309, wherein each compute engine operates in accordance
31 with a multi-stage compute engine pipeline having a plurality of compute engine pipeline stages.
32

1 [0281] EC311) The system of EC309, wherein the instruction execution controller comprises
2 an instruction sequencer implemented using one or more of microcode, PLAs, one or more counters,
3 and a gate-level state machine.

4
5 [0282] EC312) The system of EC309, wherein each compute engine further comprises a
6 register file, an instruction decoder, an instruction cache, and a data cache.

7
8 [0283] EC312b) The system of EC309, wherein each compute engine further comprises a
9 register file, an instruction decoder, an instruction buffer, and a data buffer.

10
11 [0284] EC313) The system of EC300, wherein:
12 each compute engine is configured to perform a predefined set of basic operations in response
13 to receiving a corresponding basic instruction selected from a predefined native
14 instruction set of codes; and further comprising
15 a training workload comprising
16 a first set of machine codes selected from the native instruction set for performing a
17 mapping of at least a part of a neuron onto the compute engine of the
18 processor element, the mapping comprising managing at least one partial-
19 neuron weight,
20 a second set of machine codes selected from the native instruction set for performing
21 a forward pass to propagate activations in a forward logical direction based at
22 least in part on the at least one partial-neuron weight, the forward pass
23 initiated responsive to an input sample,
24 a third set of machine codes selected from the native instruction set for performing a
25 delta pass in a backward logical direction to generate deltas, the delta pass
26 initiated responsive to completion of the forward pass,
27 a fourth set of machine codes selected from the native instruction set for performing a
28 chain pass to calculate gradients based on the deltas, and
29 a fifth set of machine codes selected from the native instruction set for performing a
30 selective update of the at least one partial-neuron weight in accordance with a
31 predetermined learning rule and based at least in part on the deltas; and
32 wherein each compute engine comprises storage for the at least one partial-neuron weight.
33

1 **[0285]** EC313a) The system of EC313, wherein each basic instruction is performed in
2 accordance with the task specifier of a respective fabric packet of the fabric packets.
3

4 **[0286]** EC313b) The system of EC313, wherein the fabric comprises a 2D array of the
5 processor elements comprising a first, second, third, and fourth physical directions, the first and
6 second physical directions being collinear and opposite, the third and fourth physical directions being
7 collinear and opposite, the first and third physical directions being orthogonal, and the forward logical
8 direction is in the first physical direction and the backward logical direction is in the second physical
9 direction.
10

11 **[0287]** EC313c) The system of EC313, wherein the training workload further comprises a
12 sixth set of machine codes selected from the native instruction set for performing a nonlinear
13 activation function.
14

15 **[0288]** EC313d) The system of EC313c, wherein the nonlinear activation function is
16 selected from the group consisting of sigmoid, tanh, and ReLU.
17

18 **[0289]** EC314) The system of EC313, wherein the mapping is in accordance with initializing
19 the fabric to implement a partitioning of a neural network into a plurality of layers, the neuron is a first
20 neuron of a plurality of neurons of the neural network, the first neuron is comprised in a first layer of
21 the plurality of layers, and each of the plurality of neurons is mapped in a distributed manner across a
22 plurality of the processor elements of the fabric.
23

24 **[0290]** EC315) The system of EC314, wherein the mapping is in accordance with each input
25 sample of a training set completing all of the passes for each layer in the same amount of time.
26

27 **[0291]** EC315b) The system of EC314, wherein the mapping is in accordance with each
28 input sample of a training set completing all of the passes for each layer within a same predetermined
29 amount of time.
30

31 **[0292]** EC315c) The system of EC314, wherein the mapping is in accordance with each
32 input sample of a training set completing all of the passes for each layer within a same time period
33 determined in real time.
34

1 [0293] EC316) The system of EC314, wherein the plurality of layers operates as a logical
2 fabric pipeline comprising logical fabric pipeline stages, each logical fabric pipeline stage comprising
3 completion of all of the passes for each layer, the completion for each layer taking a time step
4 comprising the same amount of time.

5

6 [0294] EC316b) The system of EC314, wherein each of the plurality of layers operates as a
7 logical fabric pipeline stage of a respective logical fabric pipeline of each of the passes, the
8 completion for each layer taking a time step comprising the same amount of time.

9

10 [0295] EC317) The system of EC314, wherein as each input sample of a training set streams
11 through at least a first plurality of the processor elements across the plurality of layers, the neuron
12 weights are selectively updated in the first plurality of the processor elements across the plurality of
13 layers.

14

15 [0296] EC317b) The system of EC316, wherein as each input sample of a training set
16 streams through at least a first plurality of the processor elements across the plurality of layers, the
17 neuron weights are selectively updated in the first plurality of the processor elements across the
18 plurality of layers, and the streaming and updating is ongoing for each time step over a plurality of
19 time steps.

20

21 [0297] EC318) The system of EC316, further comprising a digital clock, and wherein the
22 time step is an integral multiple of a clock-cycle of the digital clock.

23

24 [0298] EC318b) The system of EC316, further comprising a digital clock, and wherein the
25 time step is a variable amount of time.

26

27 [0299] EC318c) The system of EC318 or EC318b, wherein the time step is determined in
28 real-time.

29

- 1 **[0300]** EC319) The system of EC314, further comprising:
2 means for operating each compute engine in accordance with a multi-stage compute engine
3 pipeline having a plurality of compute engine pipeline stages, a compute engine
4 machine cycle comprising the time to complete each compute engine pipeline stage, a
5 compute engine pipeline cycle comprising the time to complete the plurality of
6 compute engine pipeline stages;
7 wherein the compute engine machine cycle comprises a first multiple of a clock-cycle of a
8 digital clock;
9 wherein the plurality of layers operates as a logical fabric pipeline comprising logical fabric
10 pipeline stages, each logical fabric pipeline stage comprising completion of all of the
11 passes for each layer, a time step comprising the time to complete each logical fabric
12 pipeline stage; and
13 wherein the time step comprises a second multiple of the compute engine pipeline cycle.
14
- 15 **[0301]** EC320) The system of EC319, wherein the first multiple is one.
16
- 17 **[0302]** EC321) The system of EC319, wherein the second multiple is in the hundreds to
18 thousands.
19
- 20 **[0303]** EC322) The system of EC316, further comprising means for selectively updating the
21 at least one partial-neuron weight within a first plurality of the processor elements in response to
22 changes in backward propagating data within the first plurality of the processor elements for each time
23 step over a plurality of time steps while forward propagation of activations are ongoing.
24
- 25 **[0304]** EC323) The system of EC316, further comprising means for selectively updating the
26 at least one partial-neuron weight each time step over a plurality of time steps.
27
- 28 **[0305]** EC323b) The system of EC317, EC317b, EC322, or EC323, wherein the selectively
29 updating is in accordance with a continuous propagation gradient descent process.
30
- 31 **[0306]** EC324) The system of EC314, wherein the neural network comprises over a thousand
32 layers.
33

1 **[0307]** EC325) The system of EC314, wherein the plurality of neurons comprises billions of
2 neurons.

3

4 **[0308]** EC325b) The system of EC314, wherein the plurality of neurons comprises millions
5 of neurons.

6

7 **[0309]** EC325c) The system of EC314, wherein the neural network comprises at least 10
8 weights per neuron for at least some of the plurality of neurons.

9

10 **[0310]** EC325d) The system of EC314, wherein the neural network comprises at least 1000
11 weights per neuron for at least some of the plurality of neurons.

12

13 **[0311]** EC326) The system of EC314, wherein the neural network comprises billions of
14 weights per layer.

15

16 **[0312]** EC326b) The system of EC314, wherein the neural network comprises millions of
17 weights per layer.

18

19 **[0313]** EC327) The system of EC314, further comprising, for each layer of the neural
20 network, means for weighting incoming activations to create partial sums that are accumulated to
21 generate output activations for the layer, and wherein the accumulated weighted partial sums represent
22 the neurons and associated synapses of the neural network.

23

24 **[0314]** EC327b) The system of EC327, wherein each weight corresponds to a synapse, each
25 partial sum corresponds to a stimulus, the accumulated weighted partial sums correspond to a total
26 stimulus, and each output activation for the layer corresponds to a neuron output.

27

28 **[0315]** EC328) The system of EC313, further comprising means for performing an iteration
29 of the training workload for each of a plurality of input samples collectively comprising a training set.

30

31 **[0316]** EC329) The system of EC328, wherein the predetermined learning rule specifies that
32 the at least one partial-neuron weight is updated after the completion of all the passes for the entire
33 training set.

34

- 1 **[0317]** EC329b) The system of EC329, wherein the predetermined learning rule is in
2 accordance with a stochastic gradient descent process.
3
- 4 **[0318]** EC329c) The system of EC329, wherein the predetermined learning rule is in
5 accordance with a mini-batch gradient descent process.
6
- 7 **[0319]** EC329d) The system of EC329, wherein the predetermined learning rule is in
8 accordance with a continuous propagation gradient descent process.
9
- 10 **[0320]** EC330) The system of EC328, further comprising means for partitioning the training
11 set into a plurality of so-called mini-batches and the predetermined learning rule specifies that the at
12 least one partial-neuron weight is updated after the completion of all the passes for the input samples
13 comprised in each of the mini-batches.
14
- 15 **[0321]** EC331) The system of EC328, further means for comprising partitioning the training
16 set into a plurality of so-called mini-batches and the predetermined learning rule specifies that the at
17 least one partial-neuron weight is updated after the completion of all the passes for each input sample
18 of each of the mini-batches.
19
- 20 **[0322]** EC331b) The system of EC331, wherein the predetermined learning rule is in
21 accordance with a continuous propagation gradient descent process.
22
- 23 **[0323]** EC332) The system of EC331, wherein the forward pass incorporates weight updates
24 within a first plurality of the processor elements while the mini-batch learning is ongoing within the
25 first plurality of the processor elements.
26
- 27 **[0324]** EC333) The system of EC313, wherein the storage is comprised in a memory local to
28 the compute engine.
29
- 30 **[0325]** EC333b) The system of EC313, wherein the storage is comprised in the compute
31 engine.
32

- 1 **[0326]** EC333b) The system of EC313, wherein the storage is a respective memory attached
2 to each compute engine.
3
- 4 **[0327]** EC334) The system of EC313, wherein the storage is enabled to store a 2D matrix
5 data structure.
6
- 7 **[0328]** EC334b) The system of EC313, wherein the storage is enabled to store a
8 multidimensional data structure.
9
- 10 **[0329]** EC334c) The system of EC313, wherein the storage is enabled to store a tensor data
11 structure comprising a dimension selected from the group consisting of 2D, 3D, 4D, 5D, and 6D.
12
- 13 **[0330]** EC335) The system of EC313, wherein each compute engine further comprises
14 storage for gradient accumulation, forward partial sums, delta partial sums, and forward pass
15 activations.
16
- 17 **[0331]** EC336) The system of EC314, wherein data propagates to a logical end of the neural
18 network during the forward pass and circulates back in a reverse logical direction during the delta and
19 chain passes.
20
- 21 **[0332]** EC337) The system of EC313, wherein the forward pass saves the activations for use
22 by the delta and chain passes.
23
- 24 **[0333]** EC338) The system of EC313, further comprising means for time sharing each
25 processor element across the forward, delta and chain passes.
26
- 27 **[0334]** EC339) The system of EC328, further comprising, for each input sample, means for
28 selectively updating the at least one partial-neuron weight in accordance with the predetermined
29 learning rule responsive to completion of the forward pass, the delta pass, and the chain pass
30 corresponding to the input sample.
31
- 32 **[0335]** EC339b) The system of EC339, wherein the predetermined learning rule is in
33 accordance with a continuous propagation gradient descent process.
34

1 **[0336]** EC340) The system of EC339, further comprising means for selectively using weight
2 information provided by the most recent selective update of the at least one partial-neuron weight for
3 each forward pass.

4
5 **[0337]** EC341) The system of EC340, further comprising means for initiating a forward pass
6 of a particular iteration of the training workload independent of whether the selective update of the at
7 least one partial-neuron weight corresponding to a prior iteration of the training workload has
8 occurred.

9
10 **[0338]** EC342) The system of EC340, further comprising means for selectively initiating a
11 forward pass of a particular iteration of the training workload independent of whether the delta pass of
12 a prior iteration of the training workload has begun.

13
14 **[0339]** EC343) The system of EC340, further comprising, in at least one of the compute
15 engines, means for performing at least a portion of a forward pass for a subsequent iteration of the
16 training workload after performing at least a portion of a forward pass for a prior iteration of the
17 training workload and before performing a portion of the selective update of the at least one partial-
18 neuron weight corresponding to the prior iteration of the training workload.

19
20 **[0340]** EC344) The system of EC340, further comprising means for selectively performing
21 the delta pass and the chain pass for each input sample based at least in part on activations that are
22 recomputed based at least in part on a first partial-neuron weight.

23
24 **[0341]** EC345) The system of EC344, wherein the first partial-neuron weight is the partial-
25 neuron weight produced by the most recent selective update.

26
27 **[0342]** EC345b) The system of EC345, wherein the recomputed activations need not be
28 stored between computations, thereby decreasing the total memory required for a given system
29 training configuration.

30
31 **[0343]** EC345c) The system of EC339, EC340, EC341, or EC342, wherein concurrent layer
32 training enables achieving a predetermined accuracy goal at a faster convergence rate, thereby
33 decreasing total training time required for a given system training configuration.

34

1 **[0344]** EC345d) The system of EC339, EC340, EC341, or EC342, wherein concurrent layer
2 training enables increased accuracy for a given total training time and system training configuration.

3

4 **[0345]** EC346) The system of EC340, further comprising, in each compute element, means
5 for selectively performing portions of a delta pass and portions of a chain pass for an input sample
6 based at least in part on activations that are recomputed based at least in part on a first partial-neuron
7 weight.

8

9 **[0346]** EC347) The system of EC346, wherein the first partial-neuron weight is the partial-
10 neuron weight produced by the most recent selective update.

11

12 **[0347]** EC400) A method comprising:
13 training a neural network comprising a plurality of ordered, connected layers;
14 wherein the order identifies for each respective layer which others of the layers are prior to the
15 respective layer and which others of the layers are subsequent to the respective layer;
16 wherein each layer comprises one or more neurons, each neuron comprising weights and
17 connected to at least one of at least one prior neuron of a prior layer, and at least one
18 subsequent neuron of a subsequent layer; and
19 wherein each neuron is implemented by one or more processing elements, each processing
20 element comprising
21 at least one coupling to a fabric the processing element being enabled to
22 communicate via the fabric via a plurality of virtual channels,
23 a first memory enabled to store instructions corresponding to at least
24 computations of the neuron,
25 a second memory enabled to store the weights, and
26 hardware execution resources enabled to execute instructions from the
27 respective first memory and access data from the respective second
28 memory.

29

1 **[0348]** EC401) The method of EC400, wherein the training comprises:
2 based on a first activation and first weights, determining a second activation;
3 based on a first delta and the first weights, determining and saving second weights,
4 based on a third activation and selected weights, determining a fourth activation, wherein the
5 selected weights are dynamically selected from the first weights and the second
6 weights; and
7 based on a second delta and the selected weights, determining and saving third weights.
8

9 **[0349]** EC402) The method of EC401, wherein the determining the second activation
10 comprises:
11 receiving the first activation via the fabric from the at least one prior neuron;
12 computing the second activation based at least in part on the first activation and first weights
13 by at least executing first instructions stored in the first memory and accessing the
14 first weights in the second memory; and
15 selectively transmitting the second activation via the fabric to the at least one subsequent
16 neuron.
17

18 **[0350]** EC403) The method of EC401, wherein the determining and saving the second
19 weights comprises:
20 receiving the first delta that is partially based on the second activation via the fabric from the
21 at least one subsequent neuron;
22 computing a first gradient based at least in part on the first delta and the second activation by
23 at least executing second instructions stored in the first memory;
24 computing the second weights based at least in part on the first gradient, a learning rule, and
25 the first weights by at least executing third instructions stored in the first memory and
26 accessing the first weights in the second memory; and
27 storing the second weights in the second memory.
28

1 **[0351]** EC404) The method of EC402, wherein the determining the fourth activation
2 comprises:
3 receiving the third activation via the fabric from the at least one prior neuron;
4 computing the fourth activation based at least in part on the third activation and the selected
5 weights by at least executing the first instructions and accessing the selected weights
6 in the second memory; and
7 selectively transmitting the fourth activation via the fabric to the at least one subsequent
8 neuron.

9
10 **[0352]** EC405) The method of EC403, wherein the determining and saving third weights
11 comprises:
12 receiving the second delta that is partially based on the fourth activation via the fabric from
13 the at least one subsequent neuron;
14 computing a second gradient based at least in part on a third delta and the fourth activation by
15 at least executing the second instructions stored in the first memory;
16 computing the third weights based at least in part on the second gradient, the learning rule and
17 the selected weights by at least executing the third instructions stored and accessing
18 the selected weights in the second memory; and
19 storing the third weights in the second memory.

20
21 **[0353]** EC406) The method of EC404, wherein the determining and saving the second
22 weights comprises:
23 receiving the first delta that is partially based on the second activation via the fabric from the
24 at least one subsequent neuron;
25 computing a first gradient based at least in part on the first delta and the second activation by
26 at least executing second instructions stored in the first memory;
27 computing the second weights based at least in part on the first gradient, a learning rule, and
28 the first weights by at least executing third instructions stored in the first memory and
29 accessing the first weights in the second memory; and
30 storing the second weights in the second memory.

31

1 **[0354]** EC407) The method of EC406, wherein the determining and saving third weights
2 comprises:
3 receiving the second delta that is partially based on the fourth activation via the fabric from
4 the at least one subsequent neuron;
5 computing a second gradient based at least in part on a third delta and the fourth activation by
6 at least executing the second instructions stored in the first memory;
7 computing the third weights based at least in part on the second gradient, the learning rule and
8 the selected weights by at least executing the third instructions stored and accessing
9 the selected weights in the second memory; and
10 storing the third weights in the second memory.

11
12 **[0355]** EC408) The method of EC403, wherein the selected weights are dynamically selected
13 in accordance with which of the first weights and the second weights was stored most recently.
14

15 **[0356]** EC409) The method of EC401, wherein the determining the fourth activation is
16 enabled to be performed after the determining the second activation and before the determining and
17 saving the second weights.
18

19 **[0357]** EC410) The method of EC404, wherein the selectively transmitting the second
20 activation and the fourth activation is selectively based upon the respective values of the second
21 activation and fourth activation.
22

23 **[0358]** EC411) The method of EC404, wherein the selectively transmitting the second
24 activation and the fourth activation is selectively based upon the respective absolute values of the
25 second activation and the fourth activation exceeding respective first and second thresholds.
26

27 **[0359]** EC412) The method of EC400, wherein at least one neuron is implemented by a
28 plurality of processing elements.
29

30 **[0360]** EC413) The method of EC405, wherein the determining the fourth activation
31 additionally comprises storing the fourth activation in the second memory and the computing the
32 second gradient additionally comprises accessing the fourth activation in the second memory.
33

- 1 **[0361]** EC414) The method of EC407, wherein the computing the second gradient
2 additionally comprises optionally recomputing the fourth activation based at least in part upon the
3 selected weights.
4
- 5 **[0362]** EC415) The method of EC407, wherein the computing the first gradient additionally
6 comprises optionally recomputing the second activation based at least in part upon the first weights.
7
- 8 **[0363]** EC416) The method of EC400, wherein each processing element is enabled to
9 perform dataflow-based processing.
10
- 11 **[0364]** EC417) The method of EC400, wherein each processing element comprises a fabric
12 router.
13
- 14 **[0365]** EC418) The method of EC400, wherein each processing element is enabled to
15 selectively communicate fabric packets with others of the processing elements.
16
- 17 **[0366]** EC419) The method of EC418, wherein each processing element is enabled to
18 perform processing in accordance with a virtual channel specifier and a task specifier of each fabric
19 packet the processing element receives.
20

21

22 **SELECTED EMBODIMENT DETAILS**

23

- 24 **[0367]** Embodiments relating to neural network training and inference, comprising deep
25 learning accelerator hardware elements and software elements are described herein (see, e.g., Figs. 1-4
26 and section “Deep Learning Accelerator Overview”). The deep learning accelerator comprises
27 hardware processing elements (see, e.g., Figs. 5-8 and section “Processing Element: Compute
28 Element and Router”). The deep learning accelerator implements and/or uses various techniques such
29 as task initiation and closeout (see, e.g., Figs. 9-12 and section “Tasks”), wavelet processing (see, e.g.,
30 Figs. 13A-15B and section “Wavelets”), task blocking and unblocking (see, e.g., Fig. 16 and section
31 “Block and Unblock”), neuron smearing (see, e.g., Figs. 17-20 and section “Neuron Smearing”),
32 fabric vectors, memory vectors, and associated data structure descriptors (see, e.g., Figs. 21A-24 and
33 section “Vectors and Data Structure Descriptors”), and instruction formats (see, e.g., Figs. 25A-25C
34 and section “Instruction Formats”). The deep learning accelerator is usable in a variety of scenarios

1 (see, e.g., Figs. 26A-27E and section “Deep Learning Accelerator Example Uses” as well as Figs.
2 28A-29 and section “Example Workload Mapping”). The deep learning accelerator is contemplated in
3 various embodiments (see, e.g., section “Other Embodiment Details”). The deep learning accelerator
4 is variously implementable (see, e.g., section “Example Implementation Techniques”).
5
6

7 DEEP LEARNING ACCELERATOR OVERVIEW

8

9 **[0368]** Fig. 1 illustrates selected details of an embodiment of a system for neural network
10 training and inference, using a deep learning accelerator, as Neural Network System 100.
11 Conceptually a neural network is trained using the deep learning accelerator. One or more results of
12 the training (e.g., weights) are then used for inferences. For example, the training comprises mapping
13 neurons of the neural network onto PEs of the deep learning accelerator. Then training data is applied
14 to the PEs. The PEs process the training data (e.g., via forward, delta, and chain passes) and update
15 weights until the training is complete. Then the weights are used for inference.
16

17 **[0369]** Referring to the figure, Deep Learning Accelerator 120 comprises FPGAs 121 and
18 PEs 122, enabled to communicate with each other, as illustrated by Coupling 123. Placement
19 Server(s) 150, (comprising CPUs 151 and CRM 152) is coupled to Connection Server(s) 160
20 (comprising CPUs 161, CRM 162, and NICs 164) via LAN 111. Connection Server(s) 160 is enabled
21 to communicate with FPGAs 121 via NICs 164 and 100Gb 112. Autonomous Vehicle 130
22 comprises CPUs 131, CRM 132, IEs 133, and Camera 135. Cell Phone 140 comprises CPUs 141,
23 CRM 142, IEs 143, and Camera 145.
24

25 **[0370]** Internet 180 provides for coupling (not explicitly illustrated) between any
26 combination of Placement Server(s) 150, Connection Server(s) 160, Autonomous Vehicle 130, and/or
27 Cell Phone 140, according to various embodiments and/or usage scenarios.
28

29 **[0371]** Dashed-arrow Placements 113 conceptually indicates placement information
30 communicated from Placement Server(s) 150 to PEs 122 (e.g., via LAN 111, Connection Server(s)
31 160 / NICs 164, 100Gb 112, FPGAs 121, and Coupling 123). In some embodiments and/or usage
32 scenarios, Placements 113 is implicit, reflected in initialization information provided to router
33 elements of PEs 122 and compute elements of PEs 122. In some embodiments and/or usage scenarios,

1 a portion of initialization information of Placements 113 is provided to FPGAs 121 to configure
2 elements of FPGAs 121 for operation with PEs 122.

3

4 [0372] Dashed-arrow Weights 114 and dashed-arrow Weights 115 conceptually indicate
5 weight information communicated from PEs 122 respectively to Autonomous Vehicle 130 and Cell
6 Phone 140 (e.g., via Coupling 123, FPGAs 121, 100Gb 112, Connection Server(s) 160 / NICs 164 and
7 Internet 180). In some embodiments and/or usage scenarios, the weight information is any one or
8 more of all or any portions of weight information as directly produced as a result of training, a sub-
9 sampling thereof, a quantization thereof, and/or other transformations thereof.

10

11 [0373] Deep Learning Accelerator 120 is enabled to perform training of neural networks,
12 such as by computing weights in response to placement information and training information received
13 via 100Gb 112. Deep Learning Accelerator 120 is further enabled to, upon training completion,
14 provide the weights as results via 100Gb 112. The weights are then usable for inference, such as in
15 Autonomous Vehicle 130 and/or in Cell Phone 140. PEs 122 comprises a relatively large number of
16 PEs (e.g., 10,000 or more) each enabled to independently perform routing and computations relating
17 to training. In some embodiments and/or usage scenarios, PEs 122 is implemented via wafer-scale
18 integration, such as respective pluralities of PEs implemented on respective dice of a single wafer.
19 FPGAs 121 is enabled to interface PEs 122 to information provided via 100Gb 112. The interfacing
20 includes conversion to/from modified Ethernet frames from/to Wavelets, as communicated on
21 Coupling 123.

22

23 [0374] Placement Server(s) 150 is enabled to programmatically determine placements of
24 neurons (e.g., as indicated by Placements 113) via one or more placement programs. The placement
25 programs are stored in CRM 152 and executed by CPUs 151. The placement information is
26 communicated to Connection Server(s) 160 via LAN 111. An example of a placement is a mapping of
27 logical neurons of a neural network onto physical memory and execution hardware resources (e.g.,
28 PEs 122).

29

30 [0375] Connection Server(s) 160 is enabled to communicate with FPGAs 121 and indirectly
31 with PEs 122 via FPGAs 121 / Coupling 123, via NICs 164 and programmed control thereof via
32 driver programs. In various embodiments and/or usage scenarios, the communication comprises
33 placement information (e.g., from Placement Server(s) 150), training information (e.g., from sources

1 not illustrated but accessible via Internet 180) and/or results of training (e.g., weights from PEs 122).
2 The driver programs are stored in CRM 162 and executed by CPUs 161.

3

4 [0376] Autonomous Vehicle 130 is enabled to use Weights 114 to perform inferences using
5 IEs 133 as programmatically controlled and/or assisted by CPUs 131 executing programs stored in
6 CRM 132. The inferences are optionally and/or selectively performed using information obtained
7 from Camera 135. For example, a car is operable as an autonomous vehicle. The car comprises
8 cameras enabled to provide video to an inference engine. The inference engine is enabled to
9 recognize objects related to navigating the car, such as traffic lanes, obstructions, and other objects.
10 The car is enabled to navigate using results of the object recognition. Any combination of the
11 providing, the recognizing, and the navigating are controlled and/or performed at least in part via one
12 or more CPUs executing programs stored in a CRM.

13

14 [0377] Cell Phone 140 is enabled to use Weights 115 to perform inferences using IEs 143 as
15 programmatically controlled and/or assisted by CPUs 141 executing programs stored in CRM 142.
16 The inferences are optionally and/or selectively performed using information obtained from Camera
17 145. For example, the cell phone is operable to post tagged photos on a social networking web site.
18 The cell phone comprises a camera enabled to provide image data to an inference engine. The
19 inference engine is enabled to tag objects (e.g., by type such as 'cat', 'dog', and so forth, or by name
20 such as 'Bob', 'Mary', and so forth) in the image. The cell phone is enabled to post the image and
21 results of the tagging to the social networking web site. Any combination of the providing, the
22 tagging, and the posting are controlled and/or performed at least in part via one or more CPUs
23 executing programs stored in a CRM.

24

25 [0378] In various embodiments and/or usage scenarios, all or any portions of weight
26 information determined via a deep learning accelerator is post-processed outside of the accelerator
27 before inference usage. For example, all or any portions of information represented by Weights 114
28 and/or Weights 115, is processed in whole or in part by Placement Server(s) 150 before inference
29 usage by Autonomous Vehicle 130 and/or Cell Phone 140. In various embodiments and/or usage
30 scenarios, an example of post-processing comprises quantizing Weights 114 and/or Weights 115 (e.g.,
31 converting from a floating-point number format to a fixed-point number format). In various
32 embodiments and/or usage models, Camera 135 and Camera 145 are respective examples of sensors
33 that provide input to IEs 133 and IEs 143. Other examples of sensors are location sensors, orientation
34 sensors, magnetic sensors, light sensors, and pressure sensors.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

[0379] CPUs 151 comprises one or more CPUs that are compatible with respective instruction set architectures. CPUs 151 is enabled to fetch and execute instructions from CRM 152 in accordance with the instruction set architectures. CPUs 161 comprises one or more CPUs that are compatible with respective instruction set architectures. CPUs 161 is enabled to fetch and execute instructions from CRM 162 in accordance with the instruction set architectures. In some embodiments, at least one of the instruction set architectures of CPUs 151 is compatible with at least one of the instruction set architectures of CPUs 161.

[0380] CPUs 131 comprises one or more CPUs that are compatible with respective instruction set architectures. CPUs 131 is enabled to fetch and execute instructions from CRM 132 in accordance with the instruction set architectures. CPUs 141 comprises one or more CPUs that are compatible with respective instruction set architectures. CPUs 141 is enabled to fetch and execute instructions from CRM 142 in accordance with the instruction set architectures. In some embodiments, at least one of the instruction set architectures of CPUs 131 is compatible with at least one of the instruction set architectures of CPUs 141. In some embodiments, any one or more of CPUs 151, CPUs 161, CPUs 131, and CPUs 141 have instruction set architectures that are compatible with each other.

[0381] At least a respective portion of each of CRM 152 and CRM 162 CRM 132, and CRM 142, is non-volatile and comprised of any one or more of flash memory, magnetic memory, optical memory, phase-change memory, and other non-volatile memory technology elements.

[0382] In various embodiments and/or usage scenarios, IEs 133 and/or IEs 143 comprise one or more inference engines enabled to use weight information as determined by Deep Learning Accelerator 120 (and indicated conceptually by Weights 114 and/or Weights 115). In various embodiments and/or usage scenarios, IEs 133 operates in conjunction with and/or under control of programs executed by CPUs 131 and stored in CRM 132. In various embodiments and/or usage scenarios, IEs 143 operates in conjunction with and/or under control of programs executed by CPUs 141 and stored in CRM 142. In various embodiments and/or usage scenarios, all or any portions of IEs 133 and/or IEs 143 are implemented via various combinations of HW and/or SW techniques. In some embodiments, all or any portions of functionality provided by IEs 133 and/or IEs 143 is implemented using techniques such as implemented by and/or associated with Deep Learning Accelerator 120. In various embodiments and/or usage scenarios, all or any portions of IEs 133 and/or

1 IEs 143 are variously implemented via techniques comprising various combinations of conventional
2 CPUs, conventional GPUs, conventional DSPs, conventional FPGAs, and specialized hardware.

3

4 [0383] In various embodiments, 100Gb 112, is variously a 100Gb Ethernet coupling for
5 sending standard Ethernet frames, a 100Gb Ethernet coupling for sending modified Ethernet frames, a
6 100GB modified Ethernet coupling for sending modified Ethernet frames, a 100Gb serial coupling of
7 other-than Ethernet technology, or some other relatively high-speed serial coupling.

8

9 [0384] In some embodiments and/or usage scenarios, Coupling 123 communicates
10 information as wavelets.

11

12 [0385] In various embodiments, LAN 111 is implemented using techniques such as Ethernet,
13 Fibre Channel, and/or other suitable interconnection technologies.

14

15 [0386] In some embodiments and/or usage scenarios, Placement Server(s) 150 and
16 Connection Server(s) 160 are implemented and/or operated as a combined element (e.g., sharing CPU,
17 CRM, and/or NIC resources), as illustrated conceptually by Combined Server(s) 110. In some
18 embodiments and/or usage scenarios, Placement Server(s) 150 and Connection Server(s) 160 are
19 coupled via Internet 180 rather than (or in addition to) LAN 111.

20

21 [0387] Fig. 2 illustrates selected details of an embodiment of software elements associated
22 with neural network training and inference, using a deep learning accelerator, as Neural Network
23 Software 200. Placement Server(s) SW 210 comprises Neuron to PE Mapping SW 212, as well as
24 other elements not illustrated, according to embodiment. In various embodiments and/or usage
25 scenarios, all or any portions of Placement Server(s) SW 210 is stored in CRM 152 and executable by
26 CPUs 151 of Fig. 1. One or more programs of Neuron to PE Mapping SW 212 enable determining
27 placements of neurons of a neural network onto specific PEs of PEs 122 of Fig. 1.

28

29 [0388] Connection Server(s) SW 220 comprises 100Gb NIC Driver 224, Training Info
30 Provider SW 225, and Weight Receiver SW 226, as well as other elements not illustrated, according to
31 embodiment. In various embodiments and/or usage scenarios, all or any portions of Connection
32 Server(s) SW 220 is stored in CRM 162 and executable by CPUs 161 of Fig. 1. One or more
33 programs of 100Gb NIC Driver 224 enable communication between Connection Server(s) 160 and
34 Deep Learning Accelerator 120, both of Fig. 1 (via NICs 164 and 100Gb 112, also of Fig. 1). One or

1 more programs of Training Info Provider SW 225 enable determination of training information for
2 application under control of 100Gb NIC Driver 224 for communication to Deep Learning Accelerator
3 120 of Fig. 1 (via NICs 164 and 100Gb 112). In various embodiments and/or usage scenarios, the
4 training information is variously determined from, e.g., non-volatile storage accessible to Connection
5 Server(s) 160 and/or Internet 180, both of Fig. 1. One or more programs of Weight Receiver SW 226
6 enable receiving weight information under control of 100Gb NIC Driver 224 as determined by Deep
7 Learning Accelerator 120 (via NICs 164 and 100Gb 112).

8

9 [0389] In various embodiments and/or usage scenarios, Misc SW on FPGAs 250
10 conceptually represents SW executed by one or more CPUs comprised in FPGAs 121 of (Fig. 1). The
11 CPUs of the FPGAs are, e.g., hard-coded during manufacturing of one or more elements of FPGAs
12 121, and/or soft-coded during initialization of one or more elements of FPGAs 121. In various
13 embodiments and/or usage scenarios, all or any portions of Misc SW on FPGAs 250 and/or a
14 representation thereof is stored in non-volatile memory comprised in FPGAs 121 and/or accessible to
15 Connection Server(s) 160. In various embodiments and/or usage scenarios, Misc SW on FPGAs 250
16 enables performing various housekeeping functions, such as relating to initialization and/or debugging
17 of PEs 122 of Fig. 1.

18

19 [0390] In various embodiments and/or usage scenarios, Task SW on PEs 260 conceptually
20 represents distributed SW executed as tasks on various PEs of PEs 122. In various embodiments
21 and/or usage scenarios, all or any portions of Task SW on PEs 260 and/or a representation thereof is
22 stored in non-volatile memory comprised in PEs 122 and/or accessible to Connection Server(s) 160.
23 In various embodiments and/or usage scenarios, Task SW on PEs 260 enables performing processing
24 of training data such as to determine weights of a neural network (e.g., via forward, delta, and chain
25 passes).

26

27 [0391] Autonomous Vehicle SW 230 comprises Video Camera SW 232, Inference Engine(s)
28 SW 233, and Navigating SW 234, as well as other elements not illustrated, according to embodiment.
29 In various embodiments and/or usage scenarios, all or any portions of Autonomous Vehicle SW 230 is
30 stored in CRM 132 and executable by CPUs 131 of Fig. 1. One or more programs of Video Camera
31 SW 232 enable controlling and/or operating Camera 135 of Fig. 1 to provide video information to
32 Inference Engine(s) SW 233. One or more programs of Inference Engine(s) SW 233 enable
33 controlling and/or operating IEs 133 of Fig. 1 to determine navigational information, such as objects
34 to avoid and/or traffic lanes to follow, from the video information. One or more programs of

1 Navigating SW 234 enable navigating Autonomous Vehicle SW 230 in response to the navigational
2 information.

3

4 [0392] Cell Phone SW 240 comprises Still Camera SW 242, Inference Engine(s) SW 243,
5 Posting SW 244, as well as other elements not illustrated, according to embodiment. In various
6 embodiments and/or usage scenarios, all or any portions of Cell Phone SW 240 is stored in CRM 142
7 and executable by CPUs 141 of Fig. 1. One or more programs of Still Camera SW 242 enable
8 controlling and/or operating Camera 145 of Fig. 1 to provide still image information to Inference
9 Engine(s) SW 243. One or more programs of Inference Engine(s) SW 243 enable controlling and/or
10 operating IEs 143 of Fig. 1 to determine tag information from the still image information. One or
11 more programs of Posting SW 244 enable posting to a social networking web site in response to the
12 still image information and/or the tag information.

13

14 [0393] In various embodiments and/or usage scenarios, any one or more of SW collections
15 Placement Server(s) SW 210, Connection Server(s) SW 220, Autonomous Vehicle SW 230, and/or
16 Cell Phone SW 240 optionally and/or selectively comprise one or more operating system elements,
17 e.g., one or more real-time operating systems, one or more non-real-time operating systems, and/or
18 one or more other control programs to coordinate elements of each respective SW collection.

19

20 [0394] Fig. 3 illustrates selected details of an embodiment of processing associated with
21 training a neural network and performing inference using the trained neural network, using a deep
22 learning accelerator, as Neural Network Training/Inference 300. As illustrated, neurons of the neural
23 network are placed, e.g., allocated and/or associated with specific PE resources in action 310. Then
24 FPGA resources are initialized in preparation for training of the neural network in action 320. Then
25 the PE resources are initialized in preparation for training of the neural network in action 330.

26

27 [0395] After the FPGA resources and PE resources are initialized in preparation for the
28 training, training data is applied to the PEs in action 340. The PE resources process the training data
29 in action 350. Then a check is made to determine if training is complete, e.g., because application of
30 the training data is complete and/or one or more completion criteria are met (such as an inference error
31 below a predetermine bound) in action 360. If not, then flow passes back to action 340 for application
32 of further training data. In some scenarios, the training does not complete and in some embodiments,
33 control instead passes to another action (not illustrated) to enable changing the neural network (e.g.,

1 adding layers of neurons, removing layers of neurons). The changed neural network is then trained in
2 accordance with actions 310, 320, 330, 340, 350, and 360.

3

4 [0396] If training is complete, then flow continues to provide weights that are results of the
5 training for use in inferences in 370. In some embodiments and/or usage scenarios, the weights are
6 quantized, e.g., transformed to an integer data format. In some embodiments and/or usage scenarios,
7 the integer data format is a reduced precision number format (e.g., 8-bit or 16-bit). The weights are
8 then provided to one or more inference engines, and used to make inferences in action 380.

9

10 [0397] In various embodiments and/or usage scenarios, the inference engines correspond to
11 one or more inference applications, e.g., text translation, optical character recognition, image
12 classification, facial recognition, scene recognition for a self-driving car, speech recognition, data
13 analysis for high energy physics, and drug discovery.

14

15 [0398] In various embodiments and/or usage scenarios, the PE resources correspond, e.g., to
16 PEs 122 of Fig. 1, and the FPGAs resources correspond, e.g., to FPGAs 121 of Fig. 1.

17

18 [0399] In various embodiments and/or usage scenarios, any one or more of all or any
19 portions of actions of Neural Network Training/Inference 300 are performed by and/or related to all or
20 any portions of any one or more elements of Neural Network System 100 of Fig. 1 and/or Neural
21 Network Software 200 of Fig. 2. For example, all or any portions of action 310 are performed by
22 Placement Server(s) 150 via execution of Neuron to PE Mapping SW 212. For another example, all
23 or any portions of action 320 are performed by Placement Server(s) 150 via execution of Neuron to
24 PE Mapping SW 212. For another example, all or any portions of action 330 are performed by
25 Placement Server(s) 150 via execution of Neuron to PE Mapping SW 212. For another example, all
26 or any portions of action 330 are performed by PEs 122 via execution of Task SW on PEs 260. For
27 another example, all or any portions of action 340 are performed by Connection Server(s) 160 via
28 execution of Training Info Provider SW 225. For another example, all or any portions of action 350
29 are performed by PEs 122 via execution of Task SW on PEs 260. For another example, all or any
30 portions of action 350 are performed by Combined Server(s) 110, Placement Server(s) 150 and/or
31 Connection Server(s) 160. For another example, all or any portions of 370 are performed by
32 Connection Server(s) 160 via execution of Weight Receiver SW 226. For another example, all or any
33 portions of action 370 are performed by FPGAs 121 via execution of Misc SW on FPGAs 250. For
34 another example, all or any portions of 380 are performed by IEs 133 such as under control of

1 Inference Engine(s) SW 233. For another example, all or any portions of action 380 are performed by
2 IEs 143 such as under control of Inference Engine(s) SW 243.

3

4 [0400] In various embodiments and/or usage scenarios, any one or more of all or any
5 portions of actions of Neural Network Training/Inference 300 are performed in conjunction with
6 communicating information between various elements of Neural Network System 100 of Fig. 1. For
7 example, various actions of Neural Network Training/Inference 300 are performed at least in part via
8 NICs 164 and 100Gb 112 communicating information between Connection Server(s) 160 and FPGAs
9 121. For another example, various actions of Neural Network Training/Inference 300 are performed
10 in conjunction with FPGAs 121 and Coupling 123 communicating information between Connection
11 Server(s) 160 and PEs 122. For another example, various actions of Neural Network
12 Training/Inference 300 performed in conjunction with any one or more of Placement Server(s) 150,
13 Connection Server(s) 160, Autonomous Vehicle 130, and Cell Phone 140 communicating information
14 as enabled at least in part by Internet 180.

15

16 [0401] Fig. 4 illustrates selected details of an embodiment of a deep learning accelerator as
17 Deep Learning Accelerator 400. Each of PE 499 elements has couplings to other of PE 499 elements.
18 Two of the PE elements (PE 497 and PE 498) are illustrated with unique identifiers, and are otherwise
19 respectively identical to a instances of PE 499. PE 497 is illustrated with identifiers for each of four
20 couplings (North coupling 430, East coupling 431 with PE 498, and South coupling 432) to others of
21 the PEs and one of the I/O FPGAs (West coupling 433), but is otherwise identical to others of the PE
22 elements illustrated. In some embodiments and/or usage scenarios, the couplings are logical and/or
23 physical. In various embodiments and/or usage scenarios, the couplings are usable to communicate
24 wavelets, backpressure information, or both. In various embodiments and/or usage scenarios, all or
25 any portions of the physical couplings are to physically adjacent PEs. In some embodiments and/or
26 usage scenarios, the PEs are physically implemented in a 2D grid. In some embodiments and/or usage
27 scenarios, the PEs are physically implemented in a 2D grid of aligned rectangles, and physically
28 adjacent PEs correspond to PEs sharing a horizontal boundary (North/South PEs with respect to each
29 other) and PEs sharing a vertical boundary (East/West PEs with respect to each other).

30

31 [0402] In some embodiments and/or usage scenarios, an array of identical instances of a
32 same ASIC is formed on a wafer, and each of the same ASICs comprises a plurality of identical
33 instances of a same PE (e.g., PE 499), forming a wafer (e.g., Wafer 412) usable in wafer-scale
34 integration techniques. In some embodiments and/or usage scenarios, a peripheral portion of the PEs

1 are coupled to I/O FPGAs 420. Example ASICs are illustrated as ASIC 410, comprising a column-
2 organized section of PEs (replicated, e.g., in a one-dimensional fashion to form a wafer), and ASIC
3 411, comprising a square-organized section or a rectangular-organized section of PEs (replicated, e.g.,
4 in a two-dimensional fashion to form a wafer). Other organizations of ASICs on a wafer are
5 contemplated.

6
7 [0403] In some embodiments and/or usage scenarios, neurons associated with layers in a
8 neural network are generally placed on PE 499 elements in a left to right fashion, with earlier layers
9 (e.g., the input layer) on the left and subsequent layers (e.g., the output layer) on the right.
10 Accordingly, data flow during training is illustrated conceptually as dashed-arrows Forward 401,
11 Delta 402, and Chain 403. During Forward 401, stimuli is applied to the input layer and activations
12 from the input layer flow to subsequent layers, eventually reaching the output layer and producing a
13 forward result. During Delta 402, deltas (e.g., differences between the forward result and the training
14 output data) are propagated in the backward direction. During Chain 403, gradients are calculated
15 based on the deltas (e.g., with respect to the weights in the neurons) as they are generated during Delta
16 402. In some embodiments and/or usage scenarios, processing for Delta 402 is substantially
17 overlapped with processing for 403.

18
19 [0404] In some embodiments and/or usage scenarios, Deep Learning Accelerator 400 is an
20 implementation of Deep Learning Accelerator 120 of Fig. 1. In some embodiments and/or usage
21 scenarios, individual PE 499 elements correspond to individual PEs of PEs 122 of Fig. 1. In some
22 embodiments and/or usage scenarios, each ASIC 410 element or alternatively each ASIC 411 element
23 corresponds to all or any portions of PEs of PEs 122 implemented as individual integrated circuits. In
24 some embodiments and/or usage scenarios, each ASIC 410 element or alternatively each ASIC 411
25 element corresponds to (optionally identical) portions of PEs 122 implemented via respective dice of a
26 wafer. In some embodiments and/or usage scenarios, I/O FPGAs 420 elements collectively
27 correspond to FPGAs 121 of Fig. 1.

28
29 [0405] In some embodiments and/or usage scenarios, the placement of neurons (e.g.,
30 associated with layers in a neural network) onto PE 499 elements is performed in whole or in part by
31 all or any portions of Placement Server(s) SW 210 of Fig. 2.

32
33

1 **PROCESSING ELEMENT: COMPUTE ELEMENT AND ROUTER**

2
3 **[0406]** Fig. 5 illustrates selected details of an embodiment of a PE as PE **500** of a deep
4 learning accelerator. PE **500** comprises Router **510** and Compute Element **520**. Router **510**
5 selectively and/or conditionally communicates wavelets between other PEs (e.g., logically adjacent
6 and/or physically adjacent PEs) and the instant PE via couplings **511 – 516**. Router **510** selectively
7 and/or conditionally communicates wavelets to the instant PE via Off Ramp **521** and communicates
8 wavelets from the instant PE via On Ramp **522**. Compute Element **520** performs computations on
9 data embodied in the wavelets according to instruction address information derivable from the
10 wavelets. The instruction address information is used to identify starting addresses of tasks embodied
11 as instructions stored in memory of the compute element.

12
13 **[0407]** In various embodiments, any one or more of **511 – 516** are omitted.

14
15 **[0408]** In some embodiments and/or usage scenarios, PE **500** is an embodiment of PE **499** of
16 Fig. 4, and/or elements of PE **500** correspond to an implementation of PE **499**. In some embodiments
17 and/or usage scenarios, North **513**, East **515**, South **516**, and West **511** correspond respectively to
18 North coupling **430**, East coupling **431**, South coupling **432**, and West coupling **433** of Fig. 4.

19
20 **[0409]** Fig. 6 illustrates selected details of an embodiment a router of a PE, as Router **600**.
21 Consider that there are a plurality of PEs, each comprising a respective router and a respective CE.
22 Router **600** is an instance of one of the respective routers. Router **600** routes wavelets, in accordance
23 with color information of the wavelets and routing configuration information, to the CE of the PE that
24 the instant router is comprised in, as well as others of the routers. The routed wavelets are variously
25 received by the instant router and/or generated by the CE of the PE that the instant router is comprised
26 in. The routing enables communication between the PEs. Stall information is communicated to
27 prevent overflowing of wavelet storage resources in Router **600**.

28
29 **[0410]** Router **600** comprises four groups of interfaces, Data In **610**, Data Out **620**, Stall Out
30 **630**, and Sources **640**. Data In **610**, Data Out **620**, Stall Out **630**, and Sources **640** respectively
31 comprise interface elements **611-617**, **621-627**, **631-637**, and **641-647**. Router **600** further comprises
32 Write Dec **651**, Out **652**, Gen Stall **656**, and Stall **657**, respectively coupled to Data In **610**, Data Out
33 **620**, Stall Out **630**, and Sources **640**. Router **600** further comprises Sources **653** comprising Src **670**

1 coupled to Gen Stall 656. Router 600 further comprises Data Queues 650, Control Info 660, and
2 Router Sched 654. Control Info 660 comprises Dest 661 and Sent 662.

3

4 [0411] Data Queues 650 is coupled to Write Dec 651 to receive incoming wavelet
5 information, and coupled to Out 652 to provide outgoing wavelet information. Data Queues 650 is
6 further coupled to Gen Stall 656 to provide data queue validity information. Router Sched 654 is
7 coupled to Control Info 660 to receive control information relevant to scheduling queued wavelets.
8 Router Sched 654 is further coupled to Stall 657 to receive stall information relevant to scheduling
9 queued wavelets. Router Sched 654 is further coupled to Out 652 to direct presentation of queued
10 wavelets on one or more of 621-627. Router Sched 654 is further coupled to Gen Stall 656 to partially
11 direct generation of stall information.

12

13 [0412] In various embodiments, each of interface elements 611-617, 621-627, 631-637, and
14 641-647 is variously implemented via passive interconnect (e.g., wire(s) without buffering), active
15 interconnect (e.g., wire(s) with selective and/or optional buffering), and coupling with logic to
16 accommodate additional functionality between one instance of Router 600 and another instance of
17 Router 600.

18

19 [0413] In some embodiments and/or usage scenarios, Router 600 is an implementation of
20 Router 510 of Fig. 5.

21

22 [0414] In some embodiments, ones of Data In 610 and ones of Data Out 620 correspond to
23 portions of West 511, Skip West 512, North 513, Skip East 514, East 515, South 516, Off Ramp 521,
24 and On Ramp 522. For example, On Ramp 617 corresponds to On Ramp 522 and Off Ramp 627
25 corresponds to Off Ramp 521. As another example, Y+ 615 comprises the portion of North 513
26 enabled to receive data, and Y+ 625 comprises the portion of North 513 enabled to transmit data.

27

28 [0415] Fig. 7 illustrates selected details of an embodiment of processing associated with a
29 router of a processing element, as Wavelet Ingress 710, Stall Info 720, and Wavelet Egress 730.
30 Conceptually, the router accepts as many wavelets as possible from ingress ports, queuing as
31 necessary and as queue space is available, and routes as many wavelets as possible to egress ports per
32 unit time (e.g., clock cycle). Wavelet Ingress 710 comprises actions 711-713 corresponding to
33 wavelet ingress from (logically and/or physically) adjacent PEs and/or an instant PE, for each
34 respective queue. Stall Info 720 comprises actions 721-723 correspond to providing stall information,

1 for each respective queue. Wavelet Egress 730 comprises actions 731-734 that correspond to wavelet
2 egress to (logically and/or physically) adjacent PEs and/or the instant PE, for each respective queue.
3 In some circumstances, in accordance with color information of a wavelet and routing configuration
4 information, Send Wavelet 734 sends a wavelet from a single queue entry to a single destination (e.g.,
5 unicast). In some circumstances, in accordance with color information of a wavelet and routing
6 configuration information, Send Wavelet 734 sends a wavelet from a single queue entry to a plurality
7 of destinations (e.g., multicast). In various embodiments and/or usage scenarios, any one or more of
8 all or any portions of actions of 710, 720, and/or 730 correspond to actions performed by and/or
9 related to all or any portions of any one or more elements of Router 600 of Fig. 6.

10

11 [0416] Fig. 8 illustrates selected details of an embodiment of a compute element of a
12 processing element, as CE 800.

13

14 [0417] In various embodiments, CE 800 is coupled via Off Ramp 820 and On Ramp 860 to a
15 router. CE 800 comprises Qdistr 824 coupled to receive wavelets via Off Ramp 820. Qdistr 824
16 coupled to transmit wavelets to Scheduling Info 896. Scheduling Info 896 comprises Qs 897, Active
17 Bits 898, and Block Bits 899.

18

19 [0418] In various embodiments, Qs 897 comprises a queue for each fabric color (e.g., to hold
20 wavelets created by other processing elements and associated with the respective color) and each local
21 color (e.g., to hold wavelets created by CE 800 and associated with the respective color), e.g., Q0
22 897.0, ..., and QN 897.N. Each one of Qs 897 (e.g., Q0 897.0) is associated with a respective one of
23 Active Bit 898 (e.g., Active Bit 0 898.0) and Block Bits 899 (e.g., Block Bit 0 899.0). Each one of
24 Active Bits 898 and each one of Block Bits 899 contain information about the respective one of Qs
25 897, e.g., Block Bit N 899.N indicates whether QN 897.N is blocked.

26

27 [0419] In various embodiments, there is variously a physical Q for each color, one or more
28 physical Qs for a predetermined subset of colors, and one or more physical Qs for a dynamically
29 determined subset of colors. In various embodiments, there is variously one or more physical Qs of a
30 same size (e.g., each enabled to hold a same number of wavelets) and one or more physical Qs of
31 differing sizes (e.g., each enabled to hold a different number of wavelets). In various embodiments,
32 there are one or more physical Qs that are variously mapped to virtual Qs, each of the virtual Qs being
33 associated with one or more colors. For example, there are N logical Qs and less than N physical Qs.
34 For another example, some of Qs 897 are enabled to hold 8 wavelets and others of Qs 897 are enabled

1 to hold 3 wavelets. In some embodiments, traffic for one or more colors associated with a particular
2 one of Qs 897 is estimated and/or measured, and the particular one of Qs 897 is enabled to hold a
3 particular number of wavelets based on the traffic.

4
5 [0420] Hash 822 is coupled to Qdistr 824 and selects a physical queue to store a wavelet,
6 based at least in part on the color of the wavelet (e.g., by applying a hash function to the color). In
7 some embodiments, the color associated with a wavelet payload is stored explicitly with the wavelet
8 payload in a queue, such that an entry in the queue holds an entire wavelet (payload with color). In
9 some embodiments, the color associated with a wavelet payload is not stored explicitly with the
10 wavelet payload in a queue, such that an entry in the queue stores a wavelet payload without storing
11 an associated color. The color of the wavelet payload is inferred, such as from the specific queue the
12 wavelet payload is stored in.

13
14 [0421] In some embodiments, one or more of Active Bits 898 and Block Bits 899 are
15 implemented as respective bit vectors with N entries, one entry for each color. In various
16 embodiments, one or more of Active Bits 898 and Block Bits 899 are implemented as respective bit
17 fields in a table comprising one entry for each color.

18
19 [0422] Picker 830 is coupled to Scheduling Info 896, RF 842, Dec 840, Base 890, PC 834, I-
20 Seq 836, and D-Seq 844. Picker 830 is enabled to select a wavelet for processing from one of Qs 897.
21 In some embodiments, Picker 830 selects a wavelet by selecting one of Qs 897, and selecting the
22 oldest wavelet in the selected queue. In some scenarios, Picker 830 selects a new wavelet for
23 processing when Dec 840 signals that a terminate instruction has been decoded. In some other
24 scenarios (e.g., an instruction accessing fabric input), Picker 830 selects a new wavelet for processing
25 from one of Qs 897 in response to a queue identifier received from D-Seq 844.

26
27 [0423] Picker 830 receives the selected wavelet from one of Qs 897 and is enabled to send
28 one or more of data and index from the selected wavelet to RF 842. In some embodiments, Qs 897 is
29 coupled to Data Path 852, and the Data Path is enabled to receive data directly from one of the Qs.
30 Picker 830 is enabled to read a base address from Base 890 and calculate an instruction address to
31 send to PC 834 and I-Seq 836. Base 890 stores a base address and is also coupled to D-Seq 844. PC
32 834 stores the address of the next instruction to fetch. In various embodiments, Base 890 and PC 834
33 are implemented as registers. In some embodiments, D-Seq 844 is enabled to read a base address

1 from Base 890 and request data at one or more addresses from Memory 854 and D-Store 848, based at
2 least in part upon the value read from Base 890.

3

4 [0424] I-Seq 836 is coupled to PC 834 and is enabled to read and modify PC 834 (e.g.,
5 increment for a sequential instruction or non-sequentially for a branch instruction). I-Seq 836 is also
6 coupled to Memory 854 and is enabled to provide an instruction fetch address to Memory 854 (e.g.,
7 based upon PC 834).

8

9 [0425] Memory 854 is further coupled to Dec 840, Data Path 852, and D-Seq 844. In
10 response to an instruction fetch address from I-Seq 836, Memory 854 is enabled to provide
11 instructions located at the instruction fetch address to Dec 840 (an instruction decoder). In various
12 embodiments, Memory 854 is enabled to provide up to three instructions in response to each
13 instruction fetch address. In some embodiments, an instruction is formatted in accordance with one or
14 more of Figs. 25A, 25B, and 25C.

15

16 [0426] Dec 840 is enabled to determine one or more characteristics of instructions, according
17 to various embodiments and/or usage scenarios. For example, Dec 840 is enabled to parse instructions
18 into an opcode (e.g., Opcode 2512 of Fig. 25A) and zero or more operands (e.g., source and/or
19 destination operands). For another example, Dec 840 is enabled to identify an instruction according to
20 instruction type (e.g., a branch instruction, or a multiply-accumulate instruction, and so forth). For yet
21 another example, Dec 840 is enabled to determine that an instruction is a specific instruction and
22 activates one or more signals accordingly.

23

24 [0427] Dec 840 is coupled to Picker 830 via Terminate 812 and is enabled to signal that one
25 of the decoded instructions is a terminate instruction that ends a task (e.g., the last instruction of the
26 instructions executed in response a task initiated in response to the selected wavelet).

27

28 [0428] In some scenarios, Dec 840 is enabled to decode a branch instruction. Examples of
29 branch instructions include: conditional branch instructions that conditionally modify PC 834 and
30 jump instructions that unconditionally modify PC 834. A branch instruction is executed by I-Seq 836
31 and optionally and/or conditionally modifies PC 834. In some scenarios, a branch instruction
32 implements software control flow (e.g., a loop) by conditionally modifying PC 834.

33

1 [0429] In response to decoding an instruction (e.g., a multiply-accumulate instruction), Dec
2 840 is enabled to transmit an opcode to Data Path 852. Dec 840 is coupled to DSRs 846 and enabled
3 to transmit one or more operand identifiers to DSRs 846. Dec 840 is also coupled to D-Seq 844 and
4 enabled to transmit one or more operand type identifiers to D-Seq 844.

5
6 [0430] DSRs 846 comprise registers that hold Data Structure Descriptors (DSDs) and is
7 coupled to and enabled to send one or more DSDs to D-Seq 844. In some embodiments, DSRs
8 comprise source DSRs, destination DSRs, extended DSRs, and stride registers. In response to
9 receiving an operand identifier from Dec 840, DSRs 846 is enabled to read the DSD specified by the
10 operand identifier, and to transmit the DSD to D-Seq 844. In various embodiments, DSRs 846 is
11 enabled to receive up to two source operand identifiers and one destination operand identifier, read
12 two source DSRs and one destination DSR, and transmit two source DSDs and one destination DSD
13 to D-Seq 844. In some embodiments, the CE is enabled to explicitly write a DSD to DSRs from
14 memory in response to load DSR instructions and the CE is enabled to explicitly write a DSD to
15 memory from DSRs in response to store DSR instructions. In some embodiments, DSRs 846 is
16 coupled to and enabled to receive data from and transmit data to Memory 854.

17
18 [0431] In some embodiments, DSRs 846 comprise three sets of DSRs: 12 DSRs for source0
19 operands (sometimes referred to as S0DSRs), 12 DSRs for source1 operands (sometimes referred to as
20 S1DSRs), and 12 DSRs for destination operands (sometimes referred to as DDSRs). In addition,
21 DSRs 846 also comprises six extended DSRs (sometimes referred to as XDSRs) and six stride
22 registers. In some embodiments, DSRs comprise 48 bits, XDSRs comprise 51 bits, and stride registers
23 comprise 15 bits. In various embodiments, respective instructions load 48 bits of data from memory
24 (e.g., D-Store 848 or Memory 854) into respective DSRs (e.g., LDS0WDS, LDS1WDS, and
25 LDDWDS instructions respectively load source0, source1, and destination DSRs). In various
26 embodiments, respective instructions store 48 bits of data from respective DSRs to memory (e.g.,
27 STS0WDS, STS1WDS, and STDWDS instructions respectively store source0, source1, and
28 destination DSRs to memory). In some embodiments, instructions (e.g., LDXDS) load data from
29 memory into XDSRs and other instructions (e.g., STXDS) store data from XDSRs to memory.
30 Instructions that move data between memory and XDSRs (e.g., LDXDS and STXDS) access 64 bits
31 of memory, and only use the lower 51 bits. In some embodiments, instructions (e.g., LDSR) load data
32 from memory into stride registers, and other instructions (e.g., STSR) store data from stride registers
33 to memory. In some embodiments, instructions that move data between memory and stride registers
34 access 16 bits of memory, and only use the lower 15 bits.

1

2 [0432] D-Seq 844 is also coupled to D-Store 848, RF 842, and Picker 830, and is enabled to
3 initiate accessing vector data at various sources in response to DSDs received from DSRs 846. In
4 some scenarios (e.g., in response to receiving a DSD describing one of a 1D memory vector, 4D
5 memory vector, and circular memory buffer), D-Seq 844 is enabled to calculate a sequence of memory
6 addresses to access (e.g., in Memory 854 and/or D-Store 848). In some other scenarios, (e.g., in
7 response to receiving a DSD describing a fabric input), D-Seq 844 is enabled to initiate reading fabric
8 data from one of Qs 897 via Picker 830. In yet other scenarios, (e.g., in response to receiving a DSD
9 describing a fabric output), D-Seq 844 is enabled to initiate transforming data into wavelet(s) and
10 transmitting wavelet(s) to fabric via On Ramp 860. In some embodiments, D-Seq 844 is enabled to
11 simultaneously access vector data at three sources (e.g., read vector data from memory, read vector
12 data from a fabric input, and write vector data to a fabric output).

13

14 [0433] In some embodiments, D-Seq 844 is enabled to access data in one or more registers in
15 RF 842 (e.g., an instruction with one or more input operands and/or one output operand). In some
16 scenarios, D-Seq 844 is enabled to request operands from registers in RF 842. In yet other scenarios,
17 D-Seq 844 is enabled to request data from a register (e.g., an index) in RF 842 as an input for
18 calculating a sequence of memory addresses to access in accordance with a DSD.

19

20 [0434] Data Path 852 is coupled to RF 842 and D-Store 848. In various embodiments, any
21 one or more of Memory 854, RF 842, Qs 897, and D-Store 848 are enabled to provide data to Data
22 Path 852 (e.g., in response to a request from D-Seq 844) and to receive data from Data Path 852 (e.g.,
23 results of operations). Data Path 852 is also coupled via On Ramp 860 to the router, and enabled to
24 send data via On Ramp 860 to the router. Data Path 852 comprises execution resources (e.g., ALUs)
25 enabled to perform operations (e.g., specified by an opcode decoded and/or provided by Dec 840,
26 according to embodiment). In some embodiments, RF 842 comprises sixteen general-purpose
27 registers sometimes referred to as GPR0-GPR15. Each of the GPRs is 16-bits wide and is enabled to
28 store integer or floating-point data.

29

30 [0435] In some embodiments, D-Store 848 is a type of memory that is smaller and more
31 efficient (e.g., lower joules per bit of data read) than Memory 854. In some embodiments, D-Store
32 848 is a type of memory of relatively lower capacity (e.g., retaining less information) and relatively
33 lower access latency and/or relatively higher throughput than Memory 854. In some scenarios, more
34 frequently used data is stored in D-Store 848, while less frequently used data is stored in Memory 854.

1 In some embodiments, D-Store **848** comprises a first address range and Memory **854** comprises a
2 second, non-overlapping address range.

3

4 **[0436]** In some embodiments and/or usage scenarios, elements of the figure correspond to an
5 implementation of Compute Element **520** of Fig. 5, and Off Ramp **820** and On Ramp **860** correspond
6 respectively to Off Ramp **521** and On Ramp **522** of Fig. 5.

7

8 **[0437]** The partitioning and coupling illustrated in Fig. 8 are illustrative only, as other
9 embodiments are contemplated with different partitioning and/or coupling. For example, in other
10 embodiments, RF **842** and DSRs **846** are combined into one module. In yet other embodiments, DSRs
11 **846** and Data Path **852** are coupled.

12

13

14 TASKS

15

16 **[0438]** Fig. 9 illustrates selected details of an embodiment of processing a wavelet for task
17 initiation as flow **900**. Conceptually, the processing comprises initiating a task by determining an
18 address to begin fetching and executing instructions of the task. The address is determined based at
19 least in part on information the wavelet comprises.

20

21 **[0439]** In some embodiments, processing a wavelet for task initiation begins (Start **901**) by
22 selecting a ready wavelet from among, e.g., one or more queues for processing (Select Ready Wavelet
23 for Task Initiation **905**). In some embodiments, the wavelet is selected based upon one or more of:
24 block/unblock state associated with each queue, active/inactive state associated with each queue,
25 color(s) of previously selected wavelets, and a scheduling algorithm.

26

27 **[0440]** After selecting the ready wavelet, the wavelet is checked to determine if the wavelet
28 is a control wavelet or a data wavelet (Control/Data? **908**). If the wavelet is a control wavelet, then a
29 starting address of a task associated with the control wavelet is calculated by adding the lower six bits
30 of the index of the wavelet to a base register (Add Lower Index Bits to Base Register to Form
31 Instruction Address **930**). If the wavelet is not a control wavelet, then the wavelet is a data wavelet.
32 The starting address of a task associated with the data wavelet is calculated by adding the base register
33 to the color of the wavelet multiplied by four (Add (Color * 4) to Base Register to Form Instruction

1 Address 920). The starting address of the task, either as calculated for a control wavelet or as
2 calculated for a data wavelet, corresponds to a starting address of instructions for the task.

3

4 [0441] Once the starting address of the instructions has been calculated, the instructions are
5 fetched from the starting instruction address (Fetch Instructions From Memory at Instruction Address
6 950). One or more of the fetched instructions are decoded and executed (Execute Fetched
7 Instruction(s) 960). Fetching and executing (as illustrated by actions 950 and 960) continue (Not
8 Terminate 961) until a Terminate instruction is executed (Terminate 962), and then processing
9 associated with the initiated task is complete (End 990). In some embodiments, a terminate
10 instruction is the last instruction associated with processing a wavelet. After the initiated task is
11 complete, flow optionally and/or selectively proceeds to process another wavelet for task initiating,
12 beginning with Start 901.

13

14 [0442] According to various usage scenarios, the executing (Execute Fetched Instruction(s)
15 960) comprises executing sequential and/or control-flow instructions, and the instruction address used
16 for fetching varies accordingly (Fetch Instructions From Memory at Instruction Address 950).

17

18 [0443] The ready wavelet selected for task initiation is comprised of a particular color. In
19 some embodiments and/or usage scenarios, once a ready wavelet has been selected for task initiation
20 (Select Ready Wavelet for Task Initiation 905), further wavelets, if any, received of the particular
21 color are consumed as operands for execution of instructions (Execute Fetched Instruction(s) 960).
22 The consuming of the wavelets comprising the particular color as operands continues until fetching
23 and executing of a terminate instruction (Terminate 962).

24

25 [0444] In some embodiments and/or usage scenarios, all or any portions of the actions of
26 flow 900 correspond conceptually to and/or are related conceptually to operations performed by
27 and/or elements of a CE of a PE, e.g., CE 800 of Fig. 8. As an example, Block Bits 899 corresponds
28 to block/unblock state associated with each queue. Active Bits 898 corresponds to active/inactive
29 state associated with each queue. As another example, portions of action 905 are performed by Picker
30 830. Picker 830 selects the oldest wavelet from one of Qs 897 that is ready (e.g., the associated one of
31 Block Bits 899 is not set and the associated one of Active Bits 898 is set), according to a scheduling
32 policy such as round-robin or pick-from-last. The wavelet selected by Picker 830 comprises a color
33 and a wavelet payload formatted in accordance with one of Fig. 13A and Fig. 13B.

34

1 [0445] As another example, action 908 is performed by elements of CE 800. If the control
2 bit of the wavelet payload (e.g., Control Bit 1320 of Fig. 13A) is asserted (determined e.g., by Picker
3 830), then the wavelet is a control wavelet. Subsequently, action 930 is performed by CE 800, such as
4 by Picker 830 adding contents of Base 890 to the 6 lowest bits of Lower Index Bits 1321.1 of Fig.
5 13A to form the instruction fetch address for instructions of the task associated with the control
6 wavelet. Picker 830 then provides the instruction fetch address to PC 834. If the control bit of the
7 wavelet payload (e.g., Control Bit 1320 of Fig. 13A) is deasserted (determined e.g., by Picker 830),
8 then the wavelet is a data wavelet. Subsequently, action 920 is performed by CE 800, such as by
9 Picker 830 adding contents of Base 890 to the color of the wavelet (e.g., corresponding to Color 1324
10 of Fig. 13A and Fig. 13B) multiplied by 4 to form the instruction fetch address for instructions of the
11 task associated with the data wavelet. Picker 830 then provides the instruction fetch address to PC
12 834.

13
14 [0446] As another example, action 950 is performed by elements of CE 800, e.g., PC 834, I-
15 Seq 836, and Memory 854. Action 960 is performed by elements of CE 800, e.g., Dec 840, D-Seq
16 844, Memory 854, RF 842, and Data Path 852, among others. Execution comprises execution of a
17 terminate instruction. An example of a terminate instruction is an instruction with a terminate bit
18 asserted. In the context of the example, when Dec 840 decodes a terminate instruction, Dec 840
19 signals Picker 830 via Terminate 812 that the wavelet is finished, and Picker 830 selects another
20 wavelet for processing, corresponding, e.g., to action 905.

21
22 [0447] In various embodiments and/or usage scenarios, all or any portions of elements of
23 Processing a Wavelet for Task Initiation 900 conceptually correspond to all or any portions of
24 executions of instructions of Task SW on PEs 260 of Fig. 2.

25
26 [0448] In various embodiments and/or usage scenarios, all or any portions of the actions
27 comprising flow 900 conceptually variously correspond to all or any portions of flow 1500 of Fig 15A
28 and/or flow 1550 of Fig. 15B. E.g., action 905 comprises all or any portions of action 1552, and
29 actions 908, 920, 930, 950, and 960 comprise all or any portions of action 1553.

30
31 [0449] Fig. 10 illustrates selected details of an embodiment of instruction processing
32 associated with a compute element of a processing element, as Instruction Processing 1000.

33

1 **[0450]** In some embodiments and/or usage scenarios, all or any portions of the actions of
2 **Instruction Processing 1000** correspond or are related conceptually to operations performed by and/or
3 elements of a CE of a PE, e.g., **CE 800** of Fig. 8.

4
5 **[0451]** Fig. 11 illustrates selected details of an embodiment of flow associated with
6 dependency management via closeouts, as **Dependency Management 1100**.

7
8 **[0452]** In some embodiments and/or usage scenarios, all or any portions of the actions of
9 **Dependency Management 1100** correspond or are related conceptually to operations performed by
10 and/or elements of PEs **122** of Fig. 1. In some embodiments and/or usage scenarios, all or any
11 portions of elements of **Dependency Management 1100** conceptually correspond to all or any portions
12 of executions of instructions of **Task SW** on PEs **260** of Fig. 2.

13
14 **[0453]** Fig. 12 illustrates selected details of an embodiment of flow associated with activation
15 accumulation and closeout, followed by partial sum computation and closeout as **Activation**
16 **Accumulation/Closeout and Partial Sum Computation/Closeout 1200**.

17
18 **[0454]** In some embodiments and/or usage scenarios, all or any portions of the actions of
19 **Activation Accumulation/Closeout and Partial Sum Computation/Closeout 1200** correspond or are
20 related conceptually to operations performed by and/or elements of PEs **122** of Fig. 1. In some
21 embodiments and/or usage scenarios, all or any portions of elements of **Activation**
22 **Accumulation/Closeout and Partial Sum Computation/Closeout 1200** conceptually correspond to all or
23 any portions of executions of instructions of **Task SW** on PEs **260**. In various embodiments and/or
24 usage scenarios, a closeout (e.g., associated with action **1210**) is an example of a control wavelet.

25
26
27 **WAVELETS**

28
29 **[0455]** Fig. 13A illustrates selected details of an embodiment of a sparse wavelet, as **Sparse**
30 **Wavelet 1301**. **Sparse Wavelet 1301** comprises **Sparse Wavelet Payload 1302** and **Color 1324**.
31 **Sparse Wavelet Payload 1302** comprises **Index 1321**, **Sparse Data 1322**, and **Control Bit 1320**. **Index**
32 **1321** comprises **Lower Index Bits 1321.1** and **Upper Index Bits 1321.2**.

33

1 [0456] In some embodiments, Sparse Data 1322 comprises a field for a 16-bit floating-point
2 number or a 16-bit integer number. In various scenarios, Sparse Data 1322 variously represents a
3 weight of a neural network, an input or stimulus of a neural network, an activation of a neural
4 network, or a partial sum of a neural network.

5

6 [0457] In some embodiments, Index 1321 comprises a 16-bit field. In some scenarios, Index
7 1321 is an integer number and is an index that explicitly indicates a specific neuron of a neural
8 network. In some embodiments, Lower Index Bits 1321.1 is 6-bits, and Upper Index Bits 1321.2 is
9 10-bits.

10

11 [0458] In some embodiments, Control Bit 1320 is 1-bit field. In some scenarios, Control Bit
12 1320 indicates whether Sparse Wavelet Payload 1302 triggers control activity or data activity. In
13 some scenarios, control activity comprises computing the last activation of a neuron and data activity
14 comprises computing activations of a neuron that are not the last activation. In some embodiments
15 and/or usage scenarios, the control activity comprises a closeout activity, such as associated with any
16 one or more of Closeout From Prior Layer 1110 and/or Closeout to Next Layer 1122 of Fig. 11, as
17 well as any one or more of Receive Activation Closeout 1204 and/or Transmit Closeout 1210 of Fig.
18 12.

19

20 [0459] In some embodiments, Color 1324 comprises a 5-bit field. In some embodiments, a
21 color corresponds to a virtual channel over a shared physical channel, such as via routing in
22 accordance with the color. In some scenarios, a color is used for a specific purpose such as sending
23 configuration information to processing elements or sending input of a neural network to a neuron that
24 is mapped to a processing element.

25

26 [0460] Fig. 13B illustrates selected details of an embodiment of a dense wavelet, as Dense
27 Wavelet 1331. Dense Wavelet 1331 comprises Dense Wavelet Payload 1332 and Color 1344. Dense
28 Wavelet Payload 1332 comprises Dense Data 1343.1, Dense Data 1343.2, and Control Bit 1340.

29

30 [0461] In some embodiments, Control Bit 1340 is a 1-bit field and is functionally identical to
31 Control Bit 1320.

32

33 [0462] In some embodiments, Color 1344 comprises a 5-bit field and is functionally identical
34 to Color 1324.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

[0463] In some scenarios, Dense Data 1343.1 and Dense Data 1343.2 comprise fields for respective 16-bit floating-point numbers or respective 16-bit integer numbers. In various scenarios, Dense Data 1343.1 and Dense Data 1343.2 variously represent weights of a neural network, inputs or stimuli of a neural network, activations of a neural network, or partial sums of a neural network. In some scenarios, Dense Data 1343.1 and Dense Data 1343.2 collectively comprise a 32-bit floating-point number (e.g., Dense Data 1343.1 comprises a first portion of a 32-bit floating-point number and Dense Data 1343.2 comprises a second portion of a 32-bit floating-point number).

[0464] In various embodiments and/or usage scenarios, usage of sparse wavelets vs. dense wavelets is variously predetermined, dynamically determined, and/or both. In various embodiments and/or usage scenarios, usage of sparse wavelets vs. dense wavelets is determined by software.

[0465] Fig. 14 illustrates selected details of an embodiment of creating and transmitting a wavelet, as Wavelet Creation Flow 1400. Actions of Wavelet Creation Flow 1400 are performed by various agents. A transmitting PE comprises a CE that performs actions 1403-1407, as illustrated by CE of Transmitting PE 1420. The transmitting PE further comprises a router that performs action 1408, as illustrated by Router of Transmitting PE 1430. A receiving PE comprises a router that performs action 1409, as illustrated by Router of Receiving PE 1440.

[0466] Creating and transmitting a wavelet begins (Start 1401) by initializing at least one transmitting PE and one or more receiving PEs, as well as any PEs comprising routers implementing fabric coupling the transmitting PEs and the receiving PEs (Initialize PEs 1402). Each of the PEs comprises a respective router (e.g., Router 510 of Fig. 5) and a respective CE (e.g., Compute Element 520 of Fig. 5). In some scenarios, initializing a PE enables the CE of the PE to perform computations and enables the router of the PE to transmit, receive, and/or forward wavelets over the fabric.

[0467] In various embodiments, a DSR holds a DSD comprising information about an operand such as location of data elements (e.g., memory, fabric input, and/or fabric output), number of the data elements (e.g., length), an address or addresses of the data elements (e.g., start address and stride in memory). For fabric output operands (e.g., wavelets sent via the fabric), the DSR comprises a color for the wavelet(s) on the fabric, a control bit, and optionally a value or location of an index.

1 [0468] In some embodiments, the CE of the transmitting PE configures a source (Set Source
2 1403). In some scenarios, the source is a source DSD describing a source operand. In various
3 embodiments, the source DSD describes one or more data elements stored in one of: cache and
4 memory. In other embodiments, the source DSD describes one or more data elements received via the
5 fabric (e.g., the data elements are payloads of wavelets arriving via the fabric). In some other
6 scenarios, the source comprises a source register (e.g., one of RF 842). In yet other scenarios, the
7 source comprises an immediate specified in an instruction.

8
9 [0469] The CE also configures a destination DSD in a destination DSR describing a fabric
10 destination operand (Set Destination (Fabric) DSR 1404). In some embodiments, the destination DSD
11 describes one or more data elements transmitted via the fabric. In various embodiments, the source
12 and the destination DSDs are configured via one or more instructions.

13
14 [0470] Subsequently, the CE fetches and decodes an instruction (e.g., FMACH, MOV, LT16)
15 comprising a destination operand specified by the DSD in the destination DSR (Fetch/Decode
16 Instruction with Destination DSR 1404.5). In some embodiments, the operand type fields of the
17 instruction specify whether an operand is specified by a DSD.

18
19 [0471] The CE reads the destination DSD from the destination DSR and any source DSDs in
20 source DSRs (Read DSR(s) 1404.6). Based on the DSDs, the CE determines the type of data
21 structure, the source of the data element(s), whether multiple data elements are read together (e.g., for
22 a SIMD operation), and a total number of data elements for each operand. In some scenarios, DSRs
23 are read for one or more of: a source0 operand, a source1 operand, and a destination operand. In some
24 embodiments and/or usage scenarios, the DSRs are read entirely or partially in parallel, and in other
25 embodiments and/or usage scenarios, the DSRs are read entirely or partially sequentially.

26
27 [0472] Then the CE of the transmitting PE reads the data elements described by the source
28 (e.g., a source DSD or a register) and creates a wavelet comprising the data elements based on the
29 destination DSD. The CE reads (e.g., from memory) the first data element(s) specified by the source
30 (Read (Next) Data Element(s) from Queue/Memory 1405). The data element(s) are used to form a
31 wavelet payload. The control bit of the wavelet payload and the color of the wavelet are specified by
32 the destination DSD. The wavelet payload and the color are provided to the router of the transmitting
33 CE (Provide Data Element(s) as Wavelet to Router 1406). In some embodiments and/or usage
34 scenarios, a single data element is used to create the payload of a sparse wavelet. In other

1 embodiments and/or usage scenarios, two data elements are used to create the payload of a dense
2 wavelet.

3

4 **[0473]** The CE of the transmitting PE determines if additional data element(s) are specified
5 by the destination DSD (More Data Elements? 1407). If additional data element(s) are specified by
6 the destination DSD, then the CE creates additional wavelet(s) via actions Read (Next) Source Data
7 Element(s) from Queue/Memory 1405, Provide Data Element(s) as Wavelet to Router 1406, and More
8 Data Elements? 1407 until no additional data element(s) are specified by the destination DSD. If no
9 additional data element(s) are specified by the destination DSD, then flow concludes (End 1410). In
10 some embodiments, the wavelets created via action 1406 are of the same color as specified by the
11 destination DSR.

12

13 **[0474]** The router of the transmitting PE transmits the wavelet(s) formed by the CE of the
14 transmitting PE in accordance with the color of the wavelet(s) (Transmit Wavelet(s) to Fabric 1408),
15 in accordance with respective colors of the wavelets. In some embodiments and/or usage scenarios,
16 the transmitting is directly to the router of the receiving PE. In some embodiments and/or usage
17 scenarios, the transmitting is indirectly to the router of the receiving PE, e.g., via one or more
18 intervening PEs acting to forward the wavelet(s) in accordance with the colors. The router of the
19 receiving PE receives the wavelet(s) in accordance with the color (Receive Wavelet(s) from Fabric
20 1409).

21

22 **[0475]** In various embodiments, action 1408 is performed asynchronously with respect to any
23 one or more of actions 1405, 1406, and 1407. For example, a plurality of wavelets is produced by
24 action 1406 before any of the produced wavelets is transmitted as illustrated by action 1408.

25

26 **[0476]** In various embodiments, Receive Wavelet(s) from Fabric 1409 corresponds in various
27 respects to Receive Wavelet at Router 1503 of Fig. 15.

28

29 **[0477]** In various embodiments and/or usage scenarios, all or any portions of any one or
30 more of elements of Wavelet Creation Flow 1400 correspond conceptually to and/or are related
31 conceptually to operations performed by and/or elements of a PE, e.g., PE 499 of Fig 4.

32

33 **[0478]** In various embodiments and/or usage scenarios, all or any portions of any one or
34 more of elements of Wavelet Creation Flow 1400 (e.g., any one or more of actions 1403-1407)

1 correspond conceptually to and/or are related conceptually to operations performed by and/or elements
2 of a compute element, such as all or any portions of a CE of a PE, e.g., Compute Element 520 of Fig.
3 5 and/or CE 800 of Fig. 8. As an example, the destination DSR (associated with Set DSR Destination
4 (Fabric) DSR 1404) is one of DSRs 846. In some scenarios, the source DSR (associated with Set
5 Source 1403) is one of DSRs 846; in other scenarios the source register (associated with Set Source
6 1403) is one of RF 842.

7

8 [0479] As another example, CE 800 as the CE of the transmitting PE performs action 1403 in
9 response to a load DSR instruction copying information from Memory 854 into the source DSR (e.g.,
10 one of DSRs 846). In various embodiments, the source DSR specifies the location of the data
11 elements as one of Memory 854, D-Store 848, and RF 842. In some scenarios, the source DSR
12 specifies an address of a first data element in Memory 854 (e.g., address 0x0008), a number of data
13 elements (e.g., nine data elements), and a stride between subsequent data elements (e.g., 12 bytes). As
14 another example, CE 800 performs action 1403 by writing data into a register of RF 842.

15

16 [0480] As another example, CE 800 as the CE of the transmitting PE performs action 1404 in
17 response to a load DSR instruction copying information from Memory 854 into the destination DSR
18 (e.g., one of DSRs 846). In various embodiments, the destination DSR specifies transformation of one
19 or more data elements into one or more wavelets and transmitted by Router 510 via a fabric-coupled
20 egress port (e.g., North 513). The destination DSR specifies a color for the wavelet(s), a control bit
21 for the wavelet(s), a number of data elements (e.g., length), and information about an index of the
22 wavelet(s). In some scenarios, the destination DSR specifies the value of the index and in other
23 scenarios the destination DSR specifies a location of the value of the index (e.g., in a register of RF
24 842).

25

26 [0481] As another example, CE 800 as the CE of the transmitting PE performs actions
27 1404.6, 1405, 1406, and 1407 in response to fetching and decoding an instruction specifying a
28 destination DSR as a destination operand (action 1404.5). In some embodiments and/or usage
29 scenarios, D-Seq 844 reads the source DSR and accesses one or two data elements specified by the
30 source DSR, e.g., from Memory 854 or D-Store 848, thereby performing action 1405. In various
31 embodiments, Memory 854 and/or D-Store 848 provide the one or two data elements to Data Path
32 852. The Data Path transforms the data into a wavelet and sends the wavelet via On Ramp 860, e.g.,
33 for storage into an element of Data Queues 650 (of Router 600 of Fig. 6), thereby performing action
34 1406. In some embodiments, On Ramp 860 comprises storage to buffer one or more wavelets. In

1 some embodiments, CE 800 of the transmitting PE reads a color from the destination DSR. Based on
2 the color, CE 800 sends the wavelet payload via On Ramp 860, e.g., for storage into an element of
3 Data Queues 650, thereby completing action 1406. In some embodiments, CE 800 of the transmitting
4 PE performs action 1407 by comparing a number of data elements specified in the destination DSR
5 (e.g., a length) against the number of data elements sent via action 1406 (e.g., tracked by a counter).
6

7 [0482] As another example, CE 800 as the CE of the transmitting PE performs action 1406.
8 The CE transforms the one or two data element(s) into a wavelet payload, according to the destination
9 DSR. In some embodiments and/or usage scenarios, the CE transforms a single data element into a
10 wavelet payload formatted in accordance with Sparse Wavelet 1301 of Fig. 13A. The single data
11 element is transformed into an instantiation of Sparse Data 1322, an index value specified by the
12 destination DSR is transformed into an instantiation of Index 1321, and a control bit from the
13 destination DSR is transformed into an instantiation of Control Bit 1320, thereby forming an
14 instantiation of Sparse Wavelet Payload 1302.
15

16 [0483] As another example, CE 800 as the CE of the transmitting PE transforms two data
17 elements into a wavelet payload formatted in accordance with Dense Wavelet 1331 of Fig. 13B. The
18 first data element is transformed into an instantiation of Dense Data 1343.1 and the second data
19 element is transformed into an instantiation of Dense Data 1343.2. The control bit from the
20 destination DSR is transformed into an instantiation of Control Bit 1340, thereby forming an
21 instantiation of Dense Wavelet Payload 1332.
22

23 [0484] In various embodiments and/or usage scenarios, all or any portions of any one or
24 more of elements of Wavelet Creation Flow 1400 (e.g., any one or more of actions 1408 and 1409)
25 correspond conceptually to and/or are related conceptually to operations performed by and/or elements
26 of a router, such as all or any portions of a router of a PE, e.g., Router 510 of Fig. 5 and/or Router 600
27 of Fig. 6.
28

29 [0485] As an example, Transmit Wavelet(s) to Fabric 1408 is performed by Router 600
30 Router of Transmitting PE 1430 as follows. Router 600 determines the destination(s) of a wavelet in
31 Data Queues 650, e.g., by reading Dest 661. For each color, Dest 661 indicates the output
32 destination(s), e.g., one or more of Data Out 620. Router 600 transmits the wavelet payload and the
33 color (collectively the wavelet) to the fabric, via Out 652 and one or more of Data Out 620. In various

1 embodiments, Router 600 of the transmitting PE performs action 1408 asynchronously with any one
2 or more of actions 1405, 1406, and 1407.

3

4 **[0486]** As another example, Receive Wavelet(s) from Fabric 1409 is performed by Router
5 600 as Router of Receiving PE 1440 as follows. Router 600 receives transmitted wavelet(s) at Data
6 Queues 650 via one of Data In 610 and Write Dec 651. The received wavelet(s) are stored in one or
7 more locations of Data Queues 650.

8

9 **[0487]** In some embodiments and/or usage scenarios, all or any portions of elements of
10 Wavelet Creation Flow 1400 conceptually correspond to all or any portions of executions of
11 instructions of Task SW on PEs 260 of Fig. 2.

12

13 **[0488]** Fig. 15A illustrates selected details of an embodiment of receiving a wavelet as
14 Wavelet Receive Flow 1500. Actions of Wavelet Receive Flow 1500 are performed by various
15 agents. A receiving PE comprises a router performing actions 1503-1506, as illustrated by Router of
16 Receiving PE 1520. The receiving PE further comprises a CE performing action 1507, as illustrated
17 by CE of Receiving PE 1530.

18

19 **[0489]** Receiving a wavelet begins (Start 1501) by initializing at least one transmitting PE
20 and one or more receiving PEs as well any PEs comprising routers implementing fabric coupling the
21 transmitting PEs and the receiving PEs (Initialize PEs 1502). Each of the PEs comprises a respective
22 router (e.g., Router 510 of Fig. 5) and a respective CE (e.g., Compute Element 520 of Fig. 5). In some
23 scenarios, initializing a PE enables the CE of the PE to perform computations and enables the router of
24 the PE to transmit, receive, and/or forward wavelets over the fabric.

25

26 **[0490]** The following description assumes there is a single receiving PE. In usage scenarios
27 where there is plurality of receiving PEs, the respective routers and CEs of each of the receiving PEs
28 perform processing in accordance with Fig. 15A.

29

30 **[0491]** The router of the receiving PE receives a wavelet 'on a color' (e.g., the wavelet
31 comprises the color) of the fabric (Receive Wavelet at Router 1503), as transmitted by the transmitting
32 PE. The router checks the destination(s) of the wavelet based on the color, e.g., by reading a
33 configuration register. If the destination(s) of the wavelet includes other PEs (To Other PE(s)? 1504),
34 then the router transmits the wavelet to the destination PE(s). The router sends the wavelet to

1 output(s) of the router (Transmit Wavelet to Output(s) 1505), and the wavelet is transmitted from the
2 output across the fabric to the destination PE(s). If the destination(s) of the wavelet does not include
3 other PEs, then the transmitting is omitted.

4

5 [0492] If the destination(s) of the wavelet do not include the local CE (For Local CE? 1506),
6 then no further action is taken (End 1510). If one of the destination(s) of the wavelet is the local CE,
7 then the router provides the wavelet to the local CE via the Off Ramp and the wavelet is written into a
8 picker queue associated with the color that the wavelet was received on (Write Wavelet to Picker
9 Queue 1507), thereby receiving the wavelet (End 1510).

10

11 [0493] In various embodiments and/or usage scenarios, all or any portions of any one or
12 more of elements of Wavelet Receive Flow 1500 (e.g., any one or more of actions 1503-1506)
13 correspond conceptually to and/or are related conceptually to operations performed by and/or elements
14 of a router, such as all or any portions of a router of a PE, e.g., Router 510 of Fig. 5 and/or Router 600
15 of Fig. 6.

16

17 [0494] As an example, Receive Wavelet at Router 1503 is performed by Router 600 as
18 Router of Receiving PE 1520 when a wavelet is received on one of Data In 610. Subsequently, To
19 Other PE(s)? 1504 and For Local CE? 1506 are performed by Router 600, using the color of the
20 wavelet to determine the destination(s) of the wavelet, e.g., by reading Dest 661. For each input color,
21 Dest 661 indicates the output destination(s), e.g., one or more of Data Out 620. If Dest 661 indicates
22 that the output includes other PEs (e.g., via one of SkipX+ 621, SkipX- 622, X+ 623, X- 624, Y+ 625,
23 and Y- 626), then the wavelet is sent to other PEs by Router Sched 654. If Dest 661 indicates that the
24 output includes the CE of the PE (e.g., Offramp 627), then the wavelet is sent to the CE by Router
25 Sched 654. The wavelet remains in one of Data Queues 650 until action 1505 is performed by
26 scheduling the wavelet (e.g., by Router Sched 654) to be sent to one or more of Data Out 620.

27

28 [0495] In various embodiments and/or usage scenarios, all or any portions of any one or
29 more of elements of Wavelet Receive Flow 1500 (e.g., action 1507) correspond conceptually to and/or
30 are related conceptually to operations performed by and/or elements of a compute element, such as all
31 or any portions of a CE of a PE, e.g., Compute Element 520 of Fig. 5 and/or CE 800 of Fig. 8. As an
32 example, Write Wavelet to Picker Queue 1507 is performed by sending the wavelet via Off Ramp 820
33 to CE 800 and writing the wavelet into one of Qs 897.

34

1 [0496] In some embodiments and/or usage scenarios, wavelets are received by the router,
2 queued, and routed to router output ports without any specific determination that a wavelet is for a
3 local CE. Instead, wavelets destined for the local CE are routed to the off ramp and are then written
4 into the picker queue. Wavelets not destined for the local CE are routed to other-than the off ramp
5 router outputs.

6

7 [0497] Fig. 15B illustrates selected details of an embodiment of consuming a wavelet as
8 Wavelet Consumption Flow 1550. Actions of Wavelet Consumption Flow 1550 are performed by a
9 CE of a PE.

10

11 [0498] Consuming a wavelet begins (Start 1551) by the picker selecting the wavelet from a
12 queue for processing (Picker Selects Wavelet for Processing 1552), and then the CE processes the
13 wavelet. The CE fetches and executes instructions associated with the wavelet (Fetch, Execute
14 Instructions 1553), thereby consuming the wavelet (End 1554). In some embodiments and/or usage
15 scenarios, fetching and executing instructions associated with the wavelet ends with fetching and
16 executing a terminate instruction.

17

18 [0499] In some embodiments, Picker Selects Wavelet for Processing 1552 is performed by
19 Picker 830 of Fig. 8. In various scenarios, Picker 830 selects one of Qs 897 that is ready (e.g., Block
20 Bits 899 and Active Bits 898 are set to certain values), according to a scheduling policy such as round-
21 robin or pick-from-last. In some embodiments, portions of Wavelet Consumption Flow 1550
22 correspond to portions of Processing a Wavelet for Task Initiation 900 of Fig. 9. As an example,
23 action 1552 corresponds to action 905. As another example, action 1553 corresponds to actions 908,
24 920, 930, 950, and 960.

25

26 [0500] In some other scenarios, the wavelet is accessed as an operand by an instruction (e.g.,
27 FMACH) executing on the CE and the wavelet is consumed by the CE during the execution of the
28 instruction, e.g., as illustrated in Fig. 23.

29

30

31 BLOCK AND UNBLOCK

32

33 [0501] Fig. 16 illustrates selected details of an embodiment of block instruction and unblock
34 instruction execution as flow 1600. Conceptually, executing a block instruction specifying a

1 particular color prevents execution of instructions associated with the particular color at least until
2 execution of an unblock instruction specifying the particular color.

3

4 **[0502]** Referring to the figure, executing an instruction begins (Start 1601) by fetching the
5 instruction from memory and decoding the instruction (Fetch, Decode Instruction 1602). If the
6 instruction decodes to a block instruction (Block Instruction? 1603), then a block operation is
7 performed (Block Color(s) 1604). The source operand of the block instruction specifies one or more
8 colors to block with respect to instruction processing associated with blocked/unblocked colors. In
9 various embodiments and/or usage scenarios, the block operation is performed by setting one or more
10 block indicators to a blocked state for the one or more colors specified by the source operand, and
11 execution is complete (End 1630). In various scenarios, the source operand variously specifies
12 blocking a single color, blocking all colors, and blocking an arbitrary plurality of colors. In
13 subsequent operation, wavelets comprised of colors with respective block indicators set to the blocked
14 state are not selected for processing.

15

16 **[0503]** If the instruction decodes to an unblock instruction (Unblock Instruction? 1610), then
17 an unblock operation is performed (Unblock Color(s) 1611). The source operand of the unblock
18 instruction specifies one or more colors to unblock with respect to instruction processing associated
19 with blocked/unblocked colors. In various embodiments and/or usage scenarios, the unblock
20 operation is performed by resetting a block indicator to an unblocked state for the one or more colors
21 specified by the source operand, and execution is complete (End 1630). In various scenarios, the
22 source operand variously specifies unblocking a single color, unblocking all colors, and unblocking an
23 arbitrary plurality of colors. In subsequent operation, wavelets comprised of colors with respective
24 block indicators set to the unblocked state are selectable for processing.

25

26 **[0504]** If the instruction decodes to an instruction that is not a block instruction and that is
27 not an unblock instruction, then the instruction is otherwise executed (Execute Instruction 1620) and
28 execution is complete (End 1630).

29

30 **[0505]** In some embodiments, if the source operand of a block operation is an immediate
31 (e.g., an 8-bit immediate), then the value of the immediate specifies the color to be blocked. If the
32 source operand is not an immediate, then all colors are blocked.

33

1 [0506] In some embodiments, the source operand of an unblock operation is an immediate
2 (e.g., an 8-bit immediate) and the value of the immediate specifies the color to be unblocked. In
3 various embodiments, an unblock operation with particular operands unblocks multiple colors.
4

5 [0507] In various embodiments and/or usage scenarios, all or any portions of any one or
6 more of elements of Block and Unblock Instruction Processing Flow 1600 correspond conceptually to
7 and/or are related conceptually to operations performed by and/or elements of a compute element,
8 such as all or any portions of a CE of a PE, e.g., Compute Element 520 of Fig. 5 and/or CE 800 of Fig.
9 8.

10
11 [0508] As an example, Block Bits 899 comprise a bit for each color (e.g., as entries in a table,
12 or as a bit-mask). The block operation (Block Color(s) 1604) is performed by setting Block Bits 899
13 to a specific blocked value (e.g., '1') for the one or more colors specified by the source operand. In
14 some embodiments, Picker 830 selects a wavelet for processing from a color where Block Bits 899
15 match an unblocked value (e.g., '0'). As another example, the unblock operation (Unblock Color(s)
16 1611) is performed by setting Block Bits 899 to a specific unblocked value (e.g., '0') for the color
17 specified by the source operand. In some embodiments, Picker 830 selects a wavelet comprising a
18 color where Block Bits 899 match an unblocked value (e.g., '0').
19

20 [0509] In some embodiments, portions of Block and Unblock Instruction Processing Flow
21 1600 correspond to portions of Processing a Wavelet for Task Initiation 900 of Fig. 9. As an example,
22 actions 1602 1603, 1604, 1610, 1611, and 1620 correspond to portions of actions 950 and 960 of Fig.
23 9.
24

25 [0510] In various embodiments and/or usage scenarios, all or any portions of elements of
26 Block and Unblock Instruction Processing Flow 1600 conceptually correspond to all or any portions
27 of executions of instructions of Task SW on PEs 260 of Fig. 2.
28
29

30 NEURON SMEARING

31
32 [0511] Fig. 17 illustrates selected details of an embodiment of a neural network as Neural
33 Network 1700. Network 1700 comprises three portions Input Layer 1710, Internal Layers 1720, and
34 Output Layer 1740. Each layer comprises a plurality of neurons. Input Layer 171, comprises neurons

1 N11 1711, N12 1712, and N13 1713. Internal Layers 1720 comprises a first layer of neurons N21
2 1721, N22 1722, N23 1723, and N24 1724, followed by a second layer of neurons N31 1731, N32
3 1732, and N33 1733. Output Layer 1740 comprises neurons N41 1741 and N42 1742.

4
5 [0512] Selected neurons (N21 1721, N22 1722, N23 1723, and N24 1724 as well as N31
6 1731 and N32 1732) and communications (1791, 1792, and 1793) between the selected neurons are
7 highlighted in the figure. The selected neurons and pathways are discussed in more detail following.

8
9 [0513] Fig. 18A illustrates selected details of a first embodiment of an allocation of
10 processing elements to neurons. Sometimes allocation of processing elements to neurons is referred to
11 as placing neurons in processing elements or alternatively placement of neurons. Like numbered
12 elements of Fig. 18A correspond to like numbered elements of Fig. 17 A first allocation of processing
13 elements to a subset of neurons of Fig. 17 (the highlighted neurons N21 1721, N22 1722, N23 1723,
14 and N24 1724 as well as N31 1731 and N32 1732) is conceptually illustrated. Vertical distance in the
15 figure indicates relative usage of computational resources of each of five processing elements PE0
16 1820, PE1 1821, PE2 1822, PE3 1823, PE4 1824, and PE5 1825.

17
18 [0514] Each of neurons N21 1721, N22 1722, N23 1723, and N24 1724 represents
19 approximately an equal amount of computational resources, e.g., M operations, K storage capacity,
20 and J bandwidth to and from the storage. Each of neurons N31 1731 and N32 1732 represents
21 approximately an equal amount of computational resources, e.g., M/2 operations, K/2 storage, and J/2
22 bandwidth. Thus, each of N31 1731 and N32 1732 represents approximately one half the
23 computational resources of each of N21 1721, N22 1722, N23 1723, and N24 1724. In various
24 embodiments, examples of computational resources comprise compute operations, storage capacity,
25 read bandwidth from storage, write bandwidth to storage, input connections from other neurons, and
26 output connections to other neurons.

27
28 [0515] In the illustrated embodiment, neuron processing is allocated such that each of the
29 foregoing neurons is allocated to an entire PE. More specifically, N21 1721 is allocated to PE0 1840,
30 N22 1722 is allocated to PE1 1841, N23 1723 is allocated to PE2 1842, N24 1724 is allocated to PE3
31 1843, N31 1731 is allocated to PE4 1844, and N32 1732 is allocated to PE5 1845. Therefore, four of
32 the six processing elements are fully subscribed (PE0 1820, PE1 1821, PE2 1822, and PE3 1823),
33 while two of the six processing elements are only one-half subscribed (PE4 1824 and PE5 1825).

34

1 [0516] Fig. 18B illustrates selected details of a second embodiment of an allocation of
2 processing elements to neurons. Like numbered elements of Fig. 18B correspond to like numbered
3 elements of Fig. 17 and Fig. 18A. A second allocation of processing elements to a subset of neurons
4 of Fig. 17 (the highlighted neurons N21 1721, N22 1722, N23 1723, and N24 1724 as well as N31
5 1731 and N32 1732) is conceptually illustrated. As in Fig. 18A, vertical distance in the figure
6 indicates relative usage of computational resources of each of five processing elements PE0 1820, PE1
7 1821, PE2 1822, PE3 1823, PE4 1824, and PE5 1825. Also as in Fig. 18A, each of N31 1731 and
8 N32 1732 represents approximately one half the computational resources of each of N21 1721, N22
9 1722, N23 1723, and N24 1724.

10

11 [0517] In the illustrated embodiment, neuron processing is allocated such that processing for
12 respective neurons is “smeared” across processing elements. Conceptually, neurons are “split” into
13 portions suitable for processing elements to be allocated to. As illustrated in the figure, neurons are
14 split and processing elements allocated so that four of the six processing elements are equally (and
15 fully) subscribed (PE0 1820, PE1 1821, PE2 1822, and PE3 1823), while two of the six processing
16 elements are completely unsubscribed and therefore available for other uses (PE4 1824, and PE5
17 1825). In some embodiments and/or usage scenarios, unsubscribed processing elements remain
18 unused and consume little or no active and/or static power (e.g., via one or more of clock gating and
19 power gating). More specifically, N21 1721 is allocated in two halves (1/2 N21 1721.1 and 1/2 N21
20 1721.2) to two respective processing elements (PE0 1820 and PE2 1822). Similarly, N22 1722 is
21 allocated in two halves (1/2 N22 1722.1 and 1/2 N22 1722.2) to two respective processing elements
22 (PE0 1820 and PE2 1822). N23 1723 is allocated in two halves (1/2 N23 1723.1 and 1/2 N23 1723.2)
23 to two respective processing elements (PE1 1821 and PE3 1823) and N24 1724 is allocated in two
24 halves (1/2 N24 1724.1 and 1/2 N24 1724.2) to two respective processing elements (PE1 1821 and
25 PE3 1823). N31 1731 is allocated in four fourths (1/4 N31 1731.1, 1/4 N31 1731.2, 1/4 N31 1731.3,
26 and 1/4 N31 1731.4) to four respective processing elements (PE0 1820, PE1 1821, PE2 1822, and PE3
27 1823). Similarly, N32 1732 is allocated in four fourths (1/4 N32 1732.1, 1/4 N32 1732.2, 1/4 N32
28 1732.3, and 1/4 N32 1732.4) to four respective processing elements (PE0 1820, PE1 1821, PE2 1822,
29 and PE3 1823). In various embodiments, neurons are split and processing elements allocated based on
30 one or more computational resources associated with the neurons. In some embodiments, neurons are
31 split and processing elements allocated based on the hardware resources available in the processing
32 elements (e.g., some neurons require specific hardware resources such as PRNGs).

33

1 **[0518]** Fig. 19 illustrates selected details of an embodiment of smearing a neuron across a
2 plurality of processing elements. The splitting results in portions of the split neuron that are then
3 smeared across processing elements. Like numbered elements of Fig. 19 correspond to like numbered
4 elements of Fig. 17, Fig. 18A, and Fig. 18B. As illustrated by Fig. 18B, N21 1721 is split into two
5 portions 1/2 N21 1721.1 and 1/2 N21 1721.2 implemented respectively by PE0 1820 and PE2 1822.
6

7 **[0519]** Conceptually, N21 1721 is considered to comprise local compute and local storage, as
8 well as inputs and outputs. Respective elements of N21 1721 are partitioned respectively. The local
9 compute of N21 is partitioned into 1/2 Local Compute 1930.1 and 1/2 Local Compute 1930.2. The
10 local storage of N21 is partitioned into 1/2 Local Storage 1940.1 and 1/2 Local Storage 1940.2. The
11 inputs of N21 are partitioned into a first half in0 1910, in1 1911 and in2 1912 as well as a second half
12 in3 1913, in4 1914, and in5 1915. The outputs of N21 are partitioned into a first half out0 1920, out1
13 1921, out2 1922 as well as a second half out3 1923, out4 1924, and out5 1925.
14

15 **[0520]** 1/2 Local Compute 1930.1, 1/2 Local Storage 1940.1, in0 1910 with in1 1911, and
16 out0 1920 are implemented by PE0 1820. 1/2 Local Compute 1930.2, 1/2 Local Storage 1940.2, in2
17 1912 with in3 1913, and out1 1921 are implemented by PE0 1822.
18

19 **[0521]** In some embodiments and/or usage scenarios, smearing a neuron across more than
20 one processing element is implemented at least in part by additional computation, additional storage,
21 and/or additional communication not otherwise performed/used by the neuron. The additional
22 computation, additional storage, and/or additional communication, enables, e.g., combining partial
23 results from the portions of the neuron into results corresponding to results of the entire neuron.
24 Additional Compute 1950.1 and Additional Storage 1960.1 are representative of additional compute
25 and additional storage for 1/2 N21 1721.1, and are implemented by PE0 1820. Additional Compute
26 1950.2 and Additional Storage 1960.2 are representative of additional compute and additional storage
27 for 1/2 N21 1721.2, and are implemented by PE0 1822.
28

29 **[0522]** Additional Communication 1970 is representative of additional communication
30 between 1/2 N21 1721.1 and 1/2 N21 1721.2, and is implemented by fabric connectivity between PE0
31 1820 and PE0 1822. In some embodiments and/or usage scenarios, all or any portions of Additional
32 Communication 1970 is representative of communications that would occur internally to a single
33 processing element if the single processing element entirely implemented N21 1721.
34

1 [0523] Fig. 20 illustrates selected details of an embodiment of communication between
2 portions of split neurons. Like numbered elements of Fig. 20 correspond to like numbered elements of
3 Fig. 17, Fig. 18A, Fig. 18B, and Fig. 19. Allocations of PE0 1820, PE1 1821, PE2 1822, and PE3
4 1823 to neuron portions are as illustrated by Fig. 18B. For clarity, only allocations specific to PE0
5 1820 and PE1 1821 are illustrated.

6
7 [0524] Wafer Portion 2000 comprises PE0 1820, PE1 1821, PE2 1822, and PE3 1823.
8 Couplings between PEs of Wafer Portion 2000 are illustrated as (coupling between adjacent PEs)
9 2040 coupling PE0 1820 and PE1 1821, 2041 coupling PE1 1821 and PE3 1823, 2043 coupling PE3
10 1823 and PE2 1822, and 2044 coupling PE2 1822 and PE0 1820. Couplings to PEs adjacent to Wafer
11 Portion 2000 are illustrated as (portion of coupling between adjacent PEs) 2050, 2051, 2052, 2053,
12 2054, 2055, 2056, and 2057. The couplings to adjacent PEs are 'portions' since in some embodiments
13 and/or usage scenarios, all or any portions of the couplings are comprised in wafer portions adjacent to
14 Wafer Portion 2000, rather than entirely in Wafer Portion 2000.

15
16 [0525] As a first example, communication portion 1791.1 conceptually represents a portion
17 of communication 1791 between N11 1711 and N21 1721 (of Fig. 17), e.g., from an input layer to an
18 internal layer, with portions of a split neuron in respective processing elements. More specifically,
19 recall that N21 1721 is split into two portions (1/2 N21 1721.1 and 1/2 N21 1721.2; see Fig. 18B).
20 Thus, communication 1791 is split into two portions. Communication portion 1791.1 is illustrative
21 specifically of the portion that is with respect to 1/2 N21 1721.1. Communication portion 1791.1 is
22 transported via (portion of coupling between adjacent PEs) 2057 between a PE adjacent to Wafer
23 Portion 2000 to PE0 1820 (allocated to 1/2 N21 1721.1). In some embodiments and/or usage
24 scenarios, communication 1791 is split into two portions, communication portion 1791.1 (illustrated)
25 and communication portion 1791.2 (not illustrated). In some embodiments and/or usage scenarios,
26 transport of communication portion 1791.1 and communication portion 1791.2 are via a same virtual
27 channel. In some embodiments and/or usage scenarios, transport of communication portion 1791.1
28 and communication portion 1791.2 are via respective unique virtual channels.

29
30 [0526] As a second example, communication portion 1792.1 conceptually represents a
31 portion of communication 1792 between N21 1721 and N31 1731 (of Fig. 17), e.g., from a first
32 internal layer to a second internal layer, with portions of split neurons in respective processing
33 elements. More specifically, recall that N21 1721 is split into two portions (1/2 N21 1721.1 and 1/2
34 N21 1721.2; see Fig. 18B). Further recall that N31 1731 is split into four portions (1/4 N31 1731.1,

1 1/4 N31 1731.2, 1/4 N31 1731.3, and 1/4 N31 1731.4; see Fig. 18B). Thus, communication 1792 is
2 split into portions. Communication portion 1792.1 is illustrative specifically of the portion that is with
3 respect to 1/2 N21 1721.1 and 1/4 N31 1731.2. Communication portion 1792.1 is transported via
4 (coupling between adjacent PEs) 2040 between PE0 1820 (allocated to 1/2 N21 1721.1) and PE1 1821
5 (allocated to 1/4 N31 1731.2). In various embodiments and/or usage scenarios, transport of
6 communication portion 1792.1 (illustrated) and, e.g., other portions (not illustrated) of communication
7 1792 are via a same virtual channel, via unique virtual channels per portion, via virtual channels per
8 portion associated with a particular neuron, and/or via virtual channels per portion associated with a
9 particular processing element.

10
11 [0527] As a third example, communication portion 1793.1 conceptually represents a portion
12 of communication 1793 between N23 1723 and N31 1731 (of Fig. 17), e.g., from a first internal layer
13 to a second internal layer, with portions of split neurons in a same processing element. More
14 specifically, recall that N23 1723 is split into two portions (1/2 N23 1723.1 and 1/2 N23 1723.2); see
15 Fig. 18B). Further recall that N31 1731 is split into four portions (1/4 N31 1731.1, 1/4 N31 1731.2,
16 1/4 N31 1731.3, and 1/4 N31 1731.4; see Fig. 18B). Thus, communication 1793 is split into portions.
17 Communication portion 1793.1 is illustrative specifically of the portion that is with respect to 1/2 N23
18 1723.1 and 1/4 N31 1731.2. Communication portion 1793.1 is transported via one or more
19 mechanisms internal to PE1 1821 (allocated to 1/2 N23 1723.1 and 1/4 N31 1731.2). E.g., PE1 1821
20 uses internal resources (such as a router) to internally feedback an output as an input, and/or to
21 internally provide an input from an output. In some embodiments and/or usage scenarios, transport of
22 communication portion 1793.1 is via a virtual channel that results in an output being used as an input,
23 and/or an input being provided from an output.

24
25 [0528] As a fourth example, communication 2060 conceptually represents all or any portions
26 of Additional Communication 1970 (of Fig. 19), e.g., communications within a neuron that is split
27 across processing elements. More specifically, communication 2060 illustrates specifically
28 communications between two of the four portions that N32 1732 is split into (1/4 N32 1732.1 and 1/4
29 N32 1732.2; see Fig. 18B). Communication 2060 is transported via (coupling between adjacent PEs)
30 2040 between PE0 1820 (allocated to 1/4 N32 1732.1) and PE1 1821 (allocated to 1/4 N32 1732.2).
31 In various embodiments and/or usage scenarios, communication 2060 is via virtual channel dedicated
32 to communication 2060, a virtual channel shared with communication 2060 and communications
33 between other portions of N32 1732, and a virtual channel shared with communication 2060 and all or
34 any portions of neurons split across processing elements.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

[0529] In some embodiments and/or usage scenarios, all or any portion of Wafer Portion **2000** comprises PEs **122** of Fig. 1. In some embodiments and/or usage scenarios, any one of PE0 **1820**, PE1 **1821**, PE2 **1822**, and PE3 **1823** correspond to PE **497** of Fig. 4. In some embodiments and/or usage scenarios, any one or more of coupling between adjacent PEs **2041**, **2042**, **2043**, and **2044** and/or portion of coupling between adjacent PEs **2050**, **2051**, **2052**, **2053**, **2054**, **2055**, **2056**, and **2057** correspond to any one or more of North coupling **430**, East coupling **431**, South coupling **432**, and West coupling **433** of Fig. 4.

[0530] Concepts relating to neuron smearing (e.g., as described with respect to and illustrated by Fig. 17, Fig. 18A, Fig. 18B, Fig. 19, and Fig. 20) are applicable to neural networks of various topologies and types, such as FCNNs, RNNs, CNNs, LSTM networks, autoencoders, deep belief networks, and generative adversarial networks.

[0531] In various embodiments and/or usage scenarios, neurons are split into same-sized portions, e.g., halves, fourths, eights, and so forth. In various embodiments and/or usage scenarios, neurons are split into different-sized portions, e.g., a first portion that is a half, and second and third portions that are respectively each fourths. In various embodiments and/or usage scenarios, neurons are split into arbitrarily-sized portions.

[0532] In various embodiments and/or usage scenarios, a multiplicity of PEs are allocated to a single neuron. In various embodiments and/or usage scenarios, a single PE is allocated to the respective entirety of a multiplicity of neurons.

[0533] In various embodiments and/or usage scenarios, allocation of PEs to neurons is entirely or partially responsive to static and/or dynamic measurements of computational and/or storage requirements. In various embodiments and/or usage scenarios, allocation of PEs to neurons is entirely or partially responsive to dimensionality of data to be processed.

[0534] In various embodiments and/or usage scenarios, dataflow as represented by directions of arrows is unidirectional (as illustrated by drawn arrowhead), bidirectional, and/or reverse-direction (against drawn arrowhead). As a specific example, in various embodiments and/or usage scenarios, communication **1792** (of Fig. 17) is representative of dataflow from N21 **1721** to N31 **1731** (e.g., during forward propagation) or in reverse from N31 **1731** to N21 **1721** (e.g., during back

1 propagation). Thus, communication portion 1792.1 and therefore communication on (portion of
2 coupling between adjacent PEs) 2057 occurs from PE0 1820 to PE1 1821 (e.g., during forward
3 propagation) and in reverse from PE1 1821 to PE0 1820 (e.g., during back propagation).
4
5

6 VECTORS AND DATA STRUCTURE DESCRIPTORS 7

8 [0535] In various embodiments and/or usages scenarios, processing of one or more vectors,
9 each vector comprising respective one or more of data elements, is performed. A vector is variously
10 read from memory (e.g., of a CE of a PE, such as Memory 854 or D-Store 848 of Fig. 8), written to
11 the memory, received from a fabric, or transmitted to the fabric. Vectors read from or written to the
12 memory are sometimes referred to as ‘memory vectors’. Vectors received from or transmitted to the
13 fabric (e.g., as wavelets) are sometimes referred to as ‘fabric vectors’. DSDs from DSRs (as well as
14 XDXDs from XDSRs) are usable to determine addressing patterns for memory vectors and accessing
15 patterns for fabric vectors.
16

17 [0536] Each element identifier in the description of Figs. 21A-E, Figs. 22A-B, and Figs. 23-
18 24 having a first digit of “8” refers to an element of Fig. 8, and for brevity is not otherwise specifically
19 identified as being an element of Fig. 8.
20

21 [0537] Fig. 21A illustrates selected details of an embodiment of a Fabric Input Data Structure
22 Descriptor (aka Fabric Input DSD), as Fabric Input Data Structure Descriptor 2100. In some
23 embodiments, Fabric Input Data Structure Descriptor 2100 describes a fabric vector received by a PE
24 from the fabric, as well as various parameters relating to processing of the fabric vector. In various
25 embodiments and/or usage scenarios, either a source0 operand or a source1 operand of an instruction
26 refers to a DSR containing an instance of a DSD in accordance with Fabric Input Data Structure
27 Descriptor 2100.
28

29 [0538] Fabric Input Data Structure Descriptor 2100 comprises Length 2101, UTID
30 (Microthread Identifier) 2102, UE (Microthread Enable) 2103, SW (SIMD Width) 2104, AC (Activate
31 Color) 2105, Term (Terminate Microthread on Control Wavelet) 2106, CX (Control Wavelet
32 Transform Enable) 2107, US (Microthread Sparse Mode) 2108, Type 2109, SS (Single Step) 2110,
33 SA (Save Address / Conditional Single Step Mode) 2111, SC (Color Specified / Normal Mode) 2112,
34 SQ (Queue Specified / Normal Mode) 2113, and CH (Color High) 2114.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

[0539] In some embodiments, Length **2101** comprises a 15-bit integer specifying the length of the vector, e.g., the number of data elements in the vector.

[0540] In some embodiments, UE (Microthread Enable) **2103** comprises a 1-bit field indicating whether, under at least some conditions, microthreading is enabled during processing of the fabric vector, sometimes referred to as the fabric vector ‘enabling microthreading’. If at least one operand (source or destination) of an instruction is a fabric vector enabling microthreading, then on either an input or output stall during processing of the instruction, processing is enabled to switch (provided sufficient microthreading resource are available) to another instruction of another task. When the stall is cleared, then processing (eventually) returns to the previously stalled instruction. An example input stall is when at least one element of an input fabric vector operands is not available. An example output stall is when there is insufficient space to buffer results associated with an element of an output fabric vector. In some scenarios, a fabric vector that does not enable microthreading is processed synchronously and stalls processing on either an input or output stall. In some scenarios, a fabric vector that enables microthreading is processed asynchronously and reduces or avoids stalling the processing element on either an input or output stall. If a fabric vector enables microthreading, then the processing element is enabled to conditionally switch to processing a different instruction (instead of stalling) and subsequently resume processing the fabric vector at a later point in time (e.g., when data is available).

[0541] In some embodiments, UTID (Microthread Identifier) **2102** comprises a 3-bit field identifying one of a plurality of microthreads and/or resources associated with one of a plurality of microthreads. The microthreads and/or the resources are associated, e.g., with a fabric vector that enables microthreading. In some embodiments, the hardware provides resources for eight microthreads. In some embodiments and/or usage scenarios, UTID **2102** identifies or partially identifies one of Qs **897**.

[0542] In some embodiments, SW (SIMD Width) **2104** comprises a 2-bit field specifying the number of operations (e.g., one, two, or four) that are, in some implementations, executed in parallel. For example, an FMACH, FADDH, FMULH or MOV16 instruction performs multiple (up to four) operations in parallel on respective operands. In some implementation, the SW field is used to determine how to parse wavelets into data versus index information. For example, when the SW field is four, then two wavelets, each having two data values (and no index values) provide four operands,

1 e.g., in parallel. Continuing with the example, when the SW field is two, then a single wavelet having
2 two data values (and no index value) provides two operands, e.g., in parallel. Continuing with the
3 example, when the SW field is one, then a single wavelet having a single data value and a single index
4 value provides a single operand.

5
6 **[0543]** In some embodiments, AC (Activate Color) **2105** comprises a 6-bit field specifying a
7 color to activate (e.g., via an activate operation). In some scenarios, when processing is complete for a
8 fabric vector that enables microthreading, the color specified by the AC field is activated and a task
9 initiated based on the activated color. The completion of processing occurs, e.g., when all elements of
10 the fabric vector have been processed, or when Term **2106** indicates to terminate upon encountering a
11 control wavelet and a control wavelet is encountered while processing the fabric vector. In some
12 embodiments, AC **2105** is enabled to specify one of: a local color and a fabric color.

13
14 **[0544]** In some embodiments, Term (Terminate Microthread on Control Wavelet) **2106**
15 comprises a 1-bit field specifying whether to terminate upon receiving a control wavelet. If the
16 wavelet at the head of the queue specified by Fabric Input Data Structure Descriptor **2100** (e.g., one of
17 Qs **897** as variously specified by various functions of any combination of UTID **2102**, SC **2112**,
18 and/or SQ **2113**, as described elsewhere herein) is a control wavelet (e.g., Control Bit **1320** of Fig.
19 **13A** or Control Bit **1340** of Fig. **13B** is set) and Term **2106** is set, then the instruction is terminated
20 and the color specified by AC **2105** is activated.

21
22 **[0545]** In some embodiments, CX (Control Wavelet Transform Enable) **2107** comprises a 1-
23 bit field specifying whether to transform control wavelets. If CX **2107** is set, then in response to
24 receiving a control wavelet in the fabric vector, bits 15:6 of the index register are set to all "1"s. In
25 some embodiments and/or usage scenarios, if bits 15:6 of the index register are all "1"s, then the
26 control bits of any output wavelets associated with an output fabric vector referencing the index
27 register are set.

28
29 **[0546]** In some embodiments, US (Microthread Sparse Mode) **2108** comprises a 1-bit field
30 specifying whether a fabric vector that enables microthreading (e.g., via the UE field) is processed in a
31 sparse mode. If US **2108** is set, then the fabric vector comprises a vector of sparse data elements and
32 respective wavelet indices of the operand described by Fabric Input Data Structure Descriptor **2100**.
33 The indices are optionally and/or selectively used for address calculation of memory operands,
34 dependent on WLI **2152** (of Fig. **21C**).

1
2 **[0547]** In some embodiments, Type **2109** comprises a 3-bit field specifying a data structure
3 type and/or how to interpret other fields of Fabric Input Data Structure Descriptor **2100**. Type **2109** is
4 “0” for all instances of Fabric Input Data Structure Descriptor **2100**.

5
6 **[0548]** In some embodiments, SS (Single Step) **2110** comprises a 1-bit field specifying
7 whether single step mode operation is enabled, under at least some conditions, for operations using the
8 DSD as an operand. In some scenarios, an instruction with one or more operands that enable single
9 step mode operates in single step mode.

10
11 **[0549]** In some embodiments, SA (Save Address / Conditional Single Step Mode) **2111**
12 comprises a 1-bit field specifying whether save address mode operation is enabled, under at least some
13 conditions, for operations using the DSD as an operand.

14
15 **[0550]** In some embodiments and/or usage scenarios, a color is activated and in response a
16 task is initiated at an address based at least in part on the color. Once initiated, the task executes. In
17 some scenarios, an input fabric vector is provided from the queue associated with the color of the
18 currently executing task. In some embodiments, SC (Color Specified, Normal Mode) **2112** comprises
19 a 1-bit field that if set, specifies that the input fabric vector is provided from a specific queue (e.g., one
20 of Qs **897**) associated with a specific fabric color. The specific fabric color is specified (e.g., as a 5-bit
21 color) as a concatenation of lower bits UTID **2102** (comprising a 3-bit field) and upper bits CH **2114**
22 (comprising a 2-bit field). In some embodiments, SQ (Queue Specified, Normal Mode) **2113**
23 comprises a 1-bit field that if set, specifies that the input fabric vector is provided from a specific
24 queue (e.g., one of Qs **897**). If SQ **2113** is set, then the input fabric vector is provided from the one of
25 Qs **897** specified by UTID **2102**.

26
27 **[0551]** Fig. 21B illustrates selected details of an embodiment of a Fabric Output Data
28 Structure Descriptor (aka Fabric Output DSD), as Fabric Output Data Structure Descriptor **2120**. In
29 some embodiments, Fabric Output Data Structure Descriptor **2120** describes a fabric vector created by
30 a PE and transmitted over the fabric, as well as various parameters relating to processing of the fabric
31 vector. In various embodiments and/or usage scenarios, a destination operand of an instruction refers
32 to a DSR containing an instance of a DSD in accordance with Fabric Output Data Structure Descriptor
33 **2120**.

34

1 [0552] Fabric Output Data Structure Descriptor 2120 comprises Length 2121, UTID
2 (Microthread Identifier) 2122, UE (Microthread Enable) 2123, SW (SIMD Width) 2124, Color 2126,
3 C (Output Control Bit) 2127, Index Low 2128.1, Type 2129, SS (Single Step) 2130, SA (Save
4 Address / Conditional Single Step Mode) 2131, WLI (Wavelet Index Select) 2132, Index High 2128.2,
5 and AC (Activate Color) 2125.

6
7 [0553] In some embodiments, the elements of Fabric Output Data Structure Descriptor 2120
8 (Length 2121, UTID 2122, UE 2123, SW 2124, SS 2130, SA 2131, and AC 2125) are respectively
9 similar in function and/or operation with respect to the elements of Fabric input Data Structure
10 Descriptor 2100 (Length 2101, UTID 2102, UE 2103, SW 2104, SS 2110, SA 2111, and AC 2105).

11
12 [0554] In some embodiments, Color 2126 comprises a 5-bit field specifying the fabric color
13 used to transmit wavelets associated with the fabric vector.

14
15 [0555] In some embodiments, C (Output Control Bit) 2127 comprises a 1-bit field specifying
16 whether a wavelet is a control wavelet. If C 2127 is set, then any wavelets created based on the DSD
17 are control wavelets (e.g., Control Bit 1320 of Fig. 13A is set).

18
19 [0556] In some embodiments, Index Low 2128.1 comprises a 3-bit field and Index High
20 2128.2 comprises a 3-bit field. The concatenation of Index Low 2128.1 and Index High 2128.2 is
21 collectively referred to as Index 2128. In some scenarios, Index 2128 is used to form an index for a
22 wavelet (e.g., Index 1321 of Fig. 13A).

23
24 [0557] In some embodiments, Type 2129 comprises a 3-bit field specifying a data structure
25 type and/or how to interpret other fields of Fabric Output Data Structure Descriptor 2120. Type 2129
26 is "0" for all instances of Fabric Output Data Structure Descriptor 2120.

27
28 [0558] In some embodiments, WLI (Wavelet Index Select) 2132 comprises a 1-bit field
29 specifying in part the index of the fabric vector. In some scenarios, if WLI 2132 is "1", then the index
30 is the value from a register (e.g., GPR4 of RF 842). In some scenarios, if WLI 2132 is "0", then the
31 index is a zero-extension to 16-bits of Index 2128.

32
33 [0559] Fig. 21C illustrates selected details of an embodiment of a 1D Memory Vector Data
34 Structure Descriptor (aka 1D Memory Vector DSD), as 1D Memory Vector Data Structure Descriptor

1 **2140.** In some embodiments, 1D Memory Vector Data Structure Descriptor **2140** describes a one-
2 dimensional memory vector stored in the memory, as well as various parameters relating to processing
3 of the memory vector. In various embodiments and/or usage scenarios, any one or more of a source0
4 operand, a source1 operand, and a destination operand of an instruction refer to respective DSRs
5 containing respective instances of DSDs in accordance with 1D Memory Vector Data Structure
6 Descriptor **2140**.

7

8 **[0560]** 1D Memory Vector Data Structure Descriptor **2140** comprises Length **2141**, Base
9 Address **2142**, Type **2149**, SS (Single Step) **2150**, SA (Save Address / Conditional Single Step Mode)
10 **2151**, WLI (Wavelet Index Select) **2152**, and Stride **2153**.

11

12 **[0561]** In some embodiments, some of the elements of 1D Memory Vector Data Structure
13 Descriptor **2140** (Length **2141**, SS **2150**, and SA **2151**) are respectively similar in function and/or
14 operation with respect to some of the elements of Fabric Input Data Structure Descriptor **2100** (Length
15 **2101**, SS **2110**, and SA **2111**). In some scenarios, if the length of the memory vector is more than 15-
16 bits, then 4D Memory Vector Data Structure Descriptor **2140** is used.

17

18 **[0562]** In some embodiments, Base Address **2142** comprises a 15-bit integer specifying the
19 base address of the memory vector.

20

21 **[0563]** In some embodiments, Type **2149** comprises a 3-bit field specifying a data structure
22 type and/or how to interpret other fields of 1D Memory Vector Data Structure Descriptor **2140**. Type
23 **2149** is “1” for all instances of 1D Memory Vector Data Structure Descriptor **2140**.

24

25 **[0564]** In some embodiments, WLI (Wavelet Index Select) **2152** comprises a 1-bit field
26 specifying in part the index of the vector. If WLI **2152** is “0”, then the index is 0. In some scenarios,
27 if WLI **2152** is “1”, then the index is the value from a register (e.g., GPR4 of RF **842**) or the index of a
28 sparse wavelet (e.g., Index **1321** of Fig. 13A).

29

30 **[0565]** In some embodiments, Stride **2153** comprises a 9-bit signed integer specifying the
31 stride of the vector. In some scenarios, Base Address **2142**, an index specified by WLI **2153**, and
32 Stride **2153** enable calculating addresses of data elements in a 1D memory vector. The address of the
33 first data element in the 1D memory vector is Base Address **2142** + the index specified by WLI **2153**.
34 The address of the next data element in the 1D vector is the address of the first data element + Stride

1 **2153.** For example, Base Address **2142** is 136, WLI **2153** is 1, GPR4 holds the value 6, Stride **2153** is
2 -2, and Length **2141** is 10, then the memory vector comprises data located at addresses {142, 140,
3 138, ..., 124}. In some scenarios, if the stride of the memory vector is more than 9-bits, then 4D
4 Memory Vector Data Structure Descriptor **2140** is used.

5
6 **[0566]** Fig. 21D illustrates selected details of an embodiment of a 4D Memory Vector Data
7 Structure Descriptor (aka 4D Memory Vector DSD), as 4D Memory Vector Data Structure Descriptor
8 **2160**. In some embodiments, 4D Memory Vector Data Structure Descriptor **2160**, in conjunction with
9 4D Memory Vector Extended Data Structure Descriptor **2240** of Fig. 22B, describe a 4-dimensional
10 memory vector stored in the memory, as well as various parameters relating to processing of the
11 memory vector. In some embodiments, 4D Memory Vector Data Structure Descriptor **2160**, in
12 conjunction with 4D Memory Vector Extended Data Structure Descriptor **2240** of Fig. 22B, describe a
13 two-dimensional or three-dimensional memory vector stored in the memory, as well as various
14 parameters relating to processing of the memory vector. In various embodiments and/or usage
15 scenarios, any one or more of a source0 operand, a source1 operand, and a destination operand of an
16 instruction refer to respective DSRs containing respective instances of DSDs in accordance with 4D
17 Memory Vector Data Structure Descriptor **2160**.

18
19 **[0567]** 4D Memory Vector Data Structure Descriptor **2160** comprises Length Lower Bits
20 **2161.1**, Base Address **2162**, Type **2169**, SS (Single Step) **2170**, SA (Save Address / Conditional
21 Single Step Mode) **2171**, WLI (Wavelet Index Select) **2172**, and Length Upper Bits **2161.2**.

22
23 **[0568]** In some embodiments, some of the elements of 4D Memory Vector Data Structure
24 Descriptor **2160** (Base Address **2162**, SS **2170**, SA **2171**, and WLI **2172**) are respectively similar in
25 function and/or operation with respect to 1D Memory Vector Data Structure Descriptor **2140** (Base
26 Address **2142**, SS **2150**, SA **2151**, and WLI **2152**).

27
28 **[0569]** In some embodiments, Lower Bits **2161.1** comprises a 15-bit field and Length Upper
29 Bits **2161.2** comprises a 9-bit field. The concatenation of Lower Bits **2161.1** and Length Upper Bits
30 **2161.2** is collectively referred to (and illustrated as) Length **2161** (a 24-bit field) interpreted in
31 conjunction with 4D Memory Vector Extended Data Structure Descriptor **2240**.

32
33 **[0570]** In some embodiments, Type **2169** comprises a 3-bit field specifying an extended DSR
34 (XDSR), storing, e.g., an extended DSD (XDSD). The XDSD specifies and describes one of: a

1 circular memory buffer (e.g., Circular Memory Buffer Extended Data Structure Descriptor **2210** of
2 Fig. 22A) and a four-dimensional memory vector (e.g., 4D Memory Vector Extended Data Structure
3 Descriptor **2240** of Fig. 22B).

4
5 **[0571]** Fig. 21E illustrates selected details of an embodiment of a Circular Memory Buffer
6 Data Structure Descriptor (aka Circular Memory Buffer DSD), as Circular Memory Buffer Data
7 Structure Descriptor **2180**. In some embodiments, Circular Memory Buffer Data Structure Descriptor
8 **2180**, in conjunction with Circular Memory Buffer Extended Data Structure Descriptor **2210**,
9 describes one of: a circular buffer of data elements stored in the memory and a FIFO of data elements
10 stored in the memory; as well as various parameters relating to processing of the data elements. In
11 various embodiments and/or usage scenarios, any one or more of a source0 operand, a source1
12 operand, and a destination operand of an instruction refer to respective DSRs containing respective
13 instances of DSDs in accordance with Circular Memory Buffer Data Structure Descriptor **2180**.

14
15 **[0572]** Circular Memory Buffer Data Structure Descriptor **2180** comprises Length **2181**,
16 Base Address **2182**, FW (FIFO Wrap Bit) **2188**, Type **2189**, SS (Single Step) **2190**, SA (Save Address
17 / Conditional Single Step Mode) **2191**, WLI (Wavelet Index Select) **2192**, and SW (SIMD Width)
18 **2184**. In some embodiments, a circular memory buffer access always has an index of zero and a stride
19 of one.

20
21 **[0573]** In some embodiments, some of the elements of Circular Memory Buffer Data
22 Structure Descriptor **2180** (Length **2181**, Base Address **2182**, SS **2190**, and SA **2191**) are respectively
23 similar in function and/or operation with respect to some of the elements of 1D Memory Vector Data
24 Structure Descriptor **2140** (Length **2141**, Base Address **2142**, SS **2150**, and SA **2151**). In some
25 embodiments, Type **2189** is similar in function and/or operation to Type **2169** of 4D Memory Vector
26 Data Structure Descriptor **2160**. In some embodiments, SW **2184** of Circular Memory Buffer Data
27 Structure Descriptor **2180** is similar in function and/or operation to SW **2104** of Fabric Input Data
28 Structure Descriptor **2100**.

29
30 **[0574]** In some embodiments, FW (FIFO Wrap Bit) **2188** comprises a 1-bit field enabling
31 distinguishing between a full FIFO and an empty FIFO. FW (FIFO Wrap Bit) **2188** is toggled when
32 an access wraps around the address range of the FIFO.

33
34 **[0575]** In some embodiments, WLI **2192** has no impact on the index of a circular buffer.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

[0576] Fig. 22A illustrates selected details of an embodiment of a Circular Memory Buffer Extended Data Structure Descriptor, as Circular Memory Buffer Extended Data Structure Descriptor 2210. Circular Memory Buffer Extended Data Structure Descriptor 2210 comprises Type 2211, Start Address 2212, End Address 2213, FIFO 2214, Push (Activate) Color 2215, and Pop (Activate) Color 2216.

[0577] In some embodiments, Type 2211 comprises a 1-bit field specifying the type of data structure. Type 2211 is “1” for all instances of Circular Memory Buffer Extended Data Structure Descriptor 2210.

[0578] In some embodiments, Start Address 2212 comprises a 15-bit field specifying the start address of the circular buffer in the memory. In some embodiments, End Address 2213 comprises a 15-bit integer specifying the end address of the circular buffer in the memory. When an address is incremented (e.g., by the stride to initiate the next access) and equals End Address 2213, the address is reset to Base Address 2212, thereby providing circular access behavior.

[0579] In some embodiments, FIFO 2214 comprises a 1-bit field specifying whether the circular buffer is a FIFO. If FIFO 2214 is “0”, then the circular buffer is not a FIFO. If FIFO 2214 is “1”, then the circular buffer is a FIFO.

[0580] In some embodiments, Push (Activate) Color 2215 and Pop (Activate) Color 2216 comprise 6-bit fields specifying colors to activate (e.g., via an activate operation). In some embodiments, Push (Activate) Color 2215 and Pop (Activate) Color 2216 are enabled to specify ones of: a local color and a fabric color.

[0581] In various embodiments, two circular memory buffer DSRs are enabled to describe a FIFO of data elements stored in a same region of the memory. A destination DSR (e.g., DDSR8) describes a write pointer of the FIFO, and a source1 DSR (e.g., S1DSR8) describes a read pointer of the FIFO. In some embodiments, destination and source1 DSRs have a same identifier. In various embodiments, only some of DSRs 846 are enabled to describe FIFOs, (e.g., DDSR8-DDSR11 and S1DSR8-S1DSR11).

1 [0582] FW (FIFO Wrap Bit) 2188 of the two DSRs enables detecting if a FIFO is full or
2 empty. When a FIFO is used as a destination, Base Address 2182 and FW 2188 of the associated
3 S1DSR is read and compared to values from the DDSR. If Base Address 2182 of the two DSRs are
4 the same, but FW 2188 are different, then the FIFO is full. When a FIFO is used as a source, Base
5 Address 2182 and FW 2188 of the associated DDSR are read and compared to values from the
6 S1DSR. If Base Address 2182 of the two DSRs are the same and FW 2188 are the same, then the
7 FIFO is empty. In some scenarios (e.g., microthreading), in response to a read accessing an empty
8 FIFO or a write accessing a full FIFO, processing is switched to an instruction in another task until the
9 FIFO is respectively not empty or not full.

10

11 [0583] Fig. 22B illustrates selected details of an embodiment of a 4D Memory Vector
12 Extended Data Structure Descriptor, as 4D Memory Vector Extended Data Structure Descriptor 2240.
13 In some embodiments, 4D Memory Vector Extended Data Structure Descriptor 2240 partially
14 describes a four-dimensional vector of data elements stored in the memory. 4D Memory Vector
15 Extended Data Structure Descriptor 2240 comprises Type 2241, Dimensions 2242, DF (Dimension
16 Format) 2243, Select Stride 1 2244.1, Select Stride 2 2244.2, Select Stride 3 2244.3, Select Stride 4
17 2244.4, and Stride 2245. In some embodiments, 4D Memory Vector Extended Data Structure
18 Descriptor 2240 comprises 51 bits.

19

20 [0584] In some embodiments, Type 2241 comprises a 1-bit field specifying the type of data
21 structure. Type 2241 is "0" for all instances of 4D Memory Vector Extended Data Structure
22 Descriptor 2240.

23

24 [0585] In some embodiments, Dimensions 2242 comprises a 20-bit field used to initialize the
25 length of the next dimension of the vector.

26

27 [0586] In some embodiments, DF (Dimension Format) 2243 comprises a 5-bit field that, in
28 conjunction with Length 2161 of Fig. 21D, specifies the length of each dimension of the N-
29 dimensional vector. Conceptually, Length 2161 is divided into 6 consecutive 4-bit nibbles and each
30 dimension is expressed using one or more of the nibbles. Bits are set in DF 2243 to indicate
31 demarcations between the dimensions in Length 2161. For example, DF 2242 is "01110" (binary),
32 indicating that the first dimension is expressed using two nibbles, e.g., bits [7:0], and represents a
33 length between 1 and 128. Similarly, the second dimension is expressed using one nibble, e.g., bits
34 [11:8], and represents a length between 1 and 4. An N-dimension vector is represented by setting (N-

1 1) bits in DF **2242**, and only the last dimension uses more than four nibbles. In some embodiments
2 and/or usage scenarios, a one-dimensional vector is described using this format, e.g., if the vector is
3 too long for Length **2141** (of Fig. 21C) to describe. In some embodiments and/or usage scenarios, a
4 two-dimensional or three-dimensional vector is described using this format.

5

6 **[0587]** In some embodiments, Select Stride 1 **2244.1** comprises a 1-bit field specifying a
7 stride for the first dimension of the vector. If Select Stride 1 **2244.1** is “0”, then the stride is 1. If
8 Select Stride 1 **2244.1** is “1”, then the stride is specified by Stride **2245**.

9

10 **[0588]** In some embodiments, Select Stride 2 **2244.2** comprises a 3-bit field and encodes a
11 stride for the second dimension of the vector. If Select Stride 2 **2244.2** is “0”, then the stride is 1. If
12 Select Stride 2 **2244.2** is “1”, then the stride is specified by Stride **2245**. If Stride Select 2 **2244.2** is 2-
13 7, then the stride is specified by a corresponding (DSR) stride register (e.g., of the six stride registers
14 of DSRs **846**.

15

16 **[0589]** In some embodiments, Select Stride 3 **2244.3** and Select Stride 4 **2244.4** comprise
17 respective 3-bit fields. In some embodiments, Select Stride 3 **2244.3** and Select Stride 4 **2244.4** are
18 respectively similar in function and/or operation with respect to the third and fourth dimension as
19 Select Stride 2 **2244.2** is with respect to the second dimension.

20

21 **[0590]** In some embodiments, Stride **2245** comprises a 15-bit field specifying a stride of the
22 vector in the memory. In some scenarios, Stride **2245** enables using a longer stride for a one-
23 dimensional vector than Stride **2153** (of Fig. 21C).

24

25 **[0591]** Fig. 23 illustrates selected details of an embodiment of accessing operands in
26 accordance with data structure descriptors, as Data Structure Descriptor Flow **2300**. In some
27 embodiments, actions of Data Structure Descriptor Flow **2300** are performed by a CE (e.g., CE **800**).

28

29 **[0592]** Accessing a source operand via a data structure descriptor begins (Start **2301**) by
30 initializing one or more DSRs of a CE of a PE with respective DSDs (Set DSR(s) **2302**) and
31 optionally initializing respective XDSDs and/or stride values of the CE ((optional) Set XDSR(s)
32 **2305**). In some embodiments, the initialized DSRs (as well as the optionally initialized XDSDs and
33 stride registers holding the stride values) are initialized by instructions that move data from memory to
34 the DSRs. Subsequently, the CE fetches and decodes an instruction (e.g., FMACH, MOV, or LT16)

1 comprising one or more operands specified by the initialized DSRs and optionally one or more
2 XDSRs and/or stride registers (Fetch/Decode Instruction with DSR(s) 2303). In some embodiments,
3 the operand type fields of the instruction specify whether an operand is specified by a DSR.

4
5 [0593] The CE reads one or more DSDs from the DSRs (Read DSR(s) 2304) and determines
6 one or more of: the type of data structure, the source of the data element(s), whether multiple data
7 elements are read together (e.g., for a SIMD operation), and the total number of data elements for each
8 operand. Depending on the determination, for each DSD read, an XDSR and one or more stride
9 registers are also optionally read ((optional) Read XDSR(s) 2306), as described with respect to Fig.
10 24. In some scenarios, DSRs are read for one or more of: a source0 operand, a source1 operand, and a
11 destination operand, and are identified by respective operand fields of the instruction obtained in
12 action 2303. In some embodiments and/or usage scenarios, any one or more of the DSRs, the XDSRs
13 and the stride registers are read entirely or partially in parallel, and in other embodiments and/or usage
14 scenarios, any one or more of the DSRs, the XDSRs and the stride registers are read entirely or
15 partially sequentially.

16
17 [0594] Based upon the DSDs obtained in action 2304 (and optional XDSRs and stride values
18 obtained in action 2306), the CE reads one or more source data element(s) from the fabric and/or
19 memory (Read (Next) Source Data Element(s) from Queue/Memory 2310). For each source
20 specified by the instruction obtained in action 2303 (e.g., each of source0 and source1), the CE reads
21 sufficient elements for an iteration of the operation specified in the instruction, and in accordance with
22 SIMD width information in the DSDs. Data element(s) from the fabric (e.g., a source data structure is
23 a fabric vector) are accessed via one or more queues of the CE. In some embodiments and/or usage
24 scenarios, the CE also reads data element(s) from registers.

25
26 [0595] After reading the source data element(s), the CE performs the operation using the data
27 element(s) as inputs (Perform (Next) Operation(s) on Data Element(s) 2311). The operation is
28 specified by the instruction obtained in action 2303 (e.g., a multiply-accumulate operation for an
29 FMACH instruction, a move operation for a MOV instruction, or a less than integer comparison for
30 LT16).

31
32 [0596] In some scenarios, the operation (e.g., a multiply-accumulate operation or a move
33 operation) produces one or more output data element(s). The CE writes the output data element(s) to
34 the fabric or the memory (Write (Next) Destination Data Element(s) to Queue/Memory 2312), based

1 upon the DSDs obtained in action **2304** (and optional XDSRs and stride values obtained in action
2 **2306**). Data element(s) sent to the fabric (e.g., the destination data structure is a fabric vector) are
3 formed into wavelets and transmitted to the fabric via the router of the PE. In some other scenarios,
4 there are no output data elements (e.g., some comparison operations).

5

6 **[0597]** After writing any results from the operation, the CE determines if there are additional
7 data element(s) to process (More Data Element(s)? **2313**). In some embodiments, the DSD specifies
8 the total number of data elements to access (e.g., the length of the vector) and the CE compares the
9 number of data element(s) that have been accessed (e.g., tracked via a counter) to the total number of
10 data element(s) specified by the length. If there are additional data element(s) to process, the CE
11 repeats actions **2310-2313** until all data element(s) have been processed and flow concludes (End
12 **2316**).

13

14 **[0598]** In various embodiments and/or usage scenarios, all or any portions of any one or
15 more of elements of Data Structure Descriptor Flow **2300** (e.g., any one or more actions of **2302-**
16 **2312**) correspond conceptually to and/or are related conceptually to operations performed by and/or
17 elements of a CE, e.g., CE **800**.

18

19 **[0599]** As an example, the source DSRs holding source DSDs (associated with Set DSR(s)
20 **2302** and Read DSR(s) **2304**) are one or more of DSRs **846** (e.g., S0DSRs, S1DSRs, DDSRs, XDSRs,
21 and stride registers). In some embodiments, CE **800** performs Set DSR(s) **2302** responsive to
22 instruction(s) that write DSDs into DSRs, e.g., LDS0WDS, LDS1WDS, LDXDS, and LDSR.

23

24 **[0600]** As another example, CE **800** performs Fetch/Decode Instruction with DSR(s) **2303**.
25 In various embodiments, PC **834** and I-Seq **836** fetch instructions from Memory **854** and Dec **840**
26 decodes fetched instructions. In some embodiments, instructions are formatted in accordance with
27 one of: Multiple Operand Instruction **2510** of Fig. 25A, One Source, No Destination Operand
28 Instruction **2520** of Fig. 25B, and Immediate Instruction **2530** of Fig. 25C. In some embodiments,
29 decoding includes detecting that an instruction operand is specified by a DSD, e.g., that the value of
30 Operand 1 Type **2514.1** is "1".

31

32 **[0601]** As another example, CE **800** performs Read DSR(s) **2304** in response to an
33 instruction with one or more operands specified by a DSR. In various embodiments, D-Seq **844** reads
34 the DSR(s) specified by the instruction obtained in action **2303** from DSRs **846**. In some

1 embodiments, DSDs read from the DSRs are formatted in accordance with one or more of: Fabric
2 Input Data Structure Descriptor **2100** of Fig. 21A, Fabric Output Data Structure Descriptor **2200** of
3 Fig. 21B, 1D Memory Vector Data Structure Descriptor **2140** of Fig. 21C, 4D Memory Vector Data
4 Structure Descriptor **2160** of Fig. 21D, and Circular Memory Buffer Data Structure Descriptor **2180** of
5 Fig. 21E. In some embodiments and/or usage scenarios, D-Seq **844**, e.g., responsive to DSDs having
6 Type **2169** or Type **2189** specifying an XDSR, performs (optional) Read XDSR(s) **2306**. In various
7 embodiments, XDSDs read from the XDSRs are formatted in accordance with one of: Circular
8 Memory Extended Buffer Data Structure Descriptor **2180** of Fig. 22A and 4D Memory Vector
9 Extended Data Structure Descriptor **2160** of Fig. 22B.

10

11 **[0602]** As another example, CE **800** performs Read (Next) Source Data Element(s) from
12 Queue/Memory **2310** based upon the source DSD(s) read in action **2304** and optionally XDSD(s) read
13 in action **2306**. In some scenarios, a source DSD specifies (e.g., via Type **2149**) that an operand
14 originates from memory, and D-Seq **844** reads data element(s) from D-Store **848** or Memory **854** at
15 address(es) specified by the DSD (e.g., based in part upon one or more of: Base Address **2142**, WLI
16 **2152**, and Stride **2153**). In some scenarios, a source DSD specifies (e.g., via Type **2109**) that an
17 operand originates from the fabric and CE **800** reads data element(s) from one of Qs **897**. In some
18 embodiments and/or usage scenarios, data elements are directly transmitted from one of Qs **897** to
19 Data Path **852**. In other embodiments and/or usage scenarios, data elements are transmitted from one
20 of Qs **897** to RF **842** and from RF to Data Path **852**. In some embodiments, the one of Qs **897** is
21 implicitly specified by portions of the DSD (e.g., one or more of: UTID **2102**, SC **2112**, and SQ
22 **2113**). In some scenarios, the CE reads from the queue associated with the color of the current task
23 (e.g., the task associated with the instruction obtained in action **2303**). In some scenarios (e.g., SQ
24 **2113** is “1”), the CE reads from a queue specified by UTID **2102**. In some scenarios (e.g., SC **2112** is
25 “1”), the CE reads from a queue associated with the color specified by UTID **2102** concatenated with
26 CH **2114**. In some scenarios, the CE reads one, two, or four data elements from the specified queue
27 based upon SW **2104**.

28

29 **[0603]** In some embodiments and/or usage scenarios, when CE **800** attempts to read more
30 data element(s) than are available in the specified queue of Qs **897**, or alternatively attempts to read
31 from an empty FIFO (e.g., as implemented in accordance with a DSD in accordance with Fig. 21E),
32 then CE **800** stalls. In some embodiments and/or usage scenarios (e.g., microthreading), Picker **830** is
33 enabled to select a different task from Qs **897** while waiting for the data element(s), thereby enabling
34 CE **800** to avoid stalling.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

[0604] As another example, CE 800 performs Perform (Next) Operation(s) on Data Element(s) 2311. In some embodiments, Data Path 852 uses the data element(s) read in action 2310 as inputs to the operation specified by the instruction obtained in action 2303. In some scenarios (e.g., a computational operation), action 2311 produces output data element(s), while in other scenarios (e.g., a comparison operation), action 2311 produces no output data element. In some embodiments, Data Path 852 is enabled to perform more than one operation simultaneously, e.g., performing two or four multiply-accumulate operations simultaneously using SIMD execution resources.

[0605] As another example, CE 800 performs Write (Next) Source Data Element(s) to Queue/Memory 2312 based upon the destination DSD read in action 2304 and optionally XDSD(s) read in action 2306. In some scenarios, the destination DSD specifies (e.g., via Type 2149) that an operand is destined for memory, and D-Seq 844 writes data element(s) to D-Store 848 or Memory 854 at address(es) specified by the destination DSD (e.g., based in part upon one or more of: Base Address 2142, WLI 2152, and Stride 2153).

[0606] In various embodiments and/or usage scenarios, portions of action 2312 (e.g., writing destination data elements to the fabric) correspond conceptually to and/or are related conceptually to Provide Data Element(s) as Wavelet to Router 1406 of Fig. 14. In some scenarios, a destination DSD specifies (e.g., via Type 2129) that an operand is sent to the fabric and CE 800 creates wavelet(s) (e.g., based in part upon Fabric Output Data Structure Descriptor 2120) from the data element(s) and transmits them via On Ramp 860 to Router 600 (of Fig. 6) to the fabric. In some scenarios, the CE transmits one, two, or four data elements as wavelets, based upon SW 2124 of the destination DSD.

[0607] In some embodiments and/or usage scenarios, when CE 800 attempts to transmit more wavelets than resources available in Router 600 (e.g., there are insufficient resources in Data Queues 650 of Fig. 6), or alternatively attempts to write to a full FIFO (e.g., as implemented in accordance with a DSD in accordance with Fig. 21E), then CE 800 stalls. In some embodiments and/or usage scenarios (e.g., microthreading), Picker 830 is enabled to select a different task from Qs 897 while waiting for more resources, thereby enabling CE 800 to avoid stalling.

[0608] As another example, CE 800 performs action 2313. In some embodiments, D-Seq 844 determines how many data element(s) have been processed (e.g., by incrementing a counter for each data element) and compares this against the length of the vector (e.g., Length 2101).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

[0609] Fig. 24 illustrates selected details of an embodiment of decoding a data structure descriptor, as Data Structure Descriptor Decode Flow 2400. In various embodiments and/or usage scenarios, Memory Data Structure Descriptor Flow 2400 is a conceptual representation of all or any portions of actions 2304, 2306, 2310, and 2312 (of Fig. 23) as performed for each DSR describing a fabric or a memory vector. In summary, Fig. 23 illustrates fetching and decoding an instruction comprising one or more operands specified by initialized DSRs, reading the DSRs to obtain and decode corresponding DSDs, reading (next) source data elements in accordance with the DSDs, performing an operation on the source data elements, writing output data elements of the operation in accordance with the DSDs, and iterating back to reading the next source data elements until complete. Fig. 24 illustrates, for fabric vectors (Fabric Vector 2410) and memory vectors (Memory Vector 2420), further details regarding decoding the DSDs obtained from the DSRs, as well as optionally reading one or more XDSRs and stride registers to obtain and decode corresponding XDSDs and stride values, to determine memory access patterns used to access data elements of the memory vectors of the instruction (e.g., any one or more of source0, source1, and destination). Conceptually, the actions illustrated in Fig. 24 are performed for each DSD obtained via action 2304 of Fig. 23. In some embodiments, actions of Memory Data Structure Descriptor Flow 2400 are performed by a CE (e.g., CE 800).

[0610] Decoding a DSD (e.g., as obtained via action 2304 of Fig. 23) begins (Start 2401) by the CE determining whether the DSD corresponds to a fabric vector (Type = Fabric? 2411), e.g., in accordance with Fig. 21A or Fig. 21B. If so, then accesses of the operand described by the DSD proceed as a fabric vector using the DSD (Access via DSD 2412), e.g., if the operand is a source (Fig. 21A), then action 2310 (of Fig. 23) reads from the fabric in accordance with the DSD, and if the operand is a destination (Fig. 21B), then action 2312 (of Fig. 23) writes to the fabric in accordance with the DSD.

[0611] If the DSD does not correspond to a fabric vector, then the DSD corresponds to a memory vector. The CE then determines whether the DSD corresponds to a 1D memory vector (Type = XDSR? 2421), e.g., in accordance with Fig. 21C. If so, then accesses of the operand described by the DSD proceed as a 1D memory vector using the DSD (Access 1D via DSD 2427). E.g., if the operand is a source, then action 2310 reads the source from the memory in accordance with a 1D memory vector described by the DSD, and if the operand is a destination, then action 2312 writes to the memory in accordance with a 1D memory vector described by the DSD. Each iteration of data

1 elements in Fig. 23 (actions **2310-2313**) advances the operand memory addresses in accordance with
2 the 1D memory vector described by the DSD.

3

4 **[0612]** If the DSD does not correspond to a 1D memory vector, then the DSD corresponds to
5 either a 4D memory vector (e.g., in accordance with Fig. 21D) or a circular buffer (e.g., in accordance
6 with Fig. 21E). The CE reads an XDSR specified by the DSD (Read XDSR Specified via DSD **2422**,
7 also conceptually corresponding to (optional) Read XDSR(s) **2306** of Fig. 23) to obtain an XDSD.
8 The XDSR is specified by Type **2169** (of Fig. 21D) or Type **2189** (of Fig. 21E).

9

10 **[0613]** The CE then determines whether the XDSD specifies a 4D memory vector (e.g., in
11 accordance with Fig. 22B). If so, then the CE optionally reads one or more stride registers
12 ((optionally) Read Stride Register(s) **2424**, also conceptually corresponding to (optional) Read
13 XDSR(s) **2306** of Fig. 23), as optionally specified by the XDSD. Accesses of the operand described
14 by the DSD, the XDSD, and any optional stride values (obtained from the stride registers) proceed as
15 a 4D memory vector using the DSD, the XDSD, and the optional stride values (Access 4D via XDSD
16 **2428**). E.g., if the operand is a source, then action **2310** reads the source from the memory in
17 accordance with the 4D memory vector, and if the operand is a destination, then action **2312** writes to
18 the memory in accordance with the 4D memory vector. Each iteration of data elements in Fig. 23
19 (actions **2310-2313**) advances the operand memory addresses in accordance with the 4D memory
20 vector described by the DSD.

21

22 **[0614]** If the XDSD does not correspond to a 4D memory vector, then the XDSD
23 corresponds to a circular buffer (e.g., in accordance with Fig. 22A). Accesses of the operand
24 described by the DSD and the XDSD proceed as a circular buffer using the DSD and the XDSD
25 (Access Circular Buffer via XDSD **2429**). E.g., if the operand is a source, then action **2310** reads the
26 source from the memory in accordance with the circular buffer, and if the operand is a destination,
27 then action **2312** writes to the memory in accordance with the circular buffer. Each iteration of data
28 elements in Fig. 23 (actions **2310-2313**) advances the operand memory addresses in accordance with
29 the circular buffer described by the DSD.

30

31 **[0615]** In various embodiments, D-Seq **844** performs Type = Fabric? **2411** and/or Type =
32 XDSD? **2421** based upon a DSD read in action **2304** (of Fig. 23). In some embodiments, a type field
33 of the DSD (e.g., Type **2109** of Fig. 21A, Type **2129** of Fig. 21B, Type **2149** of Fig. 21C, Type **2169**
34 of Fig. 21D, and Type **2189** of Fig. 21E) determines if the data structure is one of: a fabric vector (e.g.,

1 the Type = "0"), a 1D vector (e.g., the Type = "1"), and an XDSD type (e.g., the Type = "2-7"). In
2 various embodiments (e.g., the Type = "2-7"), the value of the type field specifies which XDSR of
3 DSRs 846 to read for action 2422. In some embodiments, D-Seq 844 performs action 2422 and
4 receives the XDSD from DSRs 846. In some other embodiments, DSRs 846 performs actions 2421
5 and 2422 and transmits the DSD and the XDSD to D-Seq 844.

6
7 [0616] As another example, D-Seq 844 performs Type = 4D Vector? 2423 based upon the
8 XDSD of action 2422. In some embodiments, the type field of the XDSD (e.g., Type 2211 of Fig.
9 22A or Type 2241 of Fig. 22B) read from the XDSR determines if the data structure is one of a 4D
10 vector (e.g., the XDSD Type = "0") and a circular buffer (the XDSD Type = "1").

11
12 [0617] As another example, D-Seq 844 generates memory access(es) in accordance with
13 action 2427 by computing the memory address(es) based upon the DSD (e.g., of action 2304), using
14 e.g., Base Address 2142, WLI 2152, Length 2141, and Stride 2153 of the DSD, as described
15 elsewhere herein. Similarly, D-Seq 844 generates memory access(es) in accordance with action 2428
16 by computing the memory address(es) based upon the DSD (e.g., of action 2404) and XDSD of action
17 2422 using e.g., Base Address 2162, Length 2161, WLI 2172, Stride 2245, Stride Select 1 2244.1, and
18 DF 2243 of the DSD and the XDSD, as described elsewhere herein. Similarly, D-Seq 844 generates
19 memory access(es) in accordance with action 2429 by computing the memory address(es) based upon
20 the DSD (e.g., of action 2404) and XDSD of action 2422 using e.g., Base Address 2182, Length 2181,
21 WLI 2192, Start Address 2212, and End Address 2213 of the DSD and the XDSD, as described
22 elsewhere herein.

23
24 [0618] In some embodiments, D-Seq 844 sends each computed address to one of D-Store 848
25 and Memory 854. In response to receiving a computed address, the D-Store and/or the Memory
26 accesses two bytes of data at the computed address.

27 28 29 INSTRUCTION FORMATS

30
31 [0619] Each element identifier in the description of Figs. 25A-C having a first digit of "8"
32 refers to an element of Fig. 8, and for brevity is not otherwise specifically identified as being an
33 element of Fig. 8.

34

1 **[0620]** Fig. 25A illustrates selected details of an embodiment of a multiple operand
 2 instruction, as Multiple Operand Instruction 2510. Multiple Operand Instruction 2510 is one of: a
 3 two/three source, one destination operand instruction (e.g., a multiply-add such as FMACH), a two
 4 source, no destination operand instruction (e.g., a comparison such as LT16), and a one source, one
 5 destination operand instruction (e.g., a move instruction such as MOV16).

6
 7 **[0621]** Multiple Operand Instruction 2510 comprises various fields: Instruction Type 2511,
 8 Opcode 2512, Operand 0 Encoding 2513, Operand 1 Encoding 2514, and Terminate 2515. Operand 0
 9 Encoding 2513 comprises Operand 0 Type 2513.1 and Operand 0 2513.2. Operand 1 Encoding 2514
 10 comprises Operand 1 Type 2514.1 and Operand 1 2514.2. In some embodiments, Multiple Operand
 11 Instruction 2510 comprises 20 bits.

12
 13 **[0622]** In some embodiments, the value of Instruction Type 2511 distinguishes between
 14 different types of instructions (e.g., two/three source, one destination and one source, and one
 15 destination instruction types) according to the table following. In various embodiments, the value of
 16 Opcode 2512 specifies a particular operation (e.g., multiply, add, or subtract). The length of Opcode
 17 2512 varies between different types of instructions as described in the table following.

<u>Instruction Family</u>	<u>Value of Instruction Type 2511</u>	<u>Length of Opcode 2522</u>
Two/three source, one destination	10	5 bits
Two source, no destination	1110	4 bits
One source, one destination	110	5 bits

18
 19 **[0623]** In some embodiments, Operand 0 Encoding 2513 describes a source and/or
 20 destination operand, according to the table following. In some embodiments, Operand 1 Encoding
 21 2714 describes a source operand.

<u>Instruction Family</u>	<u>Operand 0 Encoding 2513</u>	<u>Operand 1 Encoding 2514</u>
Two/three source, one destination	Source0 and destination	Source1
Two source, no destination	Source0	Source1
One source, one destination	Destination	Source1

22
 23 **[0624]** In some embodiments, Operand 0 2513.2 and Operand 1 2514.2 comprise respective
 24 4-bit fields. In some embodiments, Operand 0 Type 2513.1 and Operand 1 Type 2514.1 comprise
 25 respective 2-bit fields and respectively determine how to interpret Operand 0 2513.2 and Operand 1

1 **2514.2.** For a two/three source operand, one destination operand instruction, Operand 0 Type **2513.1**
 2 is interpreted according to the table following.

<u>Value of 2513.1</u>	<u>Operand 0 Encoding 2513</u>
0	Source0 is S0DSR[Operand 0 2513.2], destination is S0DSR[Operand 0 2513.1]
1	Source0 is S0DSR[Operand 0 2513.2], destination is DDSR[Operand 0 2513.1]
2	Source0 is GPR[Operand 0 2513.2], destination is GPR[Operand 0 2513.1]
3	Source0 is GPR[Operand 0 2513.2], destination is DDSR[Operand 0 2513.1] if Operand 1 Type 2514.1 is 0, destination is GPR[0] otherwise

3
 4 **[0625]** For example, if the value of Operand 0 Type **2513.1** is “1” and the value of Operand 0
 5 **2513.2** is “4”, then Operand 0 Encoding **2513** specifies that the source0 operand is a vector described
 6 by S0DSR[4] and the destination operand is a vector described by DDSR[4].

7
 8 **[0626]** For a two source operand, no destination operand instruction, Operand 0 Type **2513.1**
 9 is interpreted according to the table following.

<u>Value of 2513.1</u>	<u>Operand 0 Encoding 2513</u>
0	Source0 is S0DSR[Operand 0 2513.2]
1	Source0 is GPR[Operand 0 2513.2]

11
 12 **[0627]** For example, if the value of Operand 0 Type **2513.1** is “0” and the value of Operand 0
 13 **2513.2** is “4”, then Operand 0 Encoding **2513** specifies that the source0 operand is a vector described
 14 by S0DSR[4].

15
 16 **[0628]** For a one source operand, one destination operand instruction, Operand 0 Type
 17 **2513.1** is interpreted according to the table following.

<u>Value of 2513.1</u>	<u>Operand 0 Encoding 2513</u>
0	Destination is DDSR[Operand 0 2513.2]
1	Destination is GPR[Operand 0 2513.2]

18
 19 **[0629]** For example, if the value of Operand 0 Type **2513.1** is “0” and the value of Operand 0
 20 **2513.2** is “4”, then Operand 0 Encoding **2513** specifies that the destination operand is a vector
 21 described by DDSR[4].

22

1 **[0630]** For Multiple Operand Instruction **2510**, Operand 1 Type **2514.1** is interpreted
 2 according to the table following.

<u>Value of 2514.1</u>	<u>Operand 1 Encoding 2514</u>
0	Source1 is S1DSR[Operand 1 2514.2]
1	Source1 is the data in memory at the address specified by GPR[6]
2	Source1 is GPR[Operand 1 2514.2]
3	Source1 is an immediate

3
 4 **[0631]** For example, if the value of Operand 0 Type **2513.1** is “0” and the value of Operand 0
 5 **2513.2** is “4”, then Operand 0 Encoding **2513** specifies that the destination operand is a vector
 6 described by DDSR[4].
 7

8 **[0632]** In various embodiments, a source1 operand that is an immediate specifies one of:
 9 several predetermined values (e.g., 0, 1, and -1) and a pseudo-random number generated by an LFSR.
 10 For example, if the value of Operand 1 Type **2514.1** is “3” and the value of Operand 1 **2514.2** is “8”,
 11 then Operand 1 Encoding **2514** specifies a PRNG generated by an LFSR.
 12

13 **[0633]** In some embodiments, Terminate **2515** comprises a 1-bit field specifying that the
 14 instruction is the last instruction in a task. When the instruction finishes execution, the task is
 15 terminated, enabling selection and execution of a new task (e.g., via Terminate **812** and Picker **830**).
 16

17 **[0634]** Fig. 25B illustrates selected details of an embodiment of a one source, no destination
 18 operand instruction, as One Source, No Destination Instruction **2520**. One Source, No Destination
 19 Instruction **2520** comprises Instruction Type **2521**, Opcode **2522**, Operand 1 Encoding **2523**,
 20 Immediate High **2524**, and Terminate **2525**. Operand 1 Encoding **2523** describes a source operand
 21 and comprises Operand 1 Type **2523.1** and Operand 1 **2523.2**. In some embodiments, One Source,
 22 No Destination Instruction **2520** comprises 20 bits.
 23

24 **[0635]** In some embodiments, Instruction Type **2521** comprises four bits, “1111”, specifying
 25 that the instruction is a one source, no destination operand instruction, and Opcode **2522** comprises a
 26 4-bit field specifying a particular operation (e.g., block, unblock, activate, set active PRNG, data filter,
 27 conditional branch, and jump).
 28

1 [0636] In some embodiments, Immediate High 2524 comprises a 4-bit field. In some
2 scenarios, Immediate High 2524 concatenated with Operand 1 2523.2 forms an 8-bit immediate.

3

4 [0637] In some embodiments, Operand 1 Type 2523.1 comprises a 2-bit field that determines
5 how Operand 1 2523.2 is interpreted. If Operand 1 Type 2523.1 is “0”, then Operand 1 Encoding
6 2523 specifies a vector (e.g., a fabric vector of data elements from Qs 897, or a memory vector of data
7 elements in one of Memory 854 and D-Store 854) and the value of Operand 1 2523.2 identifies which
8 one of the 12 S1DSRs of DSRs 846 describe the vector. If Operand 1 Type 2523.1 is “1”, then
9 Operand 1 Encoding 2523 describes a value in memory (e.g., one of Memory 854 and D-Store 848) at
10 an 8-bit address formed by a concatenation of Immediate High 2524 with Operand 1 2523.2. If
11 Operand 1 Type 2523.1 is “2”, then Operand 1 Encoding 2523 describes a value in a register (e.g., one
12 of RF 842) identified by the value of Operand 1 2523.2. If Operand 1 Type 2523.1 is “3”, then
13 Operand 1 Encoding 2523 describes an immediate. If Opcode 2522 specifies an operation (e.g., block,
14 unblock, or activate) that operates on 16-bit integer operands, then the immediate comprises eight bits
15 and is a concatenation of Immediate High 2524 and Operand 1 2523.2.

16

17 [0638] In some embodiments, Terminate 2525 comprises a 1-bit field specifying that the
18 instruction is the last instruction in a task. When the instruction finishes execution, the task is
19 terminated, enabling selection and execution of a new task (e.g., via Terminate 812 and Picker 830. If
20 One Source, No Destination Instruction 2520 is a conditional branch, then the task is only terminated
21 if the conditional branch is not taken.

22

23 [0639] Fig. 25C illustrates selected details of an embodiment of an immediate instruction, as
24 Immediate Instruction 2530. Immediate Instruction 2530 comprises Instruction Type 2531, Opcode
25 2532, Operand 0 2533.2, and Immediate 2534. In some embodiments, Immediate Low 2534.1
26 comprises a 9-bit field and Immediate High 2534.2 comprises a 1-bit field. The concatenation of
27 Immediate Low 2534.1 and Immediate High 2534.2 is collectively referred to (and illustrated as) as
28 Immediate 2534. In some embodiments, Immediate Instruction 2520 comprises 20 bits.

29

30 [0640] In some embodiments, Instruction Type 2531 comprises a 1-bit field, “0”, specifying
31 that the instruction is an immediate instruction, and Opcode 2532 comprises a 5-bit field specifying a
32 particular operation (e.g., load source0 DSR, load source1 DSR, load destination DSR, store source0
33 DSR, store source1 DSR, and store destination DSR). In some scenarios, execution of an Immediate
34 Instruction 2530 (e.g., a load DSR instruction, and a load XDSR instruction) loads data from one of

1 Memory 854 and D-Store 848 to a DSR of DSRs 846. In other scenarios, execution of an Immediate
2 Instruction 2530 (e.g., a store DSR instruction, and a store XDSR instruction) stores data from a DSR
3 of DSRs 846 to one of Memory 854 and D-Store 848.

4
5 [0641] In some embodiments, Operand 0 2533.2 comprises a 4-bit field and Opcode 2532
6 determines how Operand 0 2533.2 is interpreted. In some scenarios (e.g., if Operand 0 2533.2
7 specifies an operation without a register operand such as a jump operation), Immediate Low 2534.1,
8 Operand 0 2533.2, and Immediate High 2534.2 are concatenated to form a 14-bit immediate. In some
9 other scenarios, Immediate 2534 is sign extended to form a 16-bit immediate. In yet other scenarios,
10 Immediate 2534 is sign extended to form a 15-bit address. In yet other scenarios, Immediate 2534 is
11 shifted one bit to the left and sign extended to form a 15-bit address (e.g., for 32-bit data).

12 13 14 DEEP LEARNING ACCELERATOR EXAMPLE USES

15
16 [0642] In various embodiments and/or usage scenarios, as described elsewhere herein, a deep
17 learning accelerator, such as a fabric of PEs (e.g., as implemented via wafer-scale integration and as
18 illustrated, for example, in Fig. 4) is usable to train a neural network, and/or to perform inferences
19 with respect to a trained neural network. The training, in some circumstances, comprises determining
20 weights of the neural network in response to training stimuli. Various techniques are usable for the
21 training, such as Stochastic Gradient Descent (SGD), Mini-Batch Gradient Descent (MBGD),
22 Continuous Propagation Gradient Descent (CPGD), and Reverse CheckPoint (RCP). Following,
23 CPGD is contrasted with other techniques, and then each of SGD, MBGD, CPGD, and RCP are
24 described in more detail.

25
26 [0643] Past deep neural network training approaches (e.g., SGD and MBGD) have used so-
27 called anchored-delta learning. That is, the delta derived weight updates have been 'anchored' or held
28 fixed until processing of all activations for a training set batch or a mini-batch are completed. In some
29 circumstances, the layer-sequential nature of anchored-delta learning resulted in high-latency
30 sequential parameter updates (including for example, weight updates), which in turn led to slow
31 convergence. In some circumstances, anchored-delta learning has limited layer-parallelism and thus
32 limited concurrency.

1 [0644] In contrast, in some circumstances, use of a continuous propagation (aka immediate-
2 delta) learning rule for deep neural network training, as taught herein, provides faster convergence,
3 decreases the latency of parameter updates, and increases concurrency by enabling layer-parallelism.
4 Deltas computed from the immediate network parameters use updated information corresponding to
5 the current parameter slope. Continuous propagation enables layer parallelism by enabling each layer
6 to learn concurrently with others without explicit synchronization. As a result, parallelization along
7 the depth of a network enables more computing resources to be applied to training. Parallelism
8 available in continuous propagation realizes up to a 10x wall clock time improvement, as compared to
9 MBGD techniques, in some usage scenarios. The continuous propagation approach also enables
10 avoiding using extra memory to store the model parameter values for multiple vectors of activations.

11

12 [0645] In some embodiments and/or usage scenarios, a neural network is trained using
13 continuous propagation of stimuli to perform SGD. In some embodiments of training via CPGD, RCP
14 enables reducing the number of activations held in memory (thus reducing the memory footprint) by
15 recomputing selected activations. In some scenarios, recomputing activations also improves the
16 accuracy of the training estimates for the weights. In training without RCP, every layer of neurons
17 receives activations during one or more forward passes, and saves the activations to re-use for
18 computations performed during the one or more backward passes associated with the forward passes
19 (e.g., the one or more delta, chain, and weight update passes associated with the forward passes). In
20 some scenarios (e.g., relatively deep neural networks), the time between saving the activations and the
21 associated backward pass is relatively long and saving all activations uses relatively more memory
22 than saving fewer than all the activations.

23

24 [0646] For example, only some of the layers of neurons (e.g., every even layer) save the
25 respective activations and the other layers discard the respective activations (e.g., every odd layer).
26 The layers with saved activations (e.g., every even layer) use the most recent weights to recompute
27 and transmit the recomputed activations to the layers that discarded activations (e.g., every odd layer).
28 In some scenarios, the recomputed activations differ from the discarded activations because the most
29 recent weights are different from the weights that were available during the forward pass (e.g., one or
30 more weight updates occurred between the forward pass and the associated backward pass). In
31 various embodiments, the number and type of layers that save and discard activations is selected to
32 optimize for the desired balance of reduced memory usage and increased computation. As one
33 example, every fourth layer saves activations and all other layers discard activations. As another

1 example, convolutional layers are selected to save activations and other layers are selected to discard
2 activations.

3

4 **[0647]** In various embodiments and/or usage scenarios, any one or more of SGD, MBGD,
5 and CPGD, with or without RCP, are implemented via one or more of: a fabric of processing elements
6 (e.g., as illustrated in Fig. 4), one or more GPUs, one or more CPUs, one or more DSPs, one or more
7 FPGAs, and one or more ASICs.

8

9 **[0648]** SGD, e.g., with back-propagation, is usable (as described elsewhere herein) for
10 training a neural network. However, learning via gradient descent is inherently sequential, because
11 each weight update uses information from a gradient measurement made after completion of a full
12 forward pass through the neural network. Further, weight updates are made during a corresponding
13 backward pass through the neural network (following and corresponding to the forward pass), and
14 thus the last weight update occurs after completion of the entire corresponding backward pass.

15

16 **[0649]** MBGD enables more parallelism than SGD by gradient averaging over a mini-batch,
17 processing several (a 'mini-batch' of) activations in parallel. However, speed of sequential updates,
18 compared to SGD, is unchanged, and weight updates, as in SGD, are completed after completion of all
19 corresponding backward passes through the neural network. As mini-batch size increases by
20 processing more activations in parallel, gradient noise is reduced. Beyond a point the reduction in
21 gradient noise, in some scenarios, results in poor generalization.

22

23 **[0650]** CPGD enables parallel processing and updating of weights in all layers of a neural
24 network, while activations propagate through the layers in a stream. Thus CPGD overcomes, in some
25 embodiments and/or usage scenarios, sequential processing limitations of SGD and MBGD.

26

27 **[0651]** RCP enables reduced memory usage via (re)computing activations that would
28 otherwise be stored, and is usable in combination with SGD, MBGD, and CPGD.

29

30 **[0652]** Pipeline flow diagrams are usable to compare and contrast various SGD, MBGD,
31 CPGD, and CPGD with RCP techniques. Information flows and concurrency in training techniques
32 are visible with the pipeline flow diagrams. Figs. 26A-D illustrate embodiments of pipeline flows for
33 layers of a neural network flow from left to right, e.g., activations enter from the left and forward pass
34 propagation of layer computations flows to the right. A gradient computation is performed in the

1 rightmost layer to begin the backward pass propagation of layer computations including weight
2 updates from right to left. Time advances from top to bottom.

3

4 [0653] Fig. 26A illustrates an embodiment of a pipeline flow for SGD. Weight updates of
5 layers of a neural network are completed after completion of a corresponding full forward pass and a
6 corresponding full backward pass through all the layers of the neural network. A next forward pass
7 begins only after completion of weight updates corresponding with an immediately preceding forward
8 pass. As illustrated, First Forward Pass 2611 is performed (from the first layer to the last layer,
9 illustrated left to right in the figure). Then First Backward Pass 2621 is performed (from the last layer
10 to the first layer, illustrated right to left in the figure). During First Backward Pass 2621, weights are
11 updated, from the last layer to the first layer. The last weight update (of the first layer) is completed as
12 First Backward Pass 2621 completes. Then Second Forward Pass 2612 is performed (using the
13 weights updated during First Backward Pass 2621), followed by Second Backward Pass 2622, during
14 which weight updates are performed.

15

16 [0654] Fig. 26B illustrates an embodiment of a pipeline flow for MBGD. A plurality of
17 activations are processed with identical weights. Coordinated quiet times are used to synchronize
18 weight updates. In some embodiments and/or usage scenarios, MBGD processing is characterized by
19 Mini-Batch Size (N) 2631, Overhead 2632, and Update Interval (U) 2633.

20

21 [0655] Unlike gradient-descent techniques (e.g., SGD and MBGD) that use a full forward
22 pass and a full backward pass through a network to compute a gradient estimate, and thus result in a
23 sequential dependency, CPGD uses a differential construction to replace the sequential dependency
24 with a continuous model that has sustained gradient generation. In some embodiments and/or usage
25 scenarios, CPGD enables layer parallelism by enabling each layer of a neural network to be trained
26 (e.g., to ‘learn’) concurrently with others of the layers without explicit synchronization. Thus,
27 parallelization along the depth of a neural network enables applying more computing resources to
28 training. In various embodiments and/or usage scenarios, CPGD provides comparable accuracy and
29 improved convergence rate expressed in epochs of training compared to other techniques.

30

31 [0656] Fig. 26C illustrates an embodiment of a pipeline flow for CPGD. CPGD processing
32 maintains a model in flux. Hidden representations and deltas enter every layer at every time step, and
33 weights update at every time step. The CPGD processing is a coordinated synchronous operation. In
34 some embodiments and/or usage scenarios, CPGD processing is characterized by Forward Pass 2651

1 and a corresponding Backward Pass 2661, respectively representing one of a number of forward
2 passes and one of a number of corresponding backward passes. In operation, respective forward
3 passes of a plurality of forward passes operate in parallel with each other, respective backward passes
4 of a plurality of backward passes operate in parallel with each other, and the pluralities of forward
5 passes and the pluralities of backward passes operate in parallel with each other. Weight updates
6 (made during backward passes) are used by forward passes and backward passes as soon as the weight
7 updates are available.

8
9 [0657] As a specific example, Forward Pass 2665 begins, and later Forward Pass 2666
10 begins. At least a portion of Forward Pass 2665 operates in parallel with at least a portion of Forward
11 Pass 2666. At least a portion of a corresponding backward pass for Forward Pass 2665 operates in
12 parallel with at least a portion of Forward Pass 2666. Further, the corresponding backward pass
13 completes at least some weight updates that are used by Forward Pass 2666, as shown by example
14 Weight Update Use 2667.

15
16 [0658] Fig. 26D illustrates an embodiment of a pipeline flow for CPGD with RCP. CPGD
17 with RCP omits saving selected activations, instead recomputing the selected activations. In some
18 embodiments and/or usage scenarios, the recomputing is performed with updated weights. Thus,
19 reverse checkpoint enables reduced memory (illustrated as reduced area covered by vertical lines
20 passing saved hidden representations forward in time) and reduces time disparity between calculated
21 hidden representations and corresponding deltas.

22
23 [0659] As a specific example, CPGD with RCP processing is characterized by Forward Pass
24 2671 and a corresponding Backward Pass 2681. A first activation is computed during the Forward
25 Pass and stored in a layer for use in the corresponding Backward Pass, as illustrated by Activation
26 Storage 2685. Activation Storage 2685 is occupied during portions of Forward Pass and Backward
27 Pass and unavailable for other uses. A specific example of memory reduction is illustrated by
28 Recomputed Activation Storage 2686. A second activation is computed during the Forward Pass, but
29 is discarded and does not require any storage. During the Backward Pass the second activation is
30 recomputed and stored in a layer for use in the Backward Pass as illustrated by Recomputed
31 Activation Storage 2686. Recomputed Activation Storage 2686 is unoccupied throughout the entire
32 Forward Pass and available for other uses (e.g., other forward passes, other backward passes), thereby
33 reducing the memory required.

34

1 [0660] Considering parallelization more generally, in some embodiments and/or usage
2 scenarios, parallelizing a computation (e.g., neural network training) spreads the computation over
3 separate computation units operating simultaneously. In a model-parallel regime, separate units
4 simultaneously evaluate a same neural network using distinct model parameters. In a data-parallel
5 regime, separate workers simultaneously evaluate distinct network inputs using the same formal
6 model parameters. Some scaling techniques use fine-grained data parallelism across layers and among
7 units in a cluster.

8
9 [0661] MBDG, in some embodiments and/or usage scenarios, improves accuracy of a
10 gradient estimate as a function of a mini-batch size, n . However, computation to perform MBDG for
11 mini-batch size n is approximately equal to computation to perform SGD for n steps. In some
12 situations, SGD for n steps is more efficient than MBDG for a mini-batch size n by approximately the
13 square root of n . Thus, higher parallelism (e.g., as in MBDG) and higher efficiency (e.g., as in SGD)
14 are sometimes mutually exclusive.

15
16 [0662] In some embodiments and/or usage scenarios, a deep neural network is a high-
17 dimensional parameterized function, sometimes expressed as a directed acyclic graph. Back
18 propagation techniques are sometimes expressed by a cyclic graph. The cycle in the graph is a
19 feedback iteration. Gradients produced by a first full network evaluation change weights used in a
20 next iteration, because the iteration is a discrete approximation of a continuous differential system.
21 The discrete approximation comprises an unbiased continuous-noise process with time-varying
22 statistics. The noise process provides regularization to enable the continuous system to model
23 phenomena observed in discrete-time learning systems. In the discrete case, regularization is provided
24 by a sampling procedure (e.g., SGD), by learning rate, and/or by other explicit mechanisms. A time-
25 dependent noise process enables using a learning-rate schedule that erases local high-frequency
26 contours in parameter space. As a correct region is approached, regularization is reduced, leading, in
27 some circumstances, to a better final solution.

28
29 [0663] CPGD, in a conceptual framework of an arbitrary feed-forward neural network,
30 expresses all nodes as functions of time and applies functional composition to formulate
31 representations in terms of internal state and stimuli the internal state is subjected to. A factorization
32 results with individual layers as systems with independent local dynamics. Two dimensions are depth
33 of the network and time evolution of parameters. In some embodiments and/or usage scenarios
34 implementing acceleration by mapping network layers to computational units separated in space, there

1 is latency communicating between the network layers. Thus there is a time delay communicating
2 between the layers. Some implementations of CPGD are synchronous implementations that account
3 for the time delays.

4
5 [0664] During CPGD processing, an activation vector and associated hidden representations
6 are combined with model parameters at different time steps during the forward pass of the activation
7 vector. The difference between model parameters at different time steps versus a same time step is not
8 detectable by the activation vector going forward. Conceptually it is as if a fixed set of parameters
9 from successive time steps were used to form an aggregate parameter state that is then used for
10 learning.

11
12 [0665] There is a choice during the backward pass (e.g., delta propagation) to use either
13 immediate parameters (e.g., weights) after updating or to retrieve historical parameters anchored to
14 when the corresponding forward pass was performed. Deltas computed from the immediate
15 parameters use updated information corresponding to a current parameter slope. Some embodiments
16 and/or usage scenarios use immediate parameters. Some embodiments and/or usage scenarios use
17 historical parameters.

18
19 [0666] Some implementations of CPGD use memory on an order similar to SGD. Reverse
20 checkpoint (as described elsewhere herein) is usable with CPGD, such as to reduce memory usage.
21 Some embodiments and/or usage scenarios of reverse checkpoint use immediate parameters (e.g.,
22 weights) to recompute activations. Some embodiments and/or usage scenarios of reverse checkpoint
23 use historical parameters to recompute activations. In some embodiments and/or usage scenarios
24 using immediate parameters to recompute activations, a time disparity between parameters used for
25 computing forward propagating activations and backward-propagating deltas is reduced in the
26 aligning wavefronts.

27
28 [0667] Continuous propagation techniques are usable in conjunction with mini-batch style
29 processing (e.g., MBGD). In some embodiments and/or usage scenarios, a subsequent batch is started
30 before an immediately preceding batch is completed, conceptually similar to asynchronous SGD.
31 Parameter inconsistency within the pipeline is limited to no more than one batch boundary.

32
33 [0668] In some embodiments and/or usage scenarios, enabling data to stream through a
34 neural network and to perform computations without a global synchronization boundary, enables

1 extracting learning information not otherwise extracted. In some embodiments and/or usage
2 scenarios, a lower learning rate dominates using larger batch sizes. In some embodiments and/or
3 usage scenarios, hidden activity and/or delta arcs are conceptually interpreted as individual vectors or
4 alternatively batch matrices. The batch matrices interpretation enables implementing techniques as
5 described herein directly on GPUs, CPUs, DSPs, FPGAs, and/or ASICs.

6
7 **[0669]** Figs. 27A-27E illustrate various aspects of forward pass and backward pass
8 embodiments in accordance with SGD, MBGD, CPGD, and RCP processing. In the figures, two
9 layers of neurons are illustrated, representing respective layers of, e.g., a portion of a deep neural
10 network. In various embodiments and/or usage scenarios, the deep neural network comprises
11 thousands or more layers and thousands or more neurons per layer. In various embodiments and/or
12 usages scenarios, the first layer is an input layer receiving activations for training from an agent
13 external to the deep neural network. In various embodiments and/or usage scenarios, the second layer
14 is an output layer where the forward pass completes, and the backward pass begins. In various
15 embodiments and/or usage scenarios, the first layer and the second layer are internal layers.

16
17 **[0670]** Fig. 27A and Fig. 27B respectively illustrate forward pass and backward pass
18 embodiments in accordance with SGD, MBGD, and CPGD, without RCP. The two layers are
19 illustrated as Previous Layer 2701 and Subsequent Layer 2702. Previous Layer 2701 comprises
20 Compute 2710 and Storage 2715. Subsequent Layer 2702 comprises Compute 2720 and Storage
21 2725. Compute 2710 and Compute 2720 are examples of compute resources and Storage 2715 and
22 Storage 2725 are examples of storage resources.

23
24 **[0671]** Figs. 27C-27E illustrate forward pass and backward pass embodiments in accordance
25 with SGD, MBGD, and CPGD, with RCP. The two layers are illustrated as Previous Layer 2703 and
26 Subsequent Layer 2704. Previous Layer 2703 comprises Compute 2730 and Storage 2735.
27 Subsequent Layer 2704 comprises Compute 2740 and Storage 2745. Compute 2730 and Compute
28 2740 are examples of compute resources and Storage 2735 and Storage 2745 are examples of storage
29 resources.

30
31 **[0672]** Like-numbered elements in Figs. 27A-27E have identical structure and operation,
32 although the compute resources produce different results dependent on differing inputs, and the
33 storage resources store and subsequently provide different values dependent on differing values
34 stored. Other embodiments are envisioned with differing compute resources and/or differing storage

1 resources usable for forward pass and backward pass computation and storage. E.g., a backward pass
2 uses a transpose weight storage not used by a forward pass. Other embodiments are envisioned with
3 differing compute and/or storage resources usable for differing forward pass and backward pass
4 implementations. E.g., an RCP-based embodiment uses an additional compute resource (not
5 illustrated) than used for forward pass or backward pass processing without RCP.

6
7 **[0673]** Regarding Fig. 27A, Compute 2710 is enabled to perform computations, such as
8 forward pass computations F 2711. Storage 2715 is enabled to store activations, such as in A 2716.
9 Storage 2715 is further enabled to store weights, such as in W 2717. Compute 2720, F 2721, Storage
10 2725, A 2726, and W 2727, are, in various embodiments and/or usage scenarios, substantially similar
11 or identical in structure and/or operation respectively to Compute 2710, F 2711, Storage 2715, A
12 2716, and W 2717.

13
14 **[0674]** In forward pass operation for SGD or MBGD, activation $A_{1,t}$ 2781 is received by
15 Previous Layer 2701 and stored in A 2716 (for later use during the backward pass). $A_{1,t}$ 2781 and a
16 weight $W_{1,t}$, previously stored in W 2717, are then processed in accordance with F 2711 to produce
17 activation $A_{2,t}$ 2782. $A_{2,t}$ 2782 is then passed to Subsequent Layer 2702. Similarly to the Previous
18 Layer, $A_{2,t}$ 2782 is received by Subsequent Layer 2702 and stored in A 2726 (for later use during the
19 backward pass). $A_{2,t}$ 2782 and a weight $W_{2,t}$, previously stored in W 2727 are then processed in
20 accordance with F 2721 to produce activation $A_{3,t}$ 2783. $A_{3,t}$ 2783 is then provided to a next
21 subsequent layer (if present) for processing, and so forth, until the forward pass is complete and the
22 backward pass commences. If Subsequent Layer 2702 is the output layer, then the forward pass is
23 completed and the backward pass corresponding to the forward pass is initiated.

24
25 **[0675]** Regarding Fig. 27B, for clarity, elements of Compute 2710 and Compute 2720
26 dedicated to forward pass processing (F 2711 and F 2721) are omitted. With respect to structure and
27 operation illustrated and described with respect to Fig. 27A, Fig. 27B illustrates that Compute 2710 is
28 further enabled to perform additional computations, such as backward pass computations B 2712, and
29 Compute 2720 is further enabled to perform additional computations, such as backward pass
30 computations B 2722. Storage 2715 is further enabled to store a computed weight, such as in W 2718,
31 and Storage 2725 is further enabled to store a computed weight, such as in W 2728. B 2722 and W
32 2728 are, in various embodiments and/or usage scenarios, substantially similar or identical in structure
33 and/or operation respectively to B 2712 and W 2718.

34

1 [0676] In backward pass operation for SGD or MBGD, delta $\Delta_{3,t}$ 2793 is received from the
2 next subsequent layer (if present) during backward pass processing. If Subsequent Layer 2702 is the
3 output layer, then Subsequent Layer 2702 computes delta $\Delta_{3,t}$ according to the delta rule, e.g., as a
4 function of the difference between the output of the Subsequent Layer (e.g., the estimated output) and
5 the training output (e.g., desired output). $\Delta_{3,t}$ 2793, the weight $W_{2,t}$ previously stored in W 2727, and
6 the activation $A_{2,t}$ previously stored in A 2726, are then processed in accordance with B 2722 (e.g., in
7 accordance with the delta rule) to produce delta $\Delta_{2,t}$ 2792 and a new weight $W_{2,t+1}$ that is then stored in
8 W 2728 for use in a next forward pass. $\Delta_{2,t}$ 2792 is then passed to Previous Layer 2701. Similarly to
9 the Subsequent Layer, delta $\Delta_{2,t}$ 2792, the weight $W_{1,t}$ previously stored in W 2717, and the activation
10 $A_{1,t}$ previously stored in A 2716, are then processed in accordance with B 2712 to produce delta $\Delta_{1,t}$
11 2791 and a new weight $W_{1,t+1}$ that is then stored in W 2718 for use in the next forward pass. $\Delta_{1,t}$ 2791
12 is then passed to a next previous layer (if present) for processing, and so forth, until the backward pass
13 is complete and a next forward pass commences. If Previous Layer 2701 is the input layer, then the
14 backward pass is complete, and the next forward pass commences.

15
16 [0677] In SGD and MBGD (and unlike CPGD), the next forward pass is delayed until the
17 previous backward pass completes, e.g., W 2717 and W 2727 are respectively updated with W 2718
18 and W 2728 after W 2717 and W 2727 have been used for a same forward pass and a same
19 corresponding backward pass. Therefore, the next forward pass is performed using weights that are
20 from the same backward pass.

21
22 [0678] Fig. 27A, in addition to illustrating SGD and MBGD forward pass processing, also
23 illustrates CPGD forward pass processing. However, operation for CPGD is different compared to
24 SGD and MBGD, in that weight updates and the next forward pass are performed as soon as possible,
25 rather than being delayed until completion of the previous backward pass. E.g., W 2717 and W 2727
26 are respectively updated with W 2718 and W 2728 as soon as possible. Therefore, the next forward
27 pass has selective access to weights from prior iterations, and thus selectively produces activations
28 differing from those produced under the same conditions by SGD and MBGD.

29
30 [0679] More specifically, in Previous Layer 2701, $A_{1,t}$ 2781 is received and stored in A 2716,
31 identically to SGD and MBGD. $A_{1,t}$ 2781 and a weight $W_{1,t,k-j}$ previously stored in W 2717 are then
32 processed in accordance with F 2711 to produce activation $A_{2,t}$ 2782. The weight $W_{1,t,k-j}$ was
33 produced and stored by a backward pass corresponding to a forward pass preceding the instant
34 forward pass by $k-j$ forward passes. $A_{2,t}$ 2782 is then passed to Subsequent Layer 2702, and similarly

1 to the Previous Layer, $A_{2,t}$ 2782 is received and stored in A 2726, identically to SGD and MBGD. $A_{2,t}$
2 2782 and a weight $W_{2,t-k}$ previously stored in W 2727 are then processed in accordance with F 2721 to
3 produce activation $A_{3,t}$ 2783. The weight $W_{2,t-k}$ was produced and stored by a backward pass
4 corresponding to a forward pass preceding the instant forward pass by k forward passes. Note that the
5 Previous Layer and the Subsequent Layer, for processing of a same forward pass, use weights from
6 different backward passes. As in SGD and MBGD, $A_{3,t}$ 2783 is then provided to a next subsequent
7 layer (if present) for processing, and so forth, until the forward pass is complete and the backward
8 pass commences. If Subsequent Layer 2702 is the output layer, then the forward pass is completed
9 and the backward pass corresponding to the forward pass is initiated. In some embodiments and/or
10 usage scenarios, the value of j is 0 and (k-j) and (k) are equal. In various embodiments and/or usage
11 scenarios, the Previous Layer and the Subsequent Layer simultaneously process one of: different
12 forward passes, different backward passes, and a forward pass and a different backward pass.

13

14 [0680] Fig. 27B, in addition to illustrating SGD and MBGD backward pass processing, also
15 illustrates CPGD backward pass processing. Processing of the backward pass in CPGD is identical to
16 that of SGD and MBGD. However, selected results (e.g., selected weights) are used earlier than in
17 SGD and MBGD. For example, $W_{1,t-k-j}$, as produced by backward pass t-k-j, and $W_{1,t-k}$, as produced
18 by backward pass t-k are used earlier than in SGD and MBGD, e.g., forward pass t.

19

20 [0681] Fig. 27C illustrates an embodiment of forward pass processing of any of SGD,
21 MBGD, and CPGD, in combination with RCP. Compute 2730 and Storage 2735, are, in various
22 embodiments and/or usage scenarios, substantially similar or identical in structure and/or operation
23 respectively to Compute 2710 and Storage 2715. Compute 2740 and Storage 2745, are, in various
24 embodiments and/or usage scenarios, substantially similar or identical in structure and/or operation
25 respectively to Compute 2720 and Storage 2725, other than omission of storage for activations A 2726
26 of Storage 2725 having no counterpart in Storage 2745.

27

28 [0682] In forward pass operation, with respect to Previous Layer 2703, activation $A_{1,t}$ 2781 is
29 received and processed in accordance with forward pass processing in Compute 2730, and stored in
30 Storage 2735 as described with respect to Fig. 27A. However, with respect to Subsequent Layer 2704,
31 activation $A_{2,t}$ 2782 is received, and processed in accordance with forward pass processing in
32 Compute 2740, but is not stored (instead it is recomputed in accordance with RCP during backward
33 pass processing).

34

1 [0683] Fig. 27D and Fig. 27E respectively illustrate first and second portions of an
2 embodiment of backward pass processing of any of SGD, MBGD, and CPGD, in combination with
3 RCP. For clarity, elements of Compute 2730 and Compute 2740 dedicated to forward pass processing
4 (F 2721) are omitted. With respect to structure and operation illustrated and described with respect to
5 Fig. 27C, Fig. 27D and Fig. 27E illustrate that Compute 2730 is further enabled to perform additional
6 computations, such as backward pass computations B 2712, and Compute 2740 is further enabled to
7 perform additional computations, such as backward pass computations B 2722. Storage 2735 is
8 further enabled to store a computed weight, such as in W 2718, and Storage 2745 is further enabled to
9 store a computed weight, such as in W 2728, as well as a recomputed activation, such as in A 2729.

10

11 [0684] In the first portion of the backward pass operation, activations not stored in the
12 corresponding forward pass are recomputed. In SGD and MBGD scenarios, the recomputed
13 activation is formulated in Previous Layer 2703 by processing the activation stored from the forward
14 pass in A 2716 and weight stored in W 2717 in accordance with F 2711 to produce activation $A'_{2,t}$
15 2784, that is then stored in A 2729 of Subsequent Layer 2704. Since SGD and MBGD delay weight
16 updates and commencement of a next forward pass until the forward pass and corresponding
17 backward pass are complete, $A'_{2,t}$ 2784 is identical to the value discarded during the forward pass, $A_{2,t}$
18 2782.

19

20 [0685] In a CPGD scenario, the recomputed activation is formulated according to the same
21 topology as the SGD and MBGD scenarios. However, CPGD performs updates without delays and
22 enables commencement of a next forward pass without regard to completion of previous backward
23 passes. Thus, a weight value stored at the time of the backward pass, e.g., in W 2717, according to
24 embodiment and/or usage scenarios, selectively differs from the weight value stored during the
25 corresponding forward pass. As a specific example, in accordance with Fig. 27C, W 2717 stored $W_{1,t}$
26 $_{k,j}$ during the forward pass. However, during the backward pass, additional weight updates have
27 occurred, e.g., corresponding to m iterations, and now W 2717 stores $W_{1,t+k-j+m}$. Therefore, $A'_{2,t}$ 2784
28 selectively differs from the value discarded during the forward pass, $A_{2,t}$ 2782.

29

30 [0686] In the second portion of backward pass operation, computation proceeds using the
31 recomputed activation. In SGD and MBGD scenarios, since the recomputed activation is identical to
32 the discarded activation (e.g., conceptually the value stored in A 2729 is identical to the value stored
33 in A 2726), the backward processing produces results that are identical to the results described with
34 respect to Fig. 27B. E.g., deltas $\Delta'_{3,t}$ 2796, $\Delta'_{2,t}$ 2795, and $\Delta'_{1,t}$ 2794 are identical, respectively, to $\Delta_{3,t}$

1 **2793**, $\Delta_{2,t}$ **2792**, and $\Delta_{1,t}$ **2791**. In the CPGD scenario, since the recomputed activation selectively
2 differs from the discarded activation, the backward processing produces results that are selectively
3 different from the results described with respect to Fig. 27B. E.g., deltas $\Delta'_{3,t}$ **2796**, $\Delta'_{2,t}$ **2795**, and
4 $\Delta'_{1,t}$ **2794** are selectively different, respectively, to $\Delta_{3,t}$ **2793**, $\Delta_{2,t}$ **2792**, and $\Delta_{1,t}$ **2791**.

5

6 **[0687]** In some embodiments and/or usage scenarios, W **2717** is distinct from W **2718** (as
7 illustrated), and in some embodiments and/or usage scenarios, W **2718** and W **2717** are a same portion
8 of storage (not illustrated), such that saving a new value in W **2718** overwrites a previously saved
9 value in W **2717**. Similarly, W **2727** is variously distinct from or the same as W **2728**. In various
10 embodiments and/or usage scenarios, A **2729** is variously implemented to use fewer memory locations
11 and/or use a same number of memory locations for a shorter time than A **2726**.

12

13 **[0688]** In various embodiments and/or usages scenarios, activations and/or weights are
14 implemented and/or represented by any one or more scalar, vector, matrix, and higher-dimensional
15 data structures. E.g., any one or more of A **2716**, A **2726**, A **2729**, W **2717**, W **2727**, W **2718**, and W
16 **2728** are enabled to store any one or more of one or more scalars, one or more vectors, one or more
17 matrices, and one or more higher-dimensional arrays.

18

19 **[0689]** In various embodiments and/or usage scenarios, one or more elements of Previous
20 Layer **2701** and Subsequent Layer **2702** are implemented by respective PEs, e.g., a portion of PE **499**
21 or similar elements of Fig. 4. E.g., PE **497** implements Previous Layer **2701** and PE **498** implements
22 Subsequent Layer **2702**. Activation $A_{2,t}$ **2782** and delta $\Delta_{2,t}$ **2792** are communicated via East coupling
23 **431**. In some embodiments and/or usage scenarios, one or more elements of Previous Layer **2701** and
24 Subsequent Layer **2702** are implemented by one or more of CPUs, GPUs, DSPs, and FPGAs.

25

26 **[0690]** In various embodiments and/or usage scenarios, all or any portions of elements of F
27 **2711**, F **2721**, B **2712**, and B **2722** conceptually correspond to all or any portions of executions of
28 instructions of Task SW on PEs **260** of Fig. 2.

29

30

31 EXAMPLE WORKLOAD MAPPING

32

33 **[0691]** Conceptually, Deep Learning Accelerator **400** (Fig. 4) is a programmable compute
34 fabric (see, e.g., Figs. 5-8 and section "Processing Element: Compute Element and Router"). For

1 example, the compute element of each PE 499 element is enabled to execute sequences of instructions
2 of tasks (such as conceptually corresponding to all or any portions of executions of instructions of
3 Task SW on PEs 260 of Fig. 2), and the router element of router element of each PE 499 is
4 configurable to route wavelets between the PEs. The programmable compute fabric enables mapping
5 of workloads onto the compute fabric in various manners. Described following is an example high-
6 level mapping of a workload to the compute fabric to illustrate various techniques and mechanisms
7 implemented by the compute fabric.

8
9 [0692] The workload is deep neural network training, implemented via SGD. The deep
10 neural network comprises a plurality of layers of neurons. The workload has three mega-phases: a
11 forward pass, a delta pass, and a chain pass. The forward pass propagates activations in a forward
12 direction. The delta pass propagates deltas in a backward direction. The chain pass calculates
13 gradients based on the deltas as the deltas are generated in the delta pass. The three mega-phases have
14 approximately a same amount of compute.

15
16 [0693] Fig. 4 illustrates an example mapping of the mega-phases to the PEs. Each layer is
17 implemented by blocks of PEs allocated from the compute fabric (aka ‘placed’) back-to-back (e.g., in
18 a horizontal dimension). Data movement propagates to the end of the fabric during the forward pass
19 (Forward 401), and then circles back in the reverse direction during the delta pass (Delta 402) and
20 chain pass (Chain 403). The placement is directed to reduce data movement since the forward pass
21 saves activations to be used by the delta pass and the chain pass. In the example, all the PEs are time
22 shared three ways between the three mega-phases, with each mega-phase using approximately a same
23 amount of compute. In some circumstances, an entire chain of PEs performing the passes operates as
24 a pipeline such that each layer is a pipe stage (taking roughly a same amount of time to complete) and
25 each activation of a mini-batch is fills the pipeline.

26
27 [0694] In some embodiments and/or usage scenarios, within a set of the PEs mapped to a
28 single one of the layers, the weights of the single layer are distributed across the PEs such that a single
29 neuron is mapped to multiple PEs. Splitting a single neuron across multiple PEs, in some
30 circumstances, provides a load balancing benefit and provides a communication partitioning benefit
31 (see, e.g., Figs. 17-20 and section “Neuron Smearing”).

32
33 [0695] Conceptually, processing proceeds as follows (see Forward 401 of Fig. 4).
34 Activations are broadcasted into the layer along the horizontal axis. Activations are received by the

1 PEs and trigger a lookup of the associated weights that are stored local to the PEs (corresponding to
2 the neurons mapped to the PEs). Only non-zero activations are broadcasted, so no compute is wasted
3 for zero activations (an example of activation sparsity harvesting). Each PE performs a local multiply
4 and accumulate of the incoming activation with all the neuron weights producing local partial sums.
5 Since the weights of each neuron are distributed to multiple PEs, partial sums are then accumulated
6 across the PEs in the vertical direction, in accordance with the neuron weight distribution. After the
7 partial sums are accumulated producing a final sum, the activation function is performed and all new
8 non-zero activations are broadcast to the next layer.

9

10 [0696] The delta pass (see Delta 402 of Fig. 4) and the chain pass (see Chain 403 of Fig. 4)
11 follow a data flow similar to that of the forward pass. In some embodiments and/or usage scenarios,
12 the delta pass and the chain pass are placed offset by one layer so the activations are stored in the same
13 layers as the weights used in the backward direction. Activations are stored by the receiving layer
14 such that in the delta pass and the chain pass, the activations are used directly without additional
15 communication. In addition to storing activations, a weight transpose is performed to implement the
16 delta pass. The weight transpose, in some embodiments and/or usage scenarios, is implemented by
17 replicating the weights, using additional memory capacity and additional communication when
18 updating the weights. In some embodiments and/or usage scenarios, the weight transpose is
19 implemented by transposing the delta broadcast in the vertical dimension.

20

21 [0697] Fig. 28A illustrates a generic operation of a matrix (m) multiplied by a vector (v).
22 Fig. 28B illustrates, in the style of Fig. 28A, various representations of memory structures used in the
23 three mega-phases in some embodiments (e.g., a fully connected neural network). In various
24 embodiments, the weight (w) and the gradient accumulation (g) data structures are two-dimensional
25 matrices. In some embodiments, the forward partial sum ($fpsum$) and delta partial sum ($\delta psum$) and
26 forward pass activations (a) are one-dimensional vectors. The two-dimensional matrices are stored in
27 memory (e.g., Memory 854 of Fig. 8) since in some embodiments and/or usage scenarios the two-
28 dimensional matrices are relatively large. In some embodiments, the one-dimensional vectors are
29 stored in higher-throughput storage (e.g., D-Store 848 of Fig. 8) to enable, usage scenarios, full
30 datapath performance for the multiply-accumulate vector operation in each of the three passes.

31

32 [0698] Fig. 29 illustrates an embodiment of tasks (see, e.g., Figs. 9-12 and section "Tasks")
33 as used in a forward pass state machine. In some embodiments and/or usage scenarios, each of the
34 PEs implements an instantiation of the state machine. In some embodiments and/or usage scenarios,

1 various portions of the state machine are implemented by respective PEs (see, e.g., Figs. 17-20 and
2 section “Neuron Smearing”). There are four tasks in the state machine: f_rxact:acc 2901,
3 f_rxact:close 2902, f_psum:prop 2903, and f_txact:tx 2904. Conceptually, activations arrive from a
4 PE to the “left” of the instant PE (corresponding to a previous layer). Incoming (non-closeout)
5 activations on the activation broadcast wire (Activations 2911) trigger f_rxact:acc 2901. The instant
6 PE executes instructions of the task, looking up (e.g., from memory local to the instant PE) the
7 weights associated with the activation and performing the local weight multiply and accumulate into
8 partial sums. Control flow dependencies exist between f_rxact:acc 2901 and f_psum:prop 2903 (Flow
9 2913). Example data structures the task references are wrow, fpsum, and fact.

10

11 [0699] An incoming activation closeout on the activation broadcast wire (Closeouts 2912)
12 triggers f_rxact:close 2902. The closeout signals the end of all activations for the current wavefront.
13 The instant PE executes instructions of the task, starting the partial sum accumulation ring with the
14 partial sums in a start list of the instant PE (Start Psums 2916). Example data structures the task
15 references are fpsum_acc_mem, and fpsum_acc_fab.

16

17 [0700] An incoming partial sum (Prop Psums 2930) triggers f_psum:prop 2903. The instant
18 PE executes instructions of the task, adding the incoming partial sum to the local partial sum of the
19 instant PE, and then forwarding the result to the next hop on the ring (Prop Psums 2931). If the
20 instant PE is the end of the ring, then the final sum is generated. In some embodiments and/or usage
21 scenarios, additional processing is performed to prevent deadlock. Example data structures the task
22 references are fpsum_acc_mem, fpsum_acc_fab, and f_txact_wake.

23

24 [0701] When there are queued activations to transmit, f_txact:tx 2904 is self-triggered (Wake
25 2914). The instant PE executes instructions of the task, de-queuing an activation and transmitting the
26 activation on the broadcast wire to the next layer (Activations 2921). When more items remain in the
27 queue, the instant PE reschedules the task (Reschedule 2915). When the queue is empty, the instant
28 PE sends a closeout wavelet to close the wavefront (Closeouts 2922).

29

30 [0702] The activations (incoming and outgoing) and the partial sums (incoming and
31 outgoing), as well as the closeout wavelets are communicated as wavelets (see, e.g., Figs. 13A-15B
32 and section “Wavelets”). In some embodiments and/or usage scenarios, one or more of the wavelets
33 correspond to one or more elements of fabric vectors as described by one or more DSDs and/or
34 XDSDs.

1

2 [0703] Data structures for the various state machines are referenced via a plurality of DSDs
3 stored in respective DSRs (see, e.g., Figs. 21A-24 and section “Vectors and Data Structure
4 Descriptors”), as described by the following table.

DSR	Data Structure Name	Description
DS1	Wrow	Weight matrix, rows
DS2	Wcol	Weight matrix, cols (points to same data as DS2)
DS3	Fpsum	Forward partial sum vector – full vector of all psums Length: number of neurons Stride: 1
DS4	fpsum_acc_mem	Forward partial sum vector – subset for psum accumulate Same data as psum but organized as 2d array Length: number of neurons in subset Stride: 1
DS5	fpsum_acc_fab	Forward partial sum vector – subset for psum accumulate Fabric type: col:ep=f_psum:prop Length: number of neurons in subset
DS6	Fact	Forward activation storage vector Length: 1 Stride: 1
DS7	fact_fab	Forward activation fabric transmit Fabric type: col:ep=f_tfact:acc Length: 1
DS8	f_tfact_wake	Self reschedule wake up wavelet Fabric type: col:ep=f_tfact:tx
DS9	fact_close_fab	Forward activation close out fabric transmit Fabric type: col:ep=f_tfact:close Length: 1

5

6 [0704] The foregoing example workload mapping is with respect to SGD. However, the
7 techniques are readily applicable to MBGD and CPGD, with and without RCP.

8

9

1 OTHER EMBODIMENT DETAILS

2
3 [0705] Embodiments and usage scenarios described with respect to Figs. 1-29 are
4 conceptually with respect to a PE comprising a CE that is programmable, e.g., that processes data
5 according to instructions. Other embodiments are contemplated with one or more of the CEs being
6 partially or entirely hardwired, e.g., that process data according to one or more fixed-circuit processing
7 elements operable without instructions. As a specific example, a particular CE comprises a hardware
8 logic unit circuit that implements all or a portion of an LSTM unit. The particular CE is comprised
9 with a router in a particular PE that is operable in a fabric with other PEs. Some of the other PEs are
10 similar to or identical to the particular PE and some of the other PEs are similar to or identical to PE
11 499 of Fig. 4.

12
13
14 EXAMPLE IMPLEMENTATION TECHNIQUES

15
16 [0706] In some embodiments, various combinations of all or any portions of operations
17 performed for and/or structure associated with any of accelerated deep learning; SGD, MBGD, CPGD
18 with and without RCP for accelerated deep learning; data structure descriptors and fabric vectors for
19 accelerated deep learning; neuron smearing for accelerated deep learning; task synchronization for
20 accelerated deep learning; dataflow triggered tasks for accelerated deep learning; a control wavelet for
21 accelerated deep learning; and/or a wavelet representation for accelerated deep learning; as well as
22 portions of a processor, microprocessor, system-on-a-chip, application-specific-integrated-circuit,
23 hardware accelerator, or other circuitry providing all or portions of the aforementioned operations, are
24 specified by a specification compatible with processing by a computer system. The specification is in
25 accordance with various descriptions, such as hardware description languages, circuit descriptions,
26 netlist descriptions, mask descriptions, or layout descriptions. Example descriptions include: Verilog,
27 VHDL, SPICE, SPICE variants such as PSpice, IBIS, LEF, DEF, GDS-II, OASIS, or other
28 descriptions. In various embodiments, the processing includes any combination of interpretation,
29 compilation, simulation, and synthesis to produce, to verify, or to specify logic and/or circuitry
30 suitable for inclusion on one or more integrated circuits. Each integrated circuit, according to various
31 embodiments, is compatible with design and/or manufacture according to a variety of techniques. The
32 techniques include a programmable technique (such as a field or mask programmable gate array
33 integrated circuit), a semi-custom technique (such as a wholly or partially cell-based integrated
34 circuit), and a full-custom technique (such as an integrated circuit that is substantially specialized),

1 any combination thereof, or any other technique compatible with design and/or manufacture of
2 integrated circuits.

3

4 **[0707]** In some embodiments, various combinations of all or portions of operations as
5 described by a computer readable medium having a set of instructions stored therein, are performed by
6 execution and/or interpretation of one or more program instructions, by interpretation and/or
7 compiling of one or more source and/or script language statements, or by execution of binary
8 instructions produced by compiling, translating, and/or interpreting information expressed in
9 programming and/or scripting language statements. The statements are compatible with any standard
10 programming or scripting language (such as C, C++, Fortran, Pascal, Ada, Java™, VBscript, and
11 Shell). One or more of the program instructions, the language statements, or the binary instructions, are
12 optionally stored on one or more computer readable storage medium elements. In various
13 embodiments, some, all, or various portions of the program instructions are realized as one or more
14 functions, routines, sub-routines, in-line routines, procedures, macros, or portions thereof.

15

1 CONCLUSION

2
3 **[0708]** Certain choices have been made in the description merely for convenience in
4 preparing the text and drawings, and unless there is an indication to the contrary, the choices should
5 not be construed per se as conveying additional information regarding structure or operation of the
6 embodiments described. Examples of the choices include: the particular organization or assignment
7 of the designations used for the figure numbering and the particular organization or assignment of the
8 element identifiers (the callouts or numerical designators, e.g.) used to identify and reference the
9 features and elements of the embodiments.

10
11 **[0709]** Various forms of the words “include” and “comprise” are specifically intended to be
12 construed as abstractions describing logical sets of open-ended scope and are not meant to convey
13 physical containment unless described explicitly (such as followed by the word “within”).

14
15 **[0710]** Although the foregoing embodiments have been described in some detail for purposes
16 of clarity of description and understanding, the invention is not limited to the details provided. There
17 are many embodiments of the invention. The disclosed embodiments are exemplary and not
18 restrictive.

19
20 **[0711]** It will be understood that many variations in construction, arrangement, and use are
21 possible consistent with the description, and are within the scope of the claims of the issued patent.
22 For example, interconnect and function-unit bit-widths, clock speeds, and the type of technology used
23 are variable according to various embodiments in each component block. The names given to
24 interconnect and logic are merely exemplary, and should not be construed as limiting the concepts
25 described. The order and arrangement of flowchart and flow diagram process, action, and function
26 elements are variable according to various embodiments. Also, unless specifically stated to the
27 contrary, value ranges specified, maximum and minimum values used, or other particular
28 specifications (such as file types; and the number of entries or stages in registers and buffers), are
29 merely those of the described embodiments, are expected to track improvements and changes in
30 implementation technology, and should not be construed as limitations.

31
32 **[0712]** Functionally equivalent techniques known in the art are employable instead of those
33 described to implement various components, sub-systems, operations, functions, routines, sub-
34 routines, in-line routines, procedures, macros, or portions thereof. It is also understood that many

1 functional aspects of embodiments are realizable selectively in either hardware (e.g., generally
2 dedicated circuitry) or software (e.g., via some manner of programmed controller or processor), as a
3 function of embodiment dependent design constraints and technology trends of faster processing
4 (facilitating migration of functions previously in hardware into software) and higher integration
5 density (facilitating migration of functions previously in software into hardware). Specific variations
6 in various embodiments include, but are not limited to: differences in partitioning; different form
7 factors and configurations; use of different operating systems and other system software; use of
8 different interface standards, network protocols, or communication links; and other variations to be
9 expected when implementing the concepts described herein in accordance with the unique engineering
10 and business constraints of a particular application.

11

12 [0713] The embodiments have been described with detail and environmental context well
13 beyond that required for a minimal implementation of many aspects of the embodiments described.
14 Those of ordinary skill in the art will recognize that some embodiments omit disclosed components or
15 features without altering the basic cooperation among the remaining elements. It is thus understood
16 that much of the details disclosed are not required to implement various aspects of the embodiments
17 described. To the extent that the remaining elements are distinguishable from the prior art,
18 components and features that are omitted are not limiting on the concepts described herein.

19

20 [0714] All such variations in design are insubstantial changes over the teachings conveyed by
21 the described embodiments. It is also understood that the embodiments described herein have broad
22 applicability to other computing and networking applications, and are not limited to the particular
23 application or industry of the described embodiments. The invention is thus to be construed as
24 including all possible modifications and variations encompassed within the scope of the claims of the
25 issued patent.

26

CLAIMS

1. A method comprising:

training a neural network comprising a plurality of ordered, connected layers;
wherein the order identifies for each respective layer which others of the layers are prior to the respective layer and which others of the layers are subsequent to the respective layer;

wherein each layer comprises one or more neurons, each neuron comprising weights, and connected to at least one of at least one prior neuron of a prior layer and at least one subsequent neuron of a subsequent layer; and

wherein each neuron is implemented by one or more processing elements, each processing element comprising

at least one coupling to a fabric, the processing element being enabled to communicate via the fabric via a plurality of virtual channels,

a first memory enabled to store instructions corresponding to at least computations of the neuron,

a second memory enabled to store the weights, and

hardware execution resources enabled to execute instructions from the respective first memory and access data from the respective second memory.

2. The method of claim 1, wherein each of the layers is a respective internal layer of the neural network, and the neural network further comprises an input layer and an output layer.

3. The method of claim 1, wherein the training comprises:
 - determining a second activation based on a first activation and first weights;
 - determining and saving second weights based on a first delta and the first weights;
 - determining a fourth activation based on a third activation and selected weights, wherein the selected weights are dynamically selected from the first weights and the second weights; and
 - determining and saving third weights based on a second delta and the selected weights.

4. The method of claim 3, wherein the determining the second activation comprises:
 - receiving the first activation via the fabric from the at least one prior neuron;
 - computing the second activation based at least in part on the first activation and the first weights by at least executing first instructions stored in the first memory and accessing the first weights in the second memory; and
 - selectively transmitting the second activation via the fabric to the at least one subsequent neuron.

5. The method of claim 4, wherein the determining the fourth activation comprises:
 - receiving the third activation via the fabric from the at least one prior neuron;
 - computing the fourth activation based at least in part on the third activation and the selected weights by at least executing the first instructions and accessing the selected weights in the second memory; and
 - selectively transmitting the fourth activation via the fabric to the at least one subsequent neuron.

6. The method of claim 5, wherein the determining and saving the second weights comprises:

receiving the first delta that is partially based on the second activation via the fabric from the at least one subsequent neuron;

computing a first gradient based at least in part on the first delta and the second activation by at least executing second instructions stored in the first memory;

computing the second weights based at least in part on the first gradient, a learning rule, and the first weights by at least executing third instructions stored in the first memory and accessing the first weights in the second memory; and

storing the second weights in the second memory.

7. The method of claim 6, wherein the determining and saving the third weights comprises:

receiving the second delta that is partially based on the fourth activation via the fabric from the at least one subsequent neuron;

computing a second gradient based at least in part on a third delta and the fourth activation by at least executing the second instructions stored in the first memory;

computing the third weights based at least in part on the second gradient, the learning rule and the selected weights by at least executing the third instructions stored in the first memory and accessing the selected weights in the second memory; and

storing the third weights in the second memory.

8. The method of claim 7, wherein the computing the second gradient additionally comprises recomputing the fourth activation based at least in part upon the selected weights.

9. An apparatus comprising:

a plurality of processing elements;

a training workload comprising a set of machine codes selected from a predefined native instruction set of codes for performing training of a neural network comprising a plurality of ordered, connected layers;

wherein the order identifies for each respective layer which others of the layers are prior to the respective layer and which others of the layers are subsequent to the respective layer;

wherein each layer comprises one or more neurons, each neuron comprising weights, and connected to at least one of at least one prior neuron of a prior layer and at least one subsequent neuron of a subsequent layer; and

wherein each neuron is implemented by one or more of the processing elements, each processing element comprising

at least one coupling to a fabric, the processing element being enabled to communicate via the fabric via a plurality of virtual channels,

a first memory enabled to store instructions corresponding to at least computations of the neuron,

a second memory enabled to store the weights, and

a compute engine enabled to perform a predefined set of basic operations in response to receiving a corresponding basic instruction selected from the predefined native instruction set of codes, the performing via execution of instructions from the respective first memory and access of data from the respective second memory.

10. The apparatus of claim 9, wherein each of the layers is a respective internal layer of the neural network, and the neural network further comprises an input layer and an output layer.

11. The apparatus of claim 9, wherein the training workload comprises respective sets of the machine codes directed to:

determining a second activation based on a first activation and first weights;

determining and saving second weights based on a first delta and the first weights;

determining a fourth activation based on a third activation and selected weights, wherein the selected weights are dynamically selected from the first weights and the second weights; and

determining and saving third weights based on a second delta and the selected weights.

12. The apparatus of claim 11, wherein the determining the second activation comprises:

receiving the first activation via the fabric from the at least one prior neuron;

computing the second activation based at least in part on the first activation and the first weights by at least executing first instructions stored in the first memory and accessing the first weights in the second memory; and

selectively transmitting the second activation via the fabric to the at least one subsequent neuron.

13. The apparatus of claim 12, wherein the determining the fourth activation comprises:

receiving the third activation via the fabric from the at least one prior neuron;
computing the fourth activation based at least in part on the third activation and the selected weights by at least executing the first instructions and accessing the selected weights in the second memory; and
selectively transmitting the fourth activation via the fabric to the at least one subsequent neuron.

14. The apparatus of claim 13, wherein the determining and saving the second weights comprises:

receiving the first delta that is partially based on the second activation via the fabric from the at least one subsequent neuron;
computing a first gradient based at least in part on the first delta and the second activation by at least executing second instructions stored in the first memory;
computing the second weights based at least in part on the first gradient, a learning rule, and the first weights by at least executing third instructions stored in the first memory and accessing the first weights in the second memory; and
storing the second weights in the second memory.

15. The apparatus of claim 14, wherein the determining and saving the third weights comprises:

receiving the second delta that is partially based on the fourth activation via the fabric from the at least one subsequent neuron;
computing a second gradient based at least in part on a third delta and the fourth activation by at least executing the second instructions stored in the first memory;

computing the third weights based at least in part on the second gradient, the learning rule and the selected weights by at least executing the third instructions stored in the first memory and accessing the selected weights in the second memory; and

storing the third weights in the second memory.

16. The apparatus of claim 15, wherein the computing the second gradient additionally comprises recomputing the fourth activation based at least in part upon the selected weights.

17. The apparatus of claim 9, wherein the apparatus is implemented via a substantially whole wafer comprising the processing elements.

18. A system comprising:

means for training a neural network comprising a plurality of ordered, connected layers;

wherein the order identifies for each respective layer which others of the layers are prior to the respective layer and which others of the layers are subsequent to the respective layer;

wherein each layer comprises one or more neurons, each neuron comprising weights, and connected to at least one of at least one prior neuron of a prior layer and at least one subsequent neuron of a subsequent layer; and

wherein each neuron is implemented by one or more processing elements, each processing element comprising

at least one coupling to a fabric, the processing element being enabled to communicate via the fabric via a plurality of virtual channels,

a first memory enabled to store instructions corresponding to at least computations of the neuron,
a second memory enabled to store the weights, and
hardware execution resources enabled to execute instructions from the respective first memory and access data from the respective second memory.

19. The system of claim 18, wherein each of the layers is a respective internal layer of the neural network, and the neural network further comprises an input layer and an output layer.

20. The system of claim 18, wherein the means for training comprises:

means for determining a second activation based on a first activation and first weights;

means for determining and saving second weights based on a first delta and the first weights;

means for determining a fourth activation based on a third activation and selected weights, wherein the selected weights are dynamically selected from the first weights and the second weights; and

means for determining and saving third weights based on a second delta and the selected weights.

21. The system of claim 20, wherein the means for determining the second activation comprises:

means for receiving the first activation via the fabric from the at least one prior neuron;

means for computing the second activation based at least in part on the first activation and the first weights implemented by at least executing first instructions stored in the first memory and accessing the first weights in the second memory; and

means for selectively transmitting the second activation via the fabric to the at least one subsequent neuron.

22. The system of claim 21, wherein the means for determining the fourth activation comprises:

means for receiving the third activation via the fabric from the at least one prior neuron;

means for computing the fourth activation based at least in part on the third activation and the selected weights implemented by at least executing the first instructions and accessing the selected weights in the second memory; and

means for selectively transmitting the fourth activation via the fabric to the at least one subsequent neuron.

23. The system of claim 22, wherein the means for determining and saving the second weights comprises:

means for receiving the first delta that is partially based on the second activation via the fabric from the at least one subsequent neuron;

means for computing a first gradient based at least in part on the first delta and the second activation implemented by at least executing second instructions stored in the first memory;

means for computing the second weights based at least in part on the first gradient, a learning rule, and the first weights implemented by at least executing third instructions stored in the first memory and accessing the first weights in the second memory; and

means for storing the second weights in the second memory.

24. The system of claim 23, wherein the means for determining and saving the third weights comprises:

means for receiving the second delta that is partially based on the fourth activation via the fabric from the at least one subsequent neuron;

means for computing a second gradient based at least in part on a third delta and the fourth activation implemented by at least executing the second instructions stored in the first memory;

means for computing the third weights based at least in part on the second gradient, the learning rule and the selected weights implemented by at least executing the third instructions stored in the first memory and accessing the selected weights in the second memory; and

means for storing the third weights in the second memory.

25. The system of claim 24, wherein the means for computing the second gradient additionally comprises means for recomputing the fourth activation based at least in part upon the selected weights.

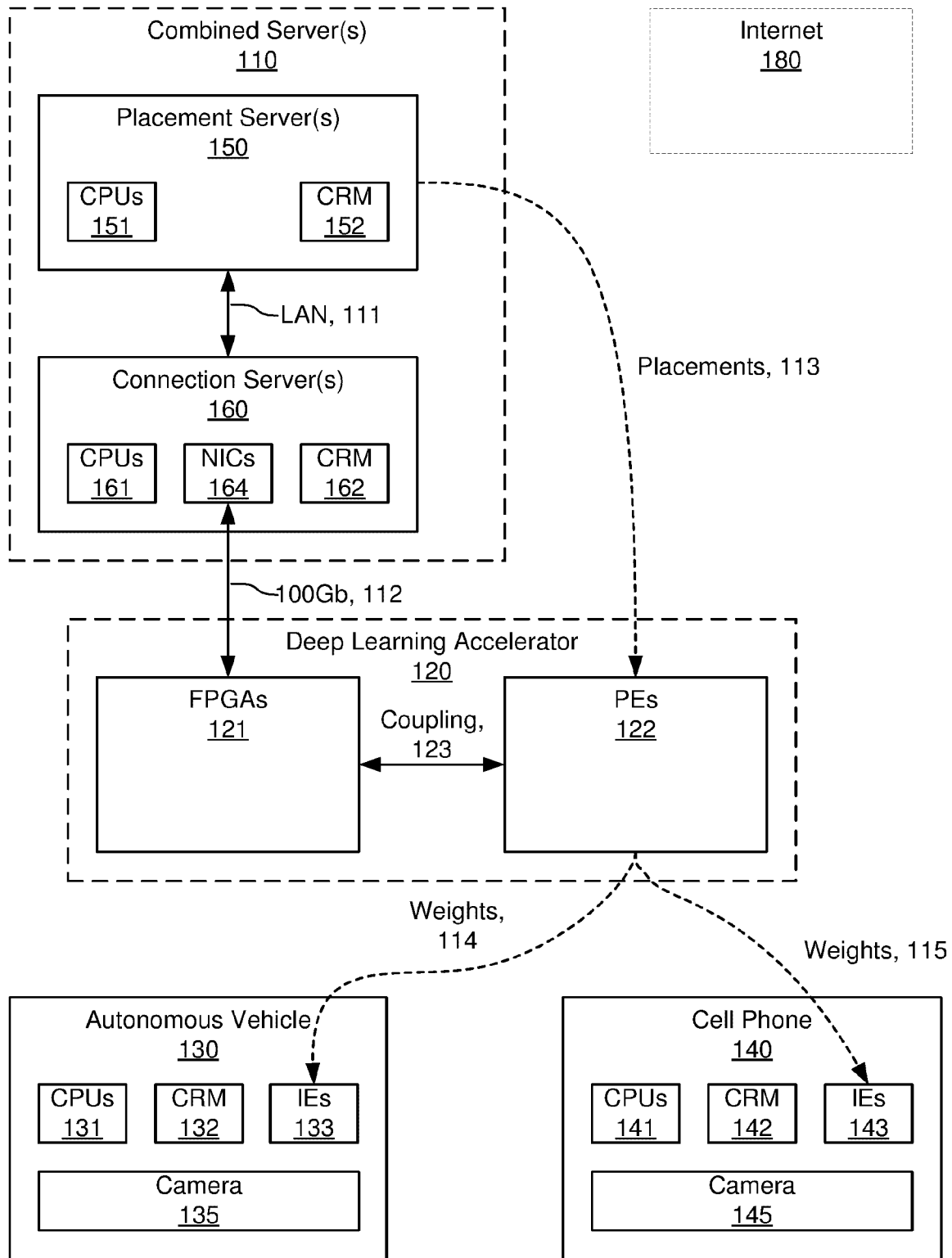


Fig. 1

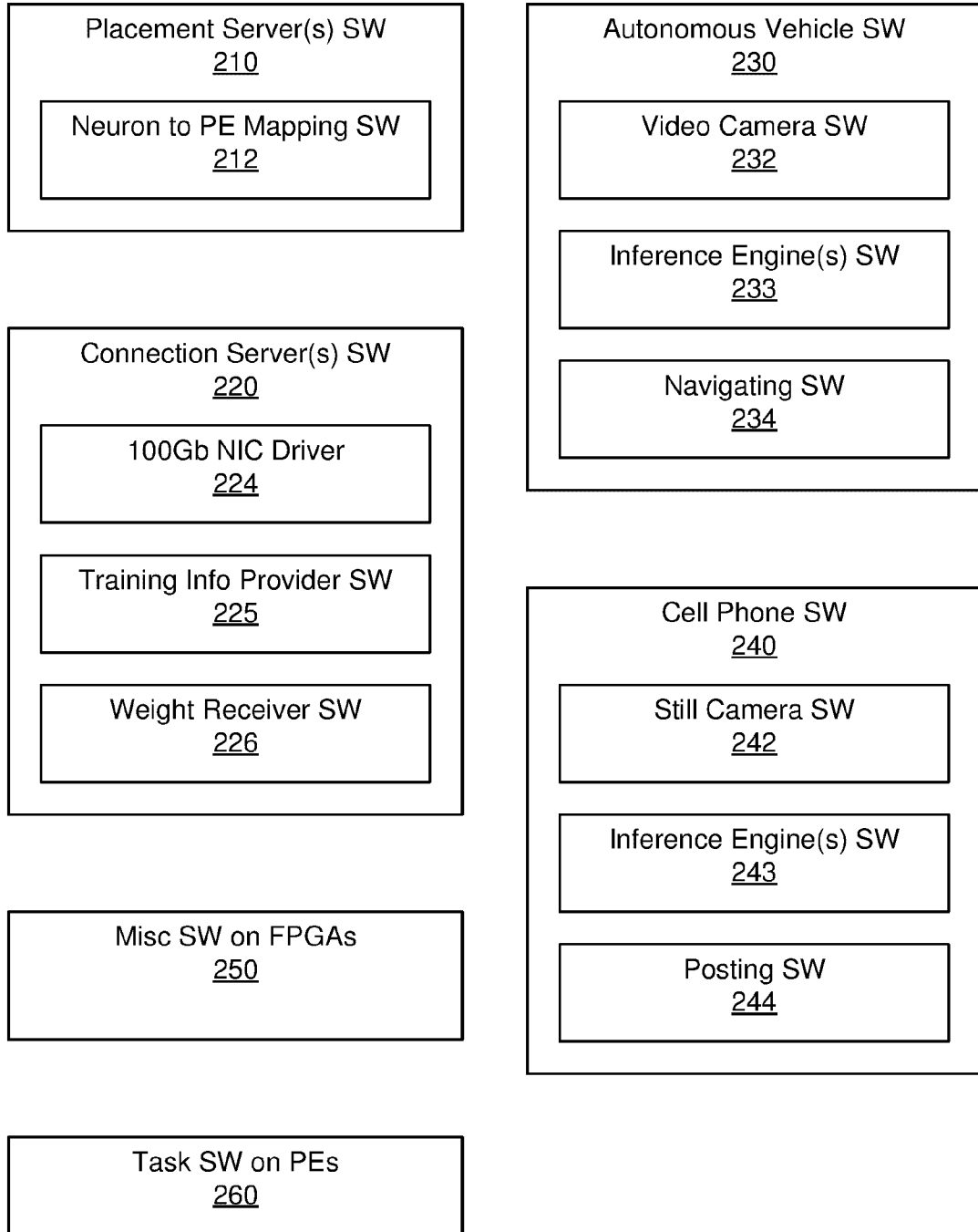


Fig. 2

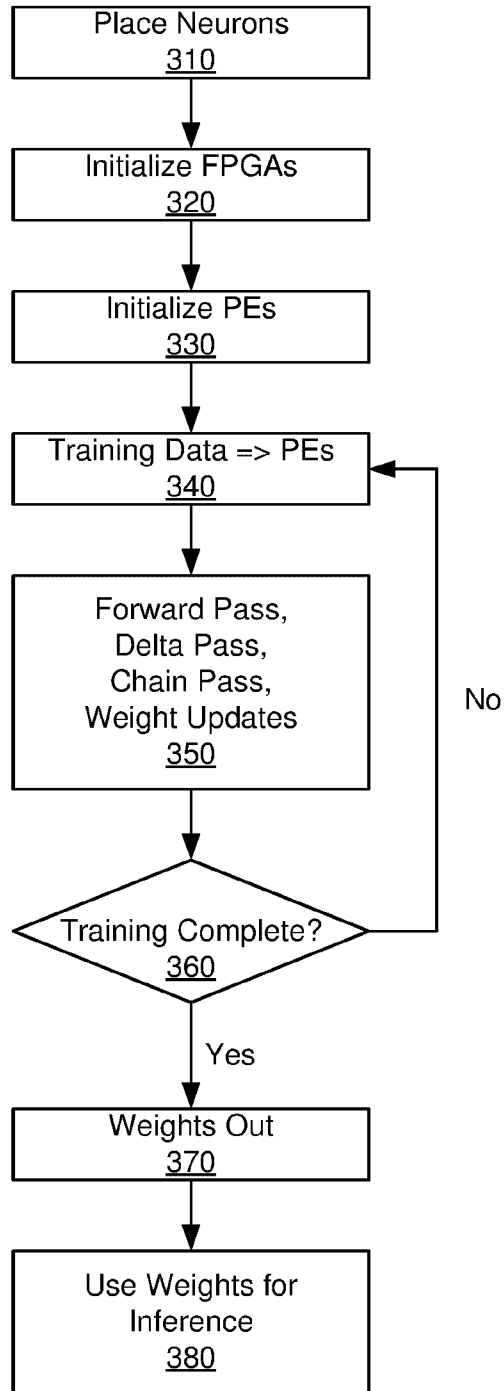


Fig. 3

Deep Learning Accelerator, 400

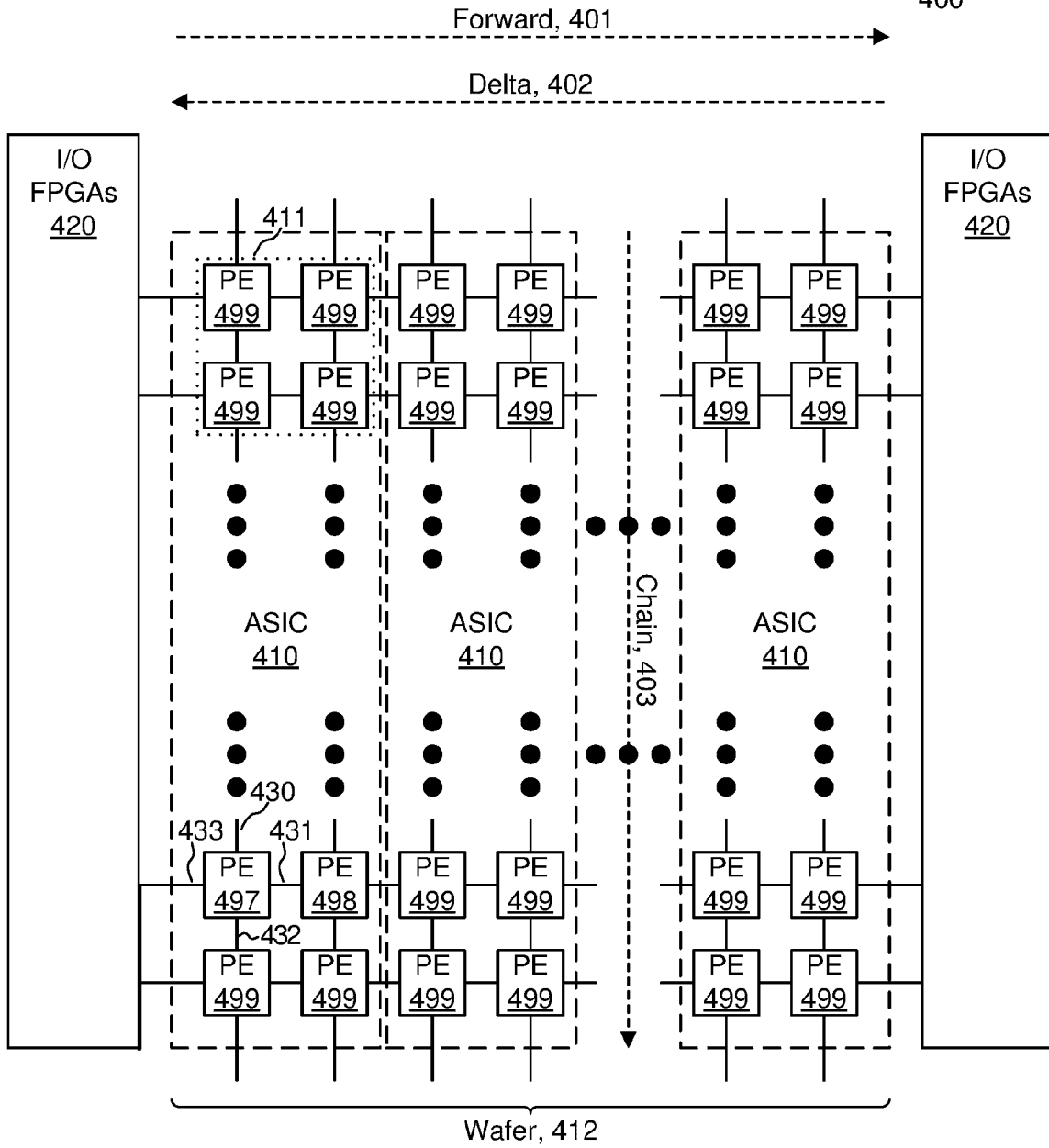


Fig. 4

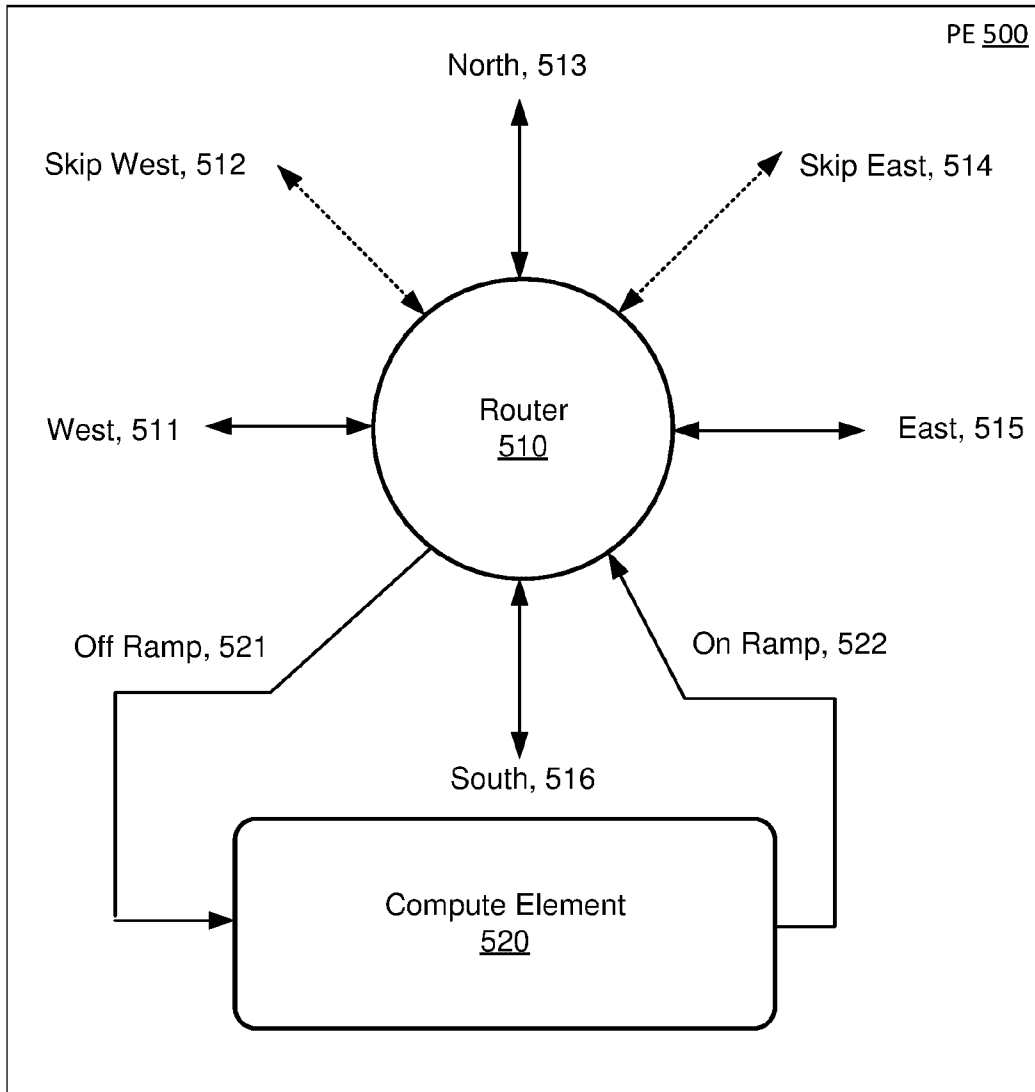


Fig. 5

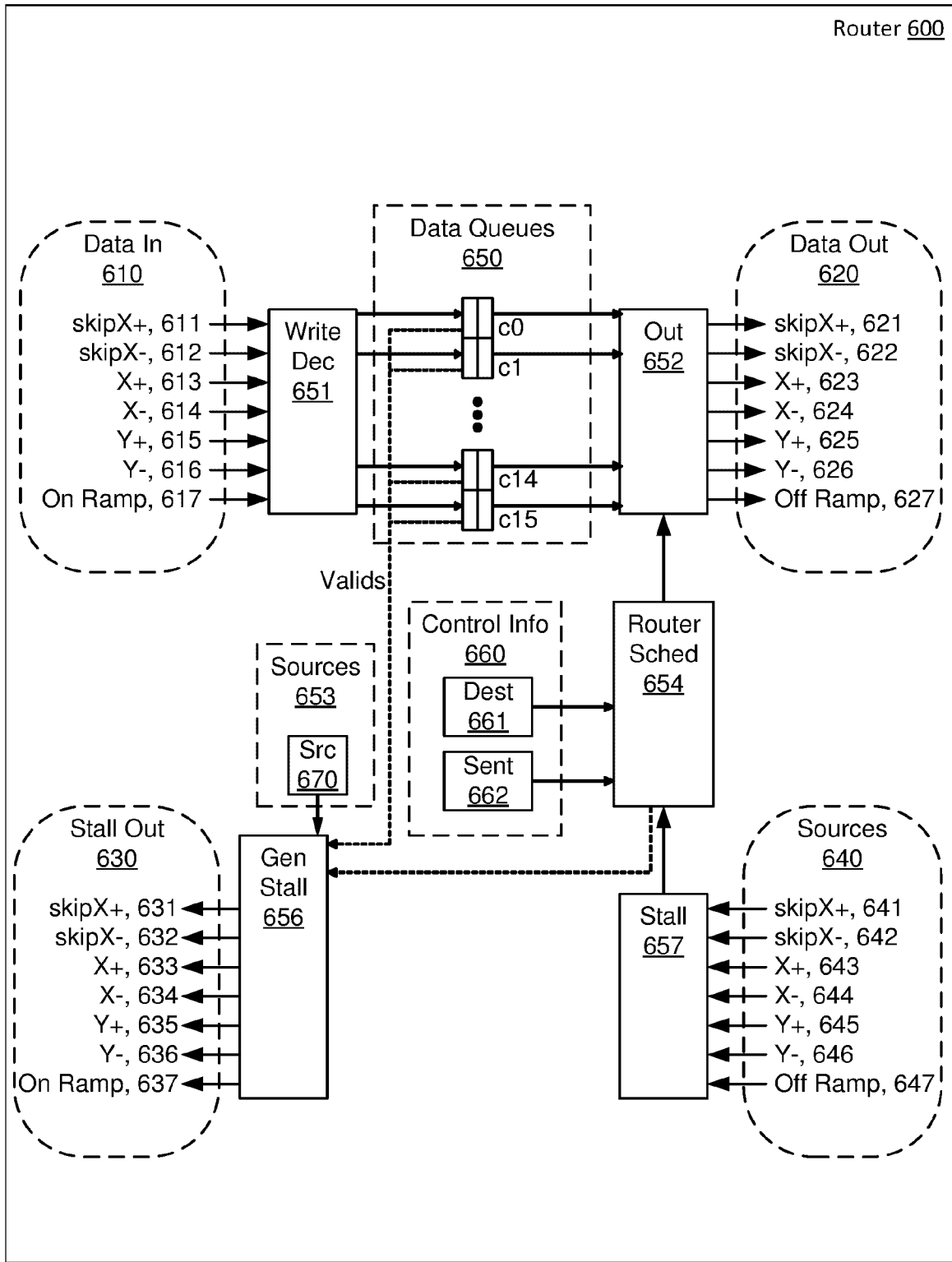


Fig. 6

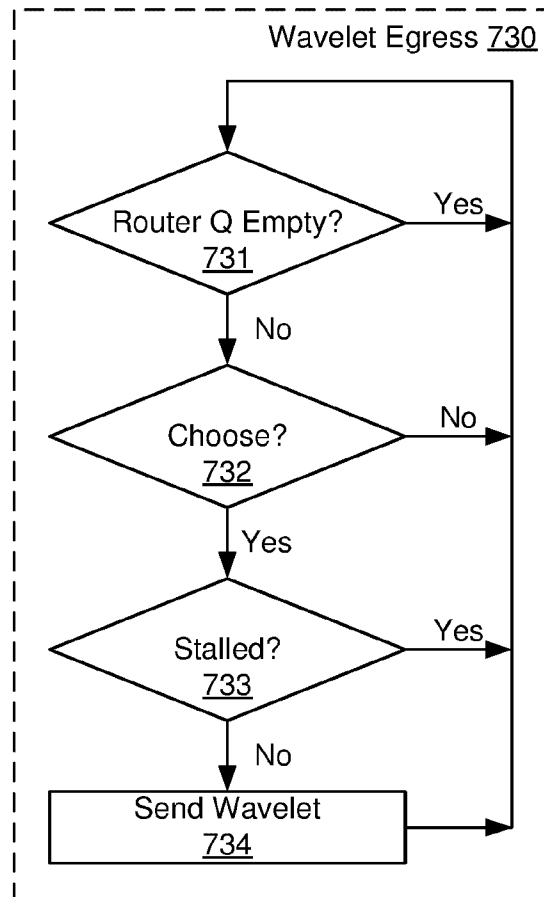
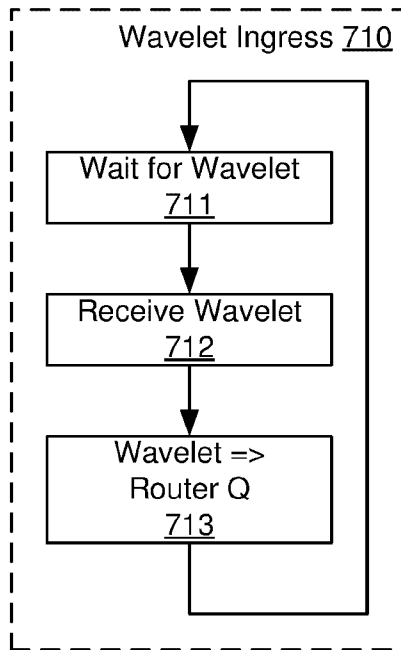
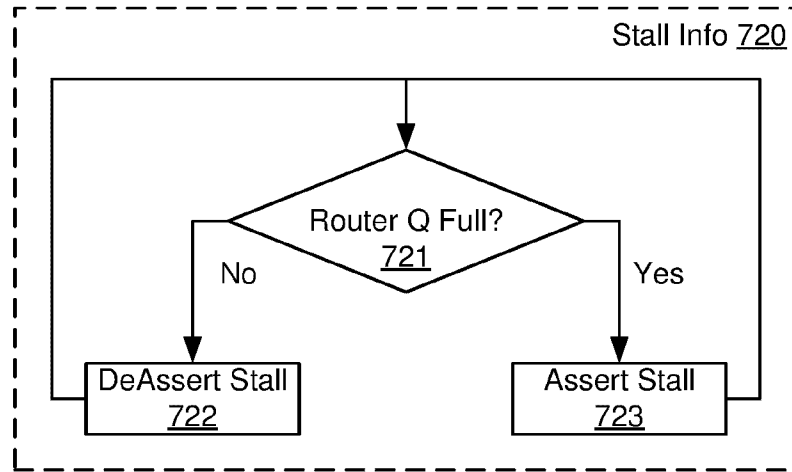


Fig. 7

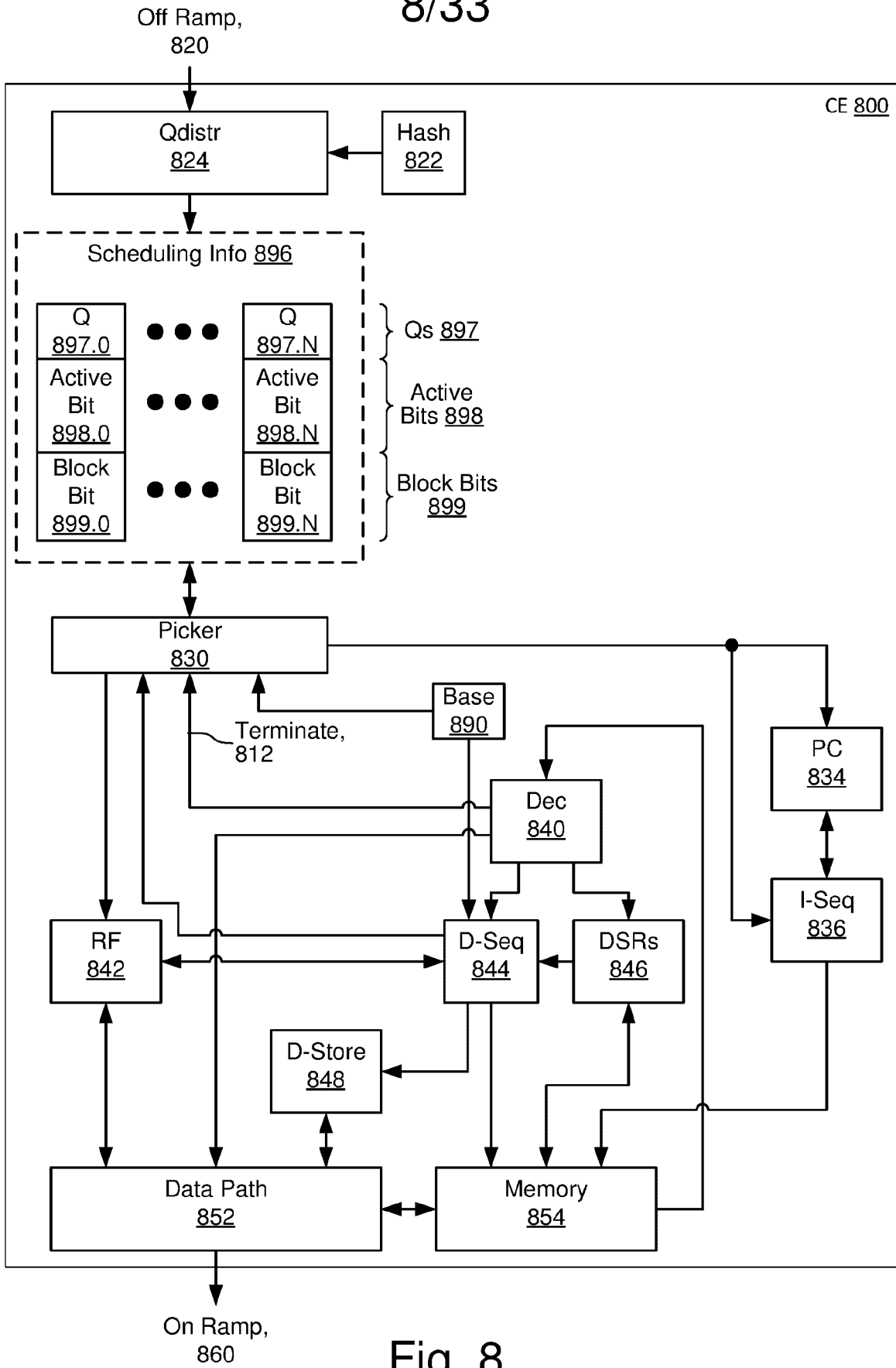


Fig. 8

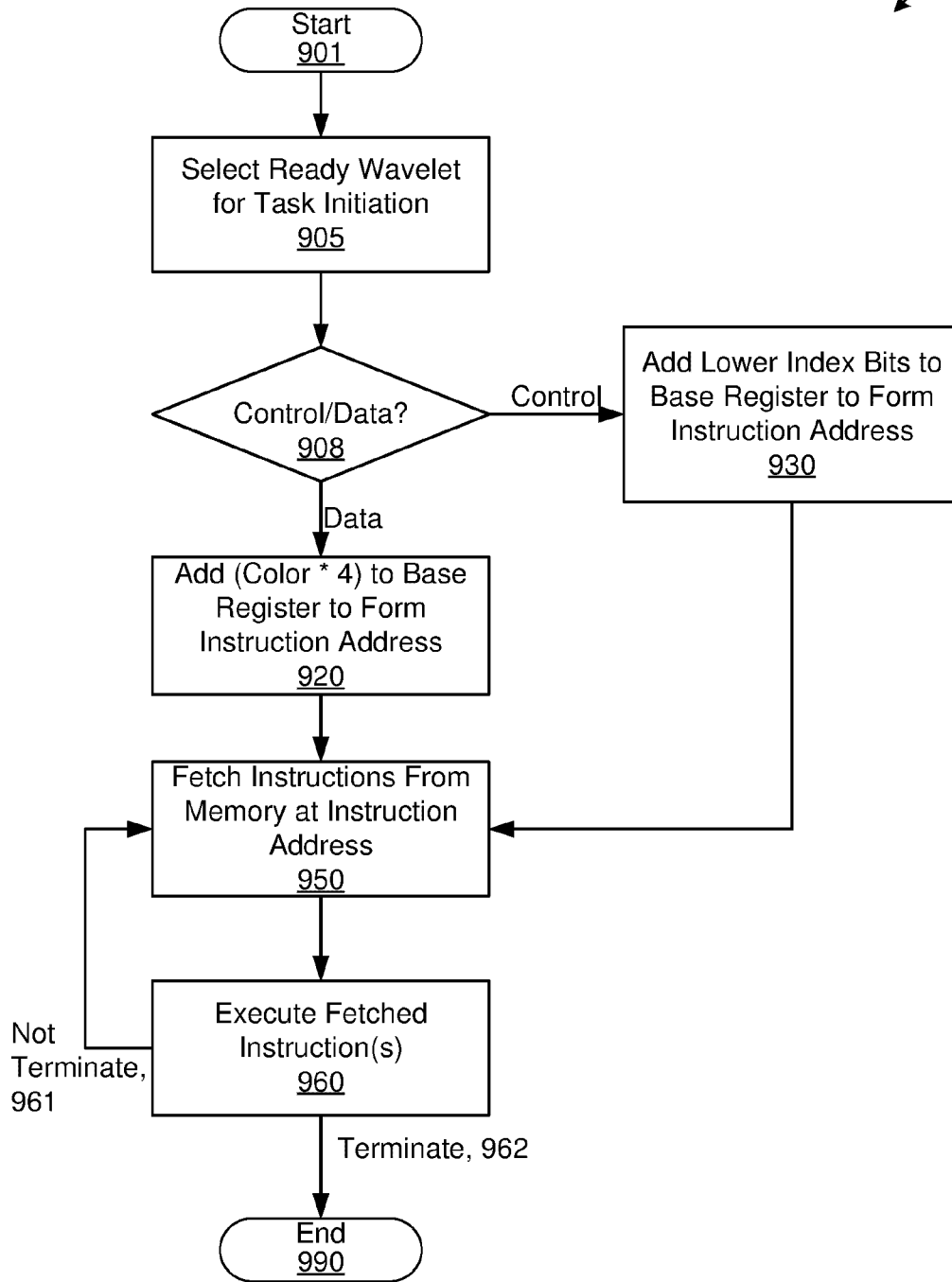


Fig. 9

10/33

1000

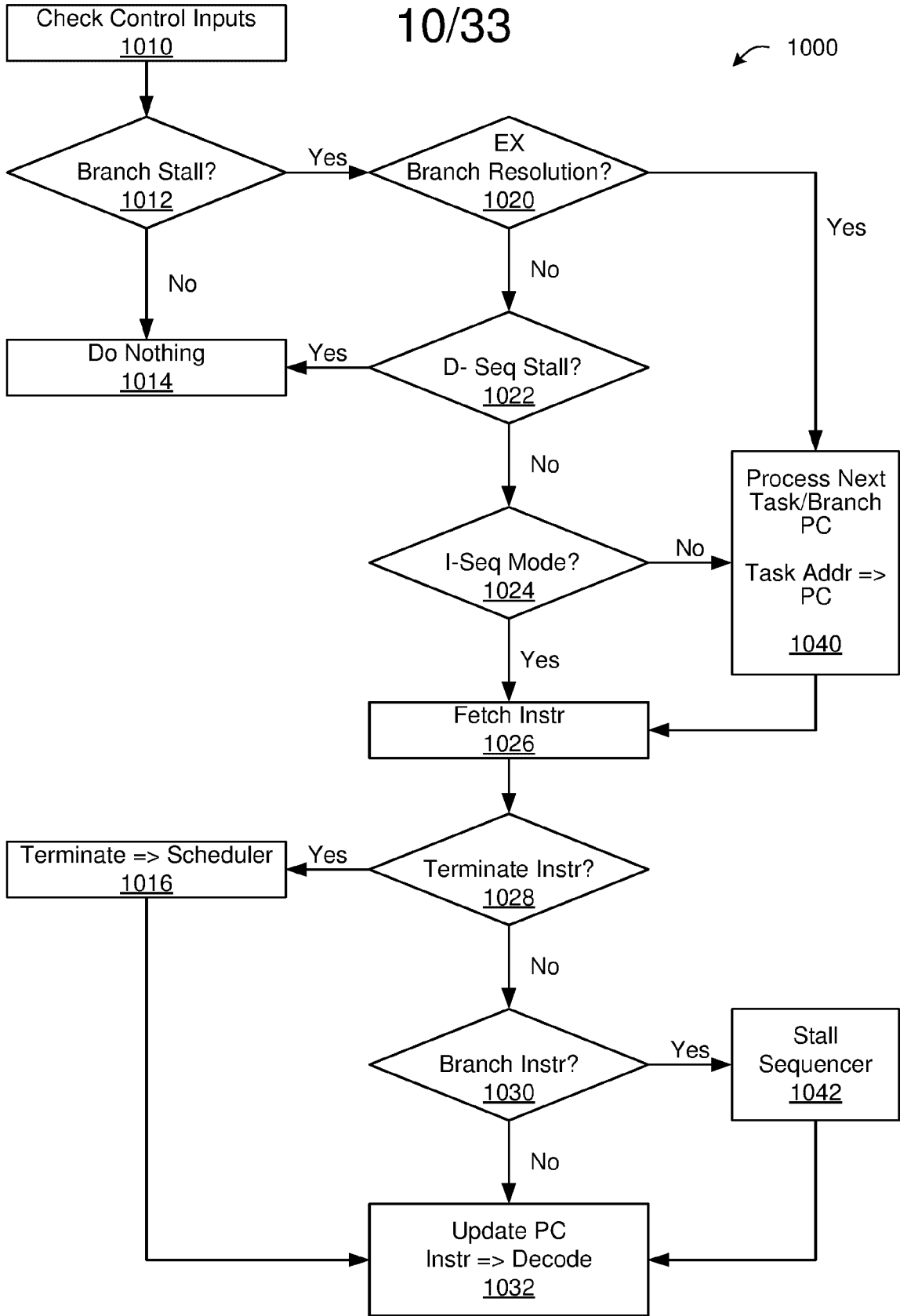


Fig. 10

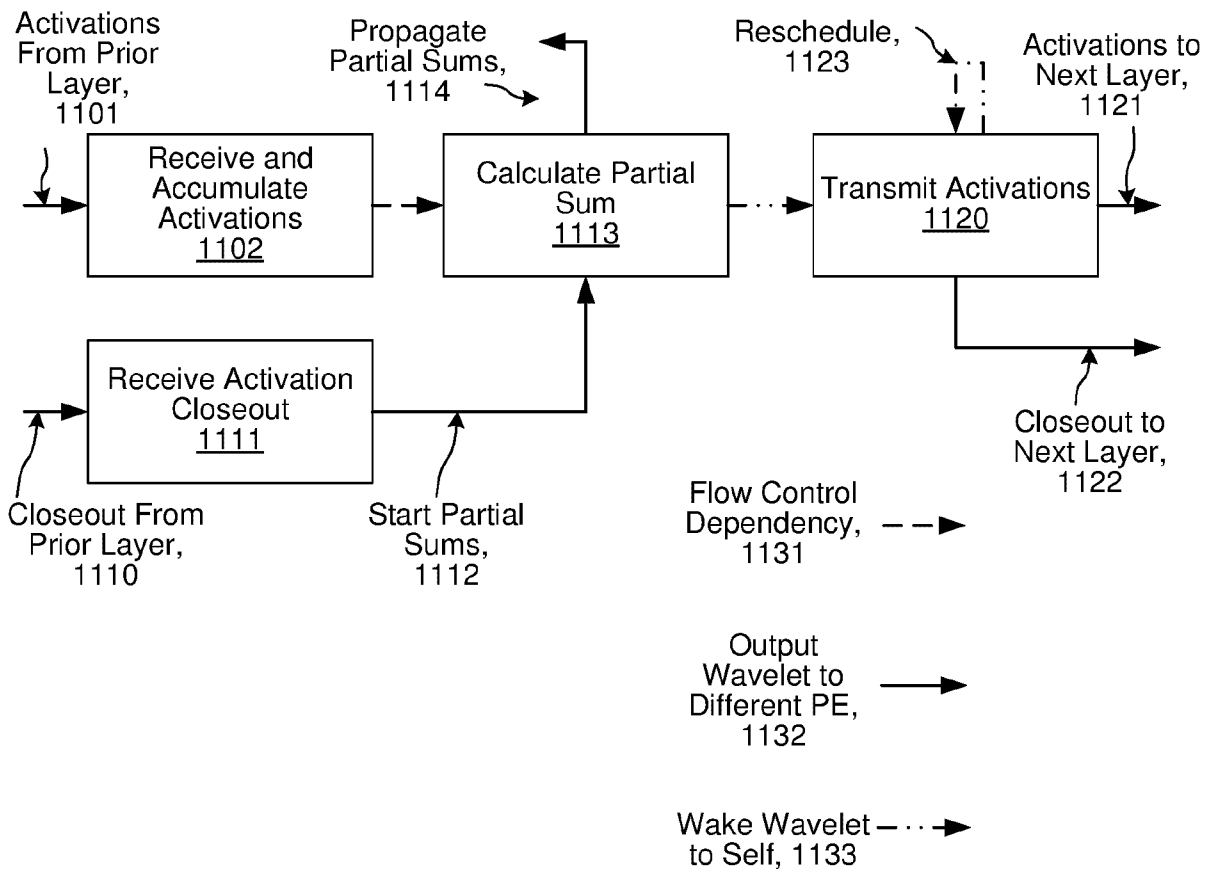


Fig. 11

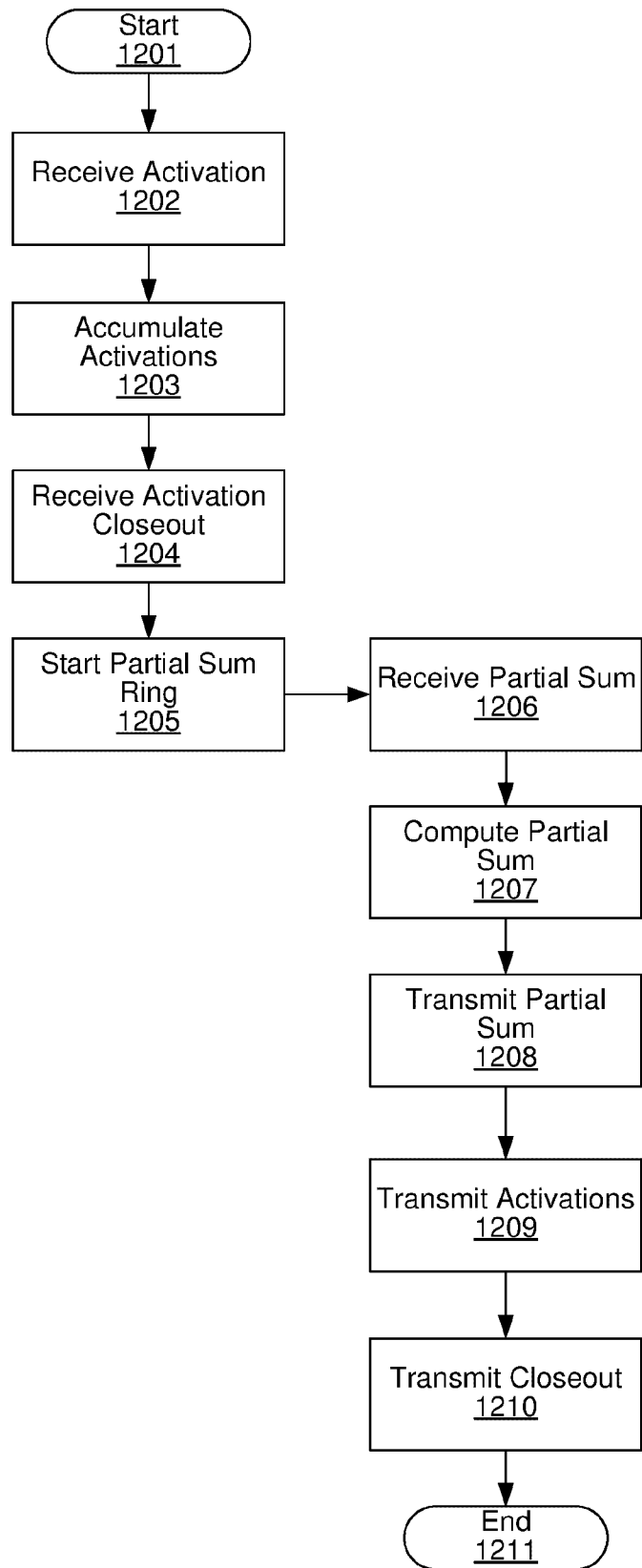


Fig. 12

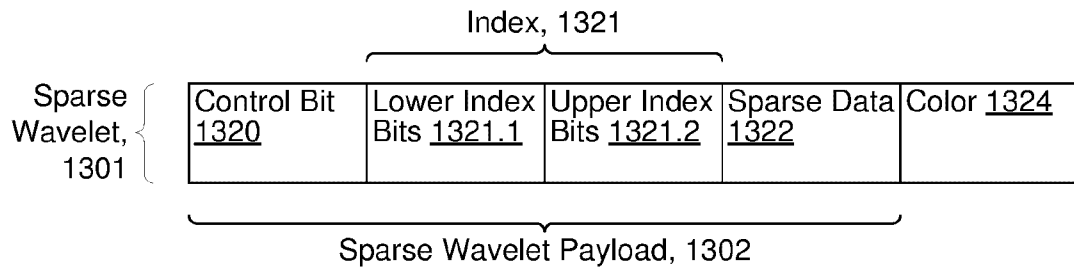


Fig. 13A

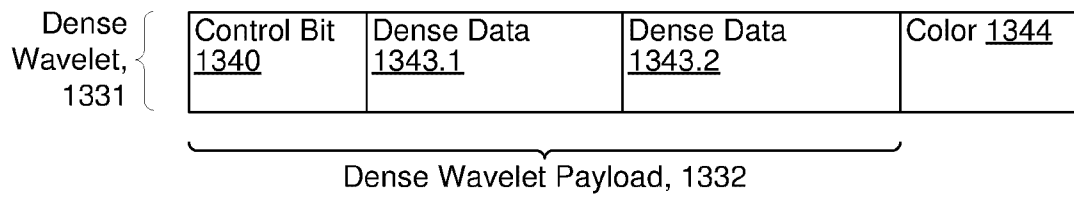


Fig. 13B

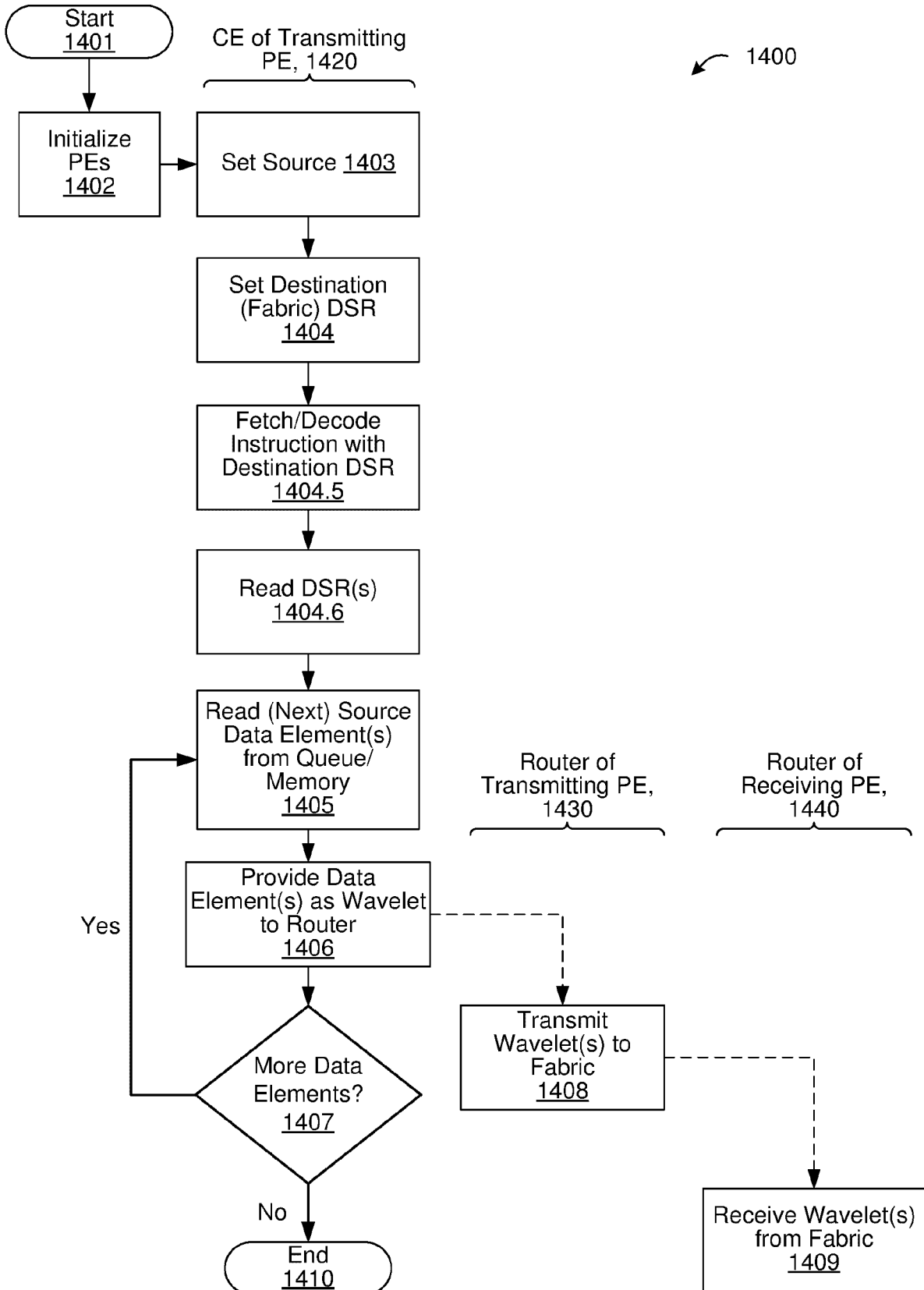


Fig. 14

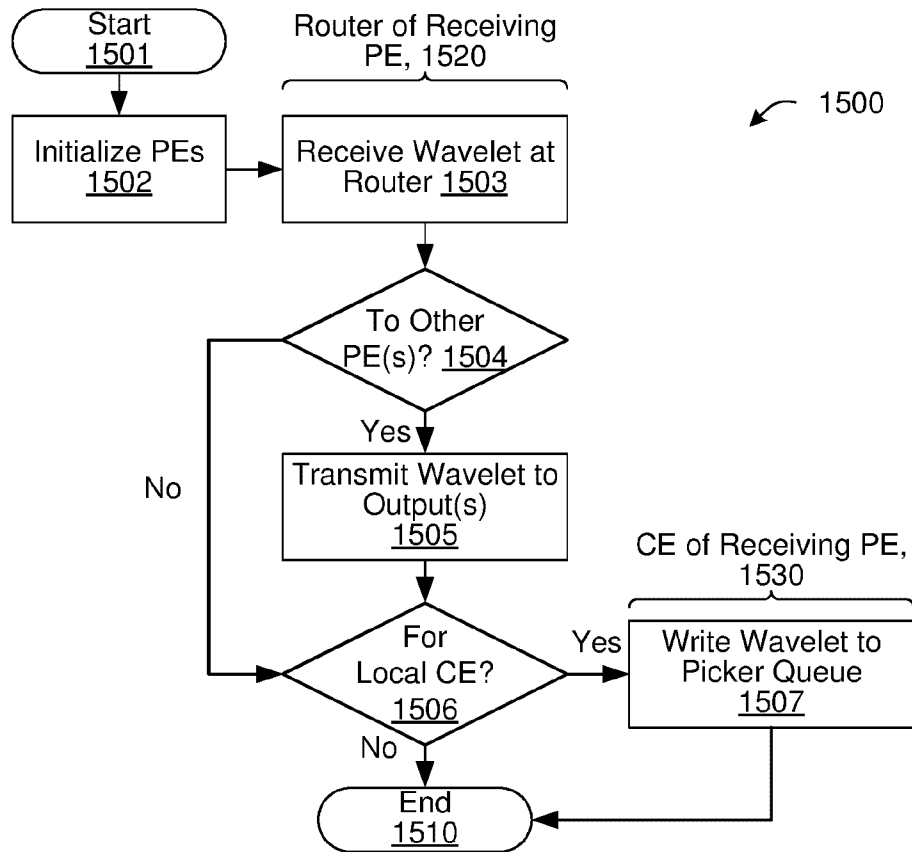


Fig. 15A

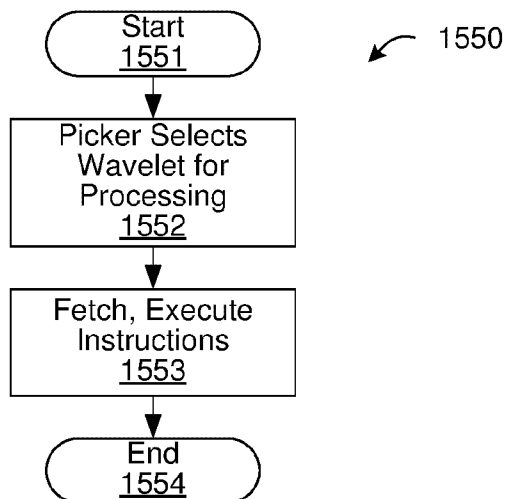


Fig. 15B

16/33

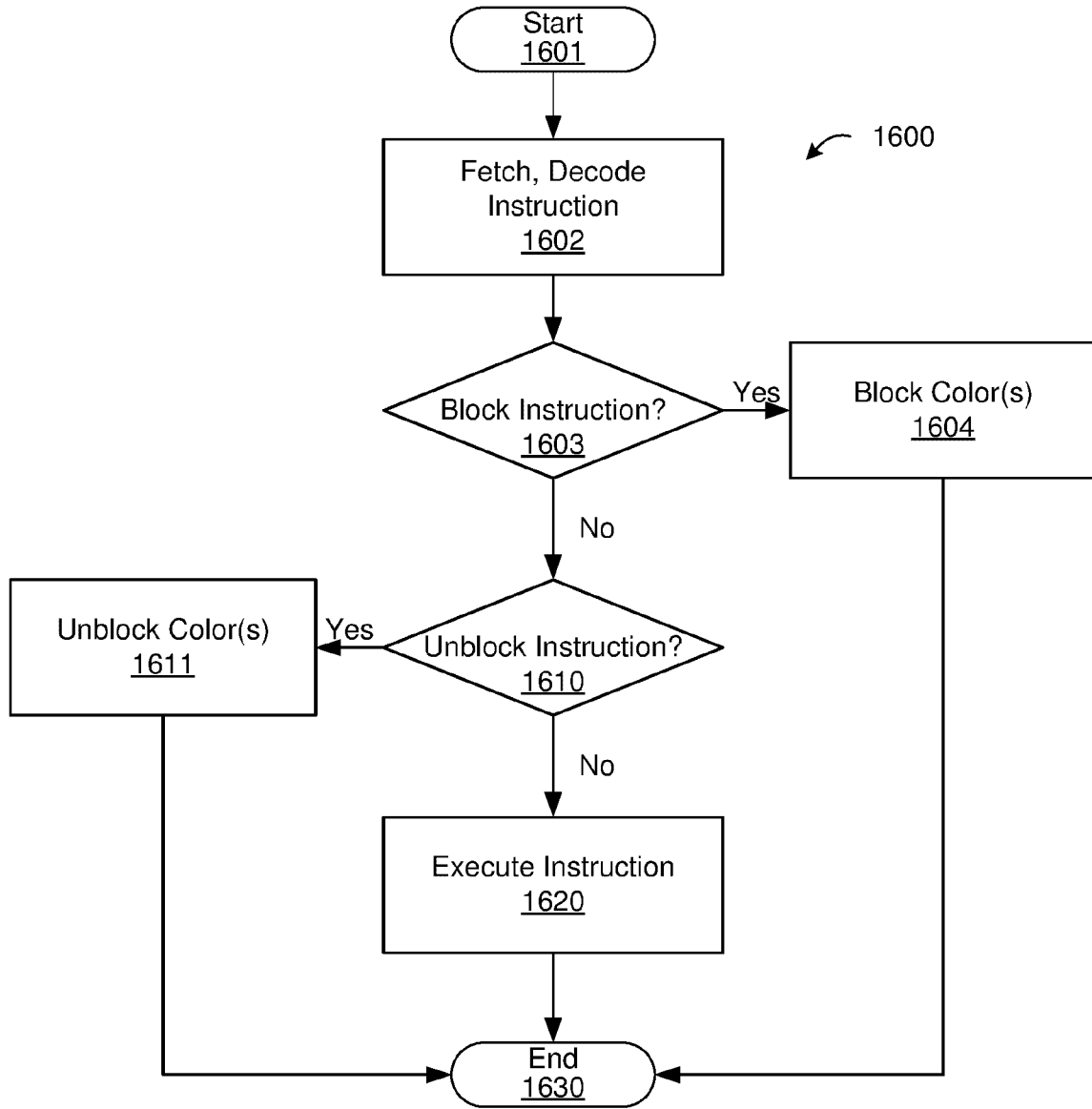


Fig. 16

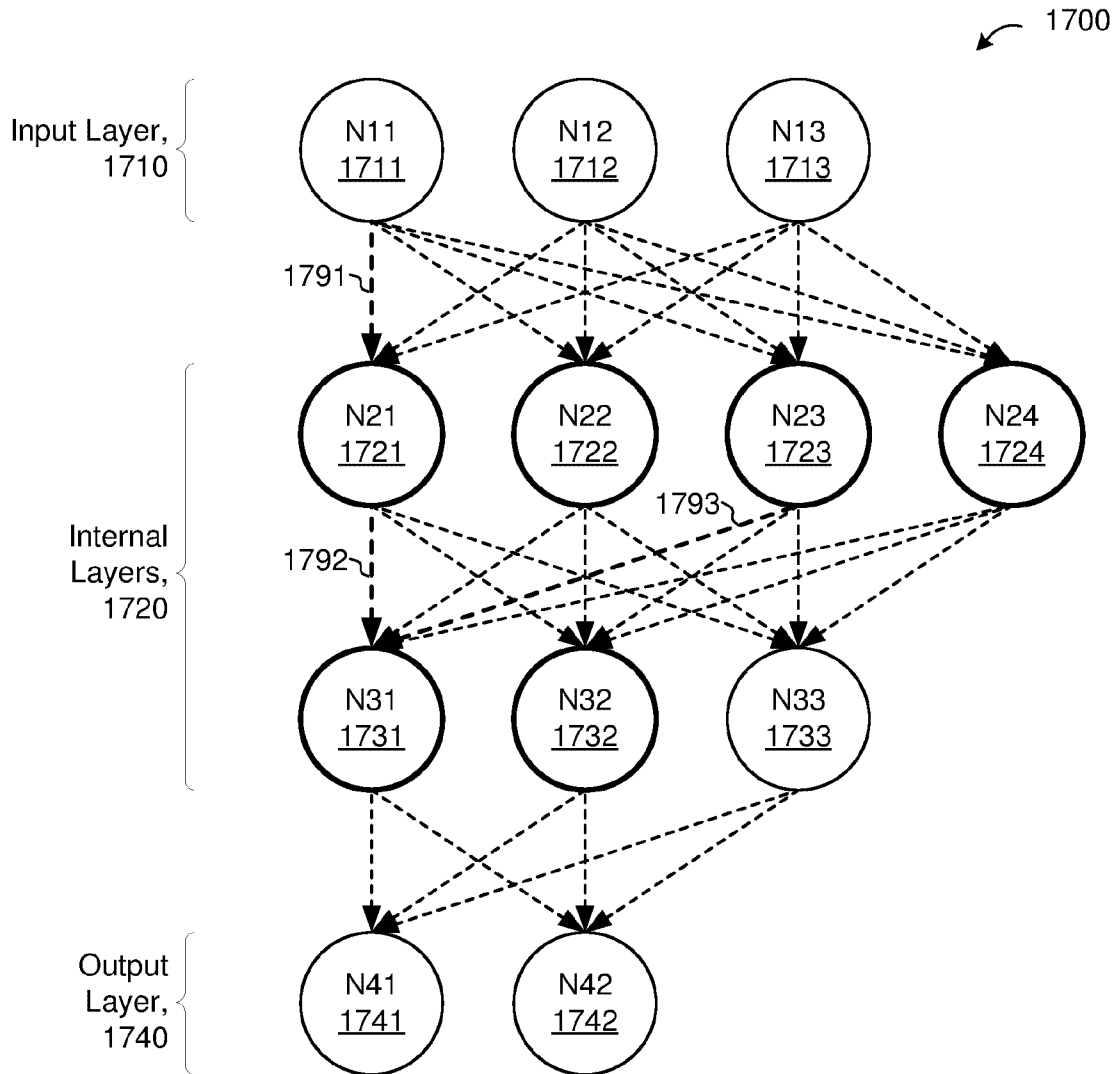


Fig. 17

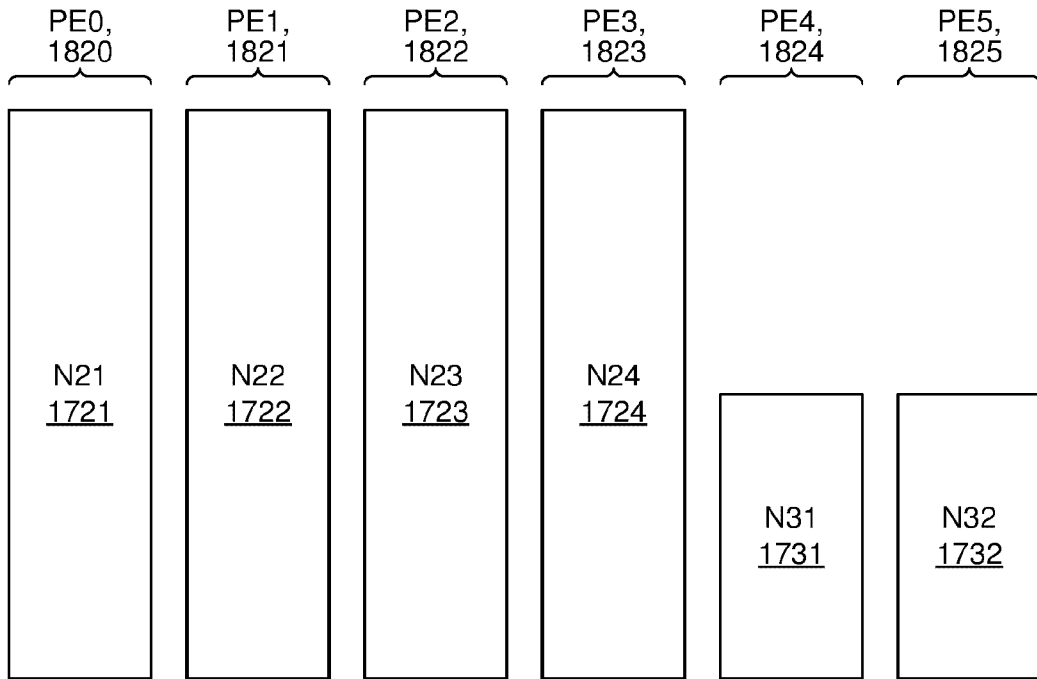


Fig. 18A

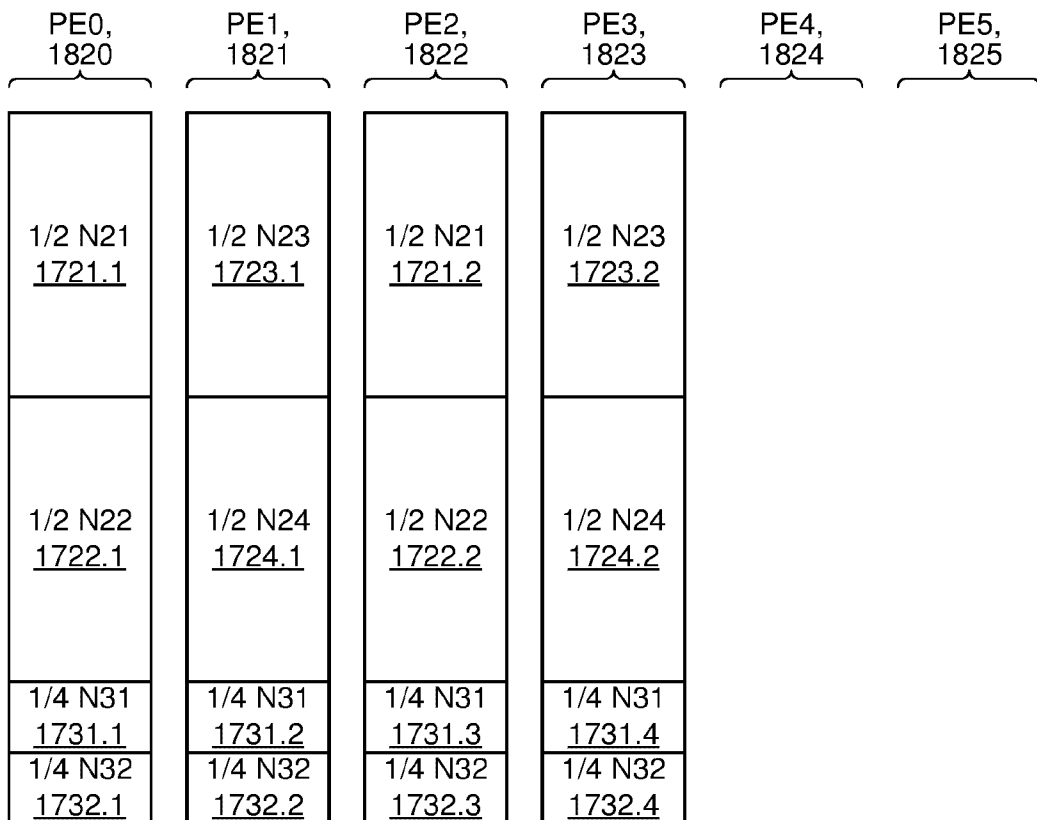


Fig. 18B

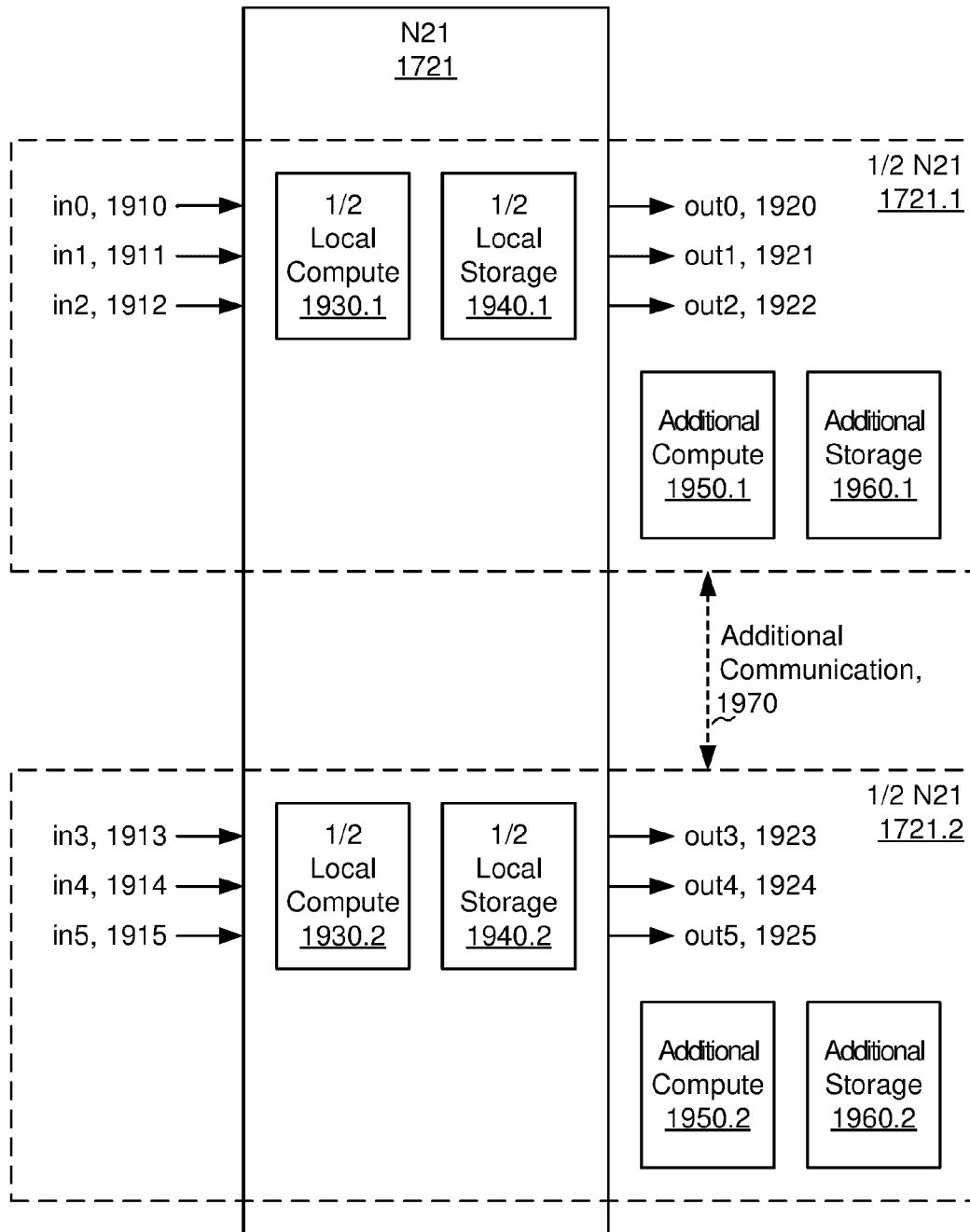


Fig. 19

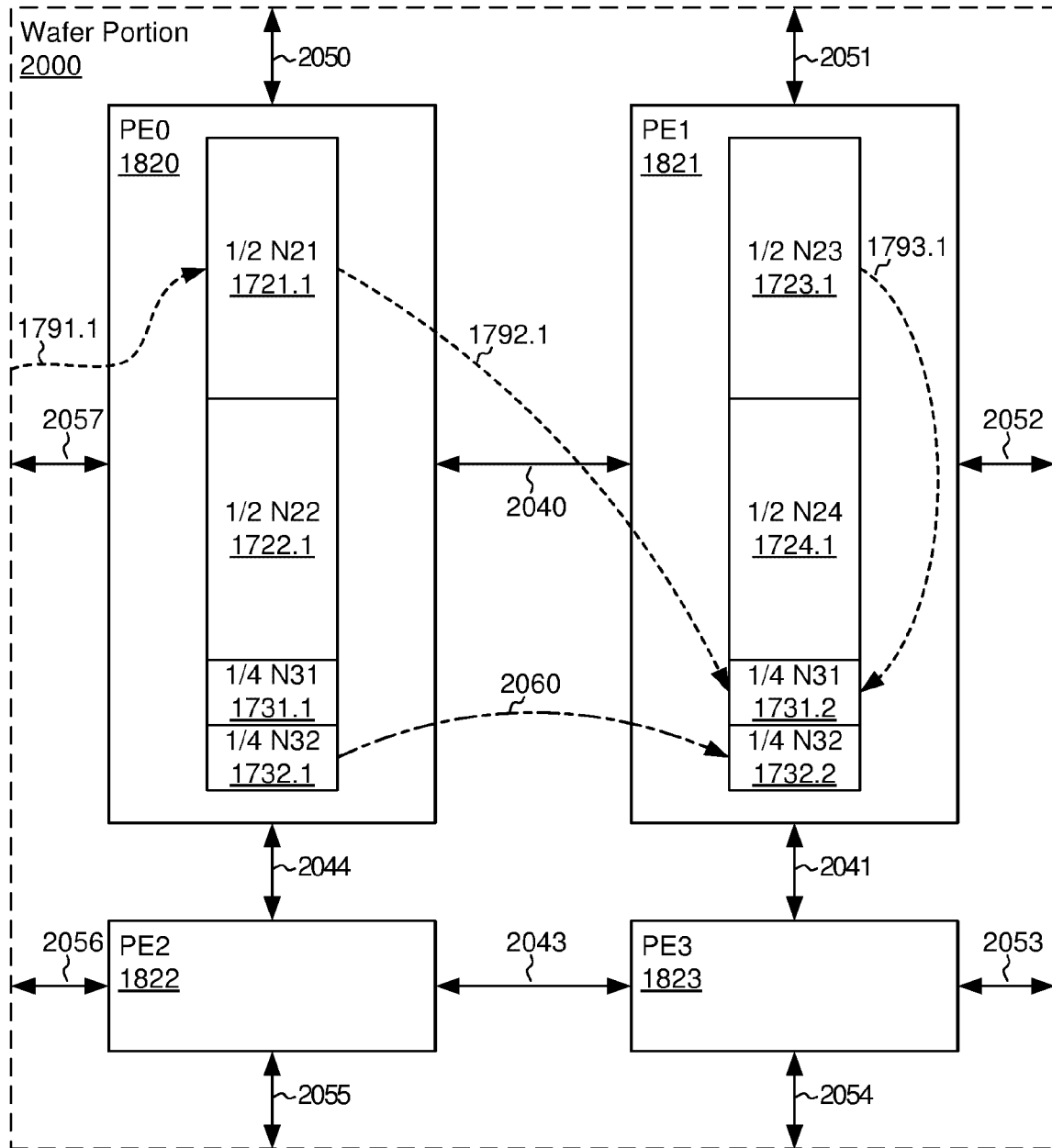


Fig. 20

21/33

CH 2114	SQ 2113	SC 2112	SA 2111	SS 2110	Type 2109	US 2108	CX 2107	Term 2106	AC 2105	SW 2104	UE 2103	UTID 2102	Length 2101
------------	------------	------------	------------	------------	--------------	------------	------------	--------------	------------	------------	------------	--------------	----------------

Fabric Input Data Structure Descriptor, 2100

Fig. 21A

AC 2125	Index High 2128.2	WLI 2132	SA 2131	SS 2130	Type 2129	Index Low 2128.1	C 2127	Color 2126	SW 2124	UE 2123	UTID 2122	Length 2121
------------	----------------------	-------------	------------	------------	--------------	---------------------	-----------	---------------	------------	------------	--------------	----------------

Fabric Output Data Structure Descriptor, 2120

Fig. 21B

Stride 2153	WLI 2152	SA 2151	SS 2150	Type 2149	Base Address 2142	Length 2141
----------------	-------------	------------	------------	--------------	----------------------	----------------

1D Memory Vector Data Structure Descriptor, 2140

Fig. 21C

Length, 2161						
Length Upper Bits 2161.2	WLI 2172	SA 2171	SS 2170	Type 2169	Base Address 2162	Length Lower Bits 2161.1

4D Memory Vector Data Structure Descriptor, 2160

Fig. 21D

SW 2184	WLI 2192	SA 2191	SS 2190	Type 2189	FW 2188	Base Address 2182	Length 2181
------------	-------------	------------	------------	--------------	------------	----------------------	----------------

Circular Memory Buffer Data Structure Descriptor, 2180

Fig. 21E

Pop Color <u>2216</u>	Push Color <u>2215</u>	FIFO <u>2214</u>	End Address <u>2213</u>	Start Address <u>2212</u>	Type <u>2211</u>
--------------------------	---------------------------	---------------------	----------------------------	------------------------------	------------------

Circular Memory Buffer Extended Data Structure Descriptor, 2210

Fig. 22A

Stride <u>2245</u>	Stride Select 4 <u>2244.4</u>	Stride Select 3 <u>2244.3</u>	Stride Select 2 <u>2244.2</u>	Stride Select 1 <u>2244.1</u>	DF <u>2243</u>	Dimensions <u>2242</u>	Type <u>2241</u>
-----------------------	----------------------------------	----------------------------------	----------------------------------	----------------------------------	-------------------	---------------------------	---------------------

4D Memory Vector Extended Data Structure Descriptor, 2240

Fig. 22B

23/33

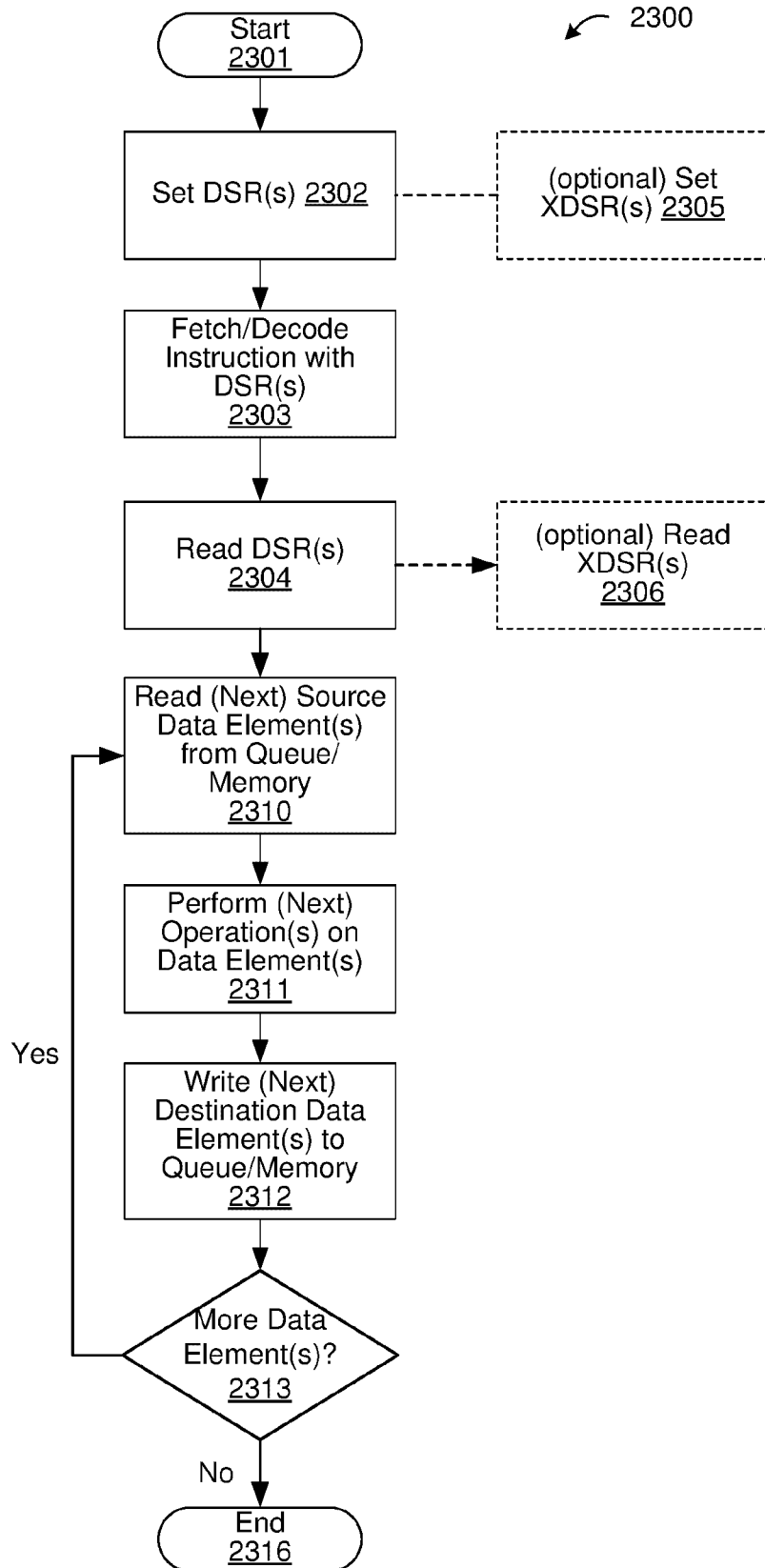


Fig. 23

24/33

2400

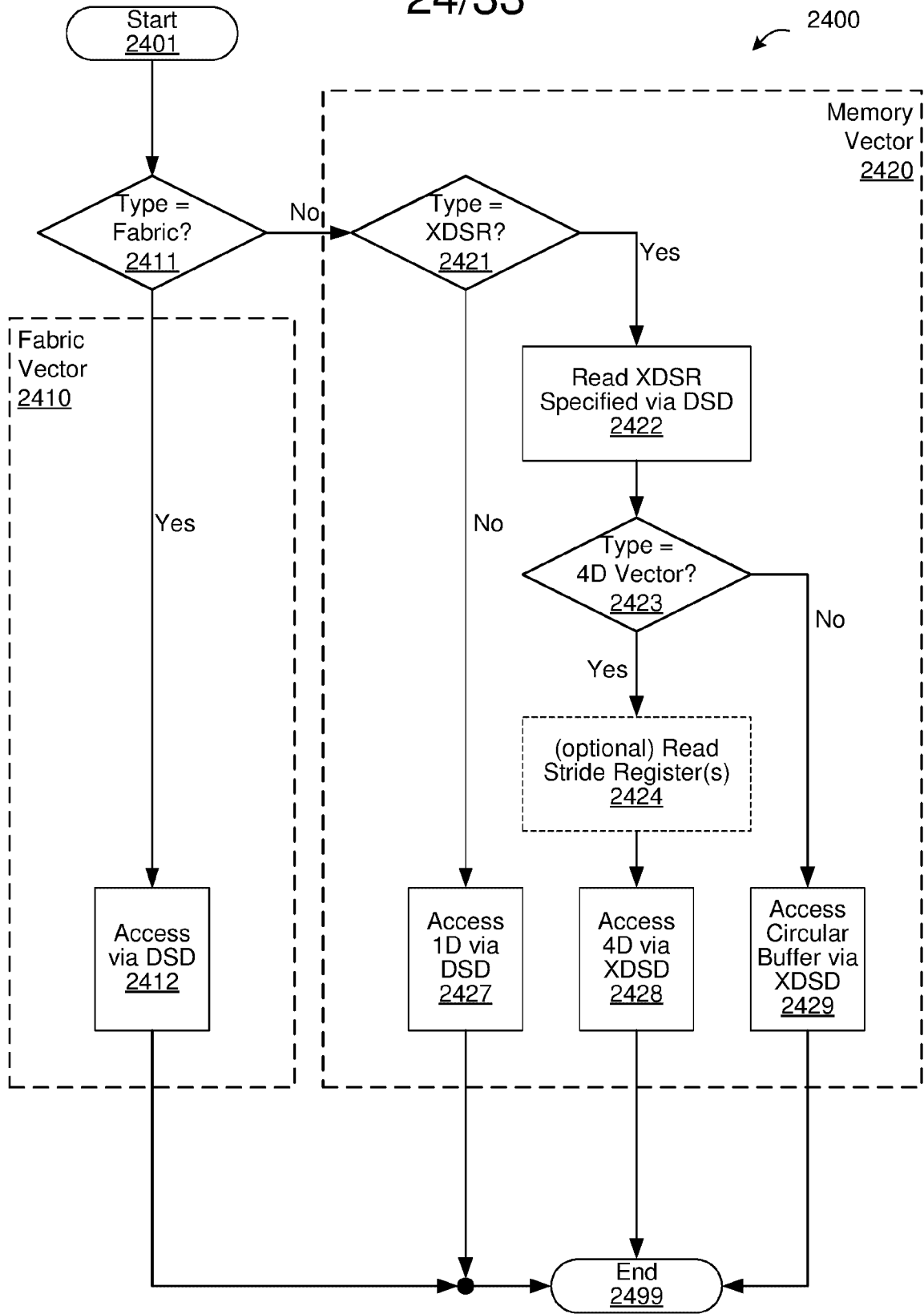


Fig. 24

25/33

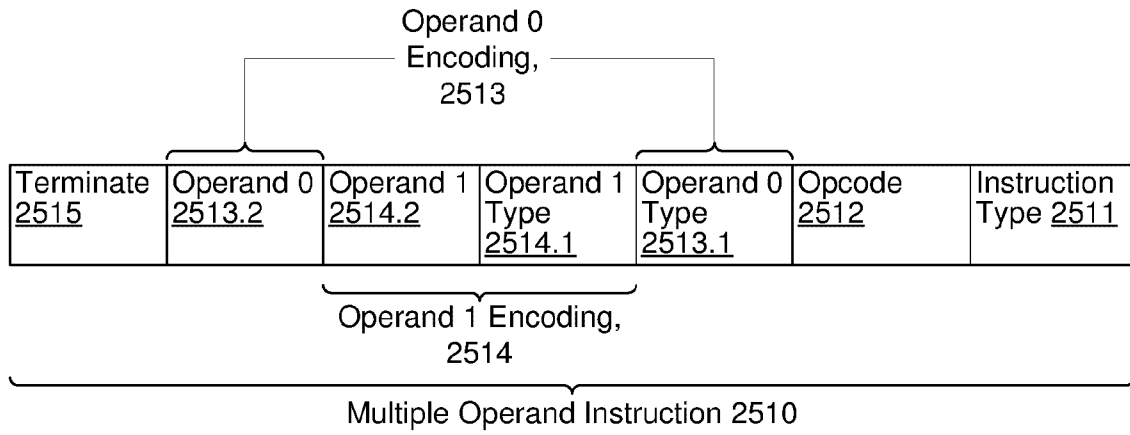


Fig. 25A

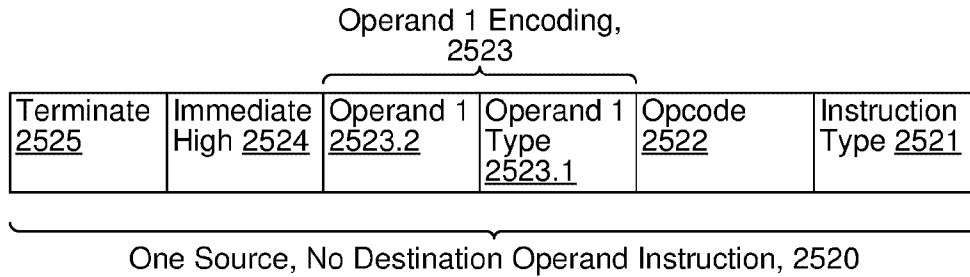


Fig. 25B

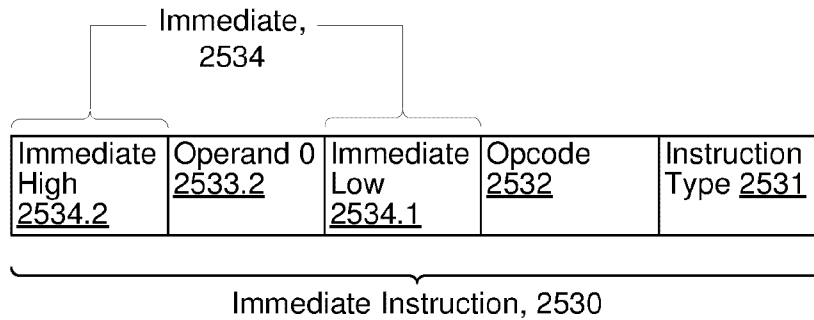


Fig. 25C

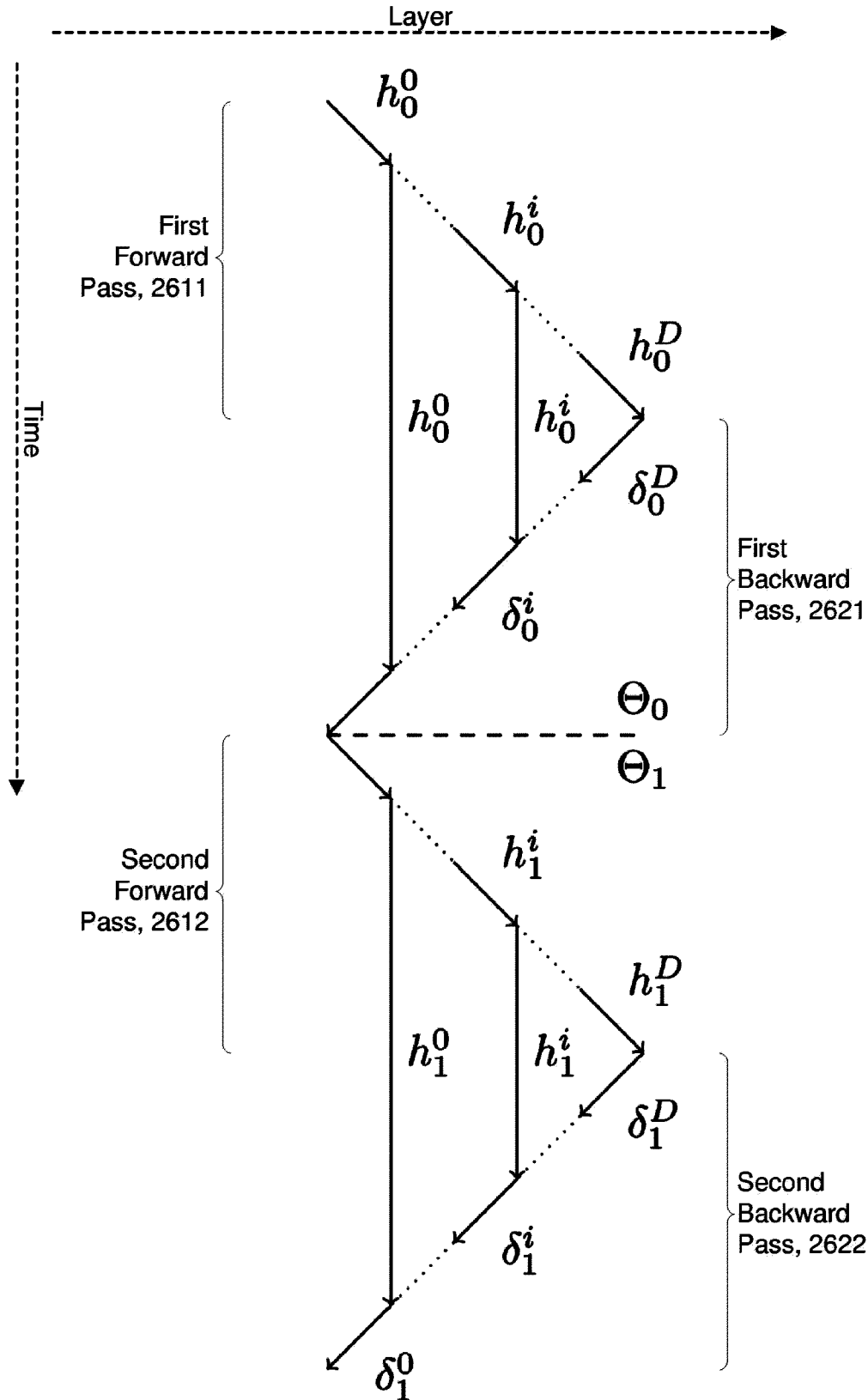


Fig. 26A

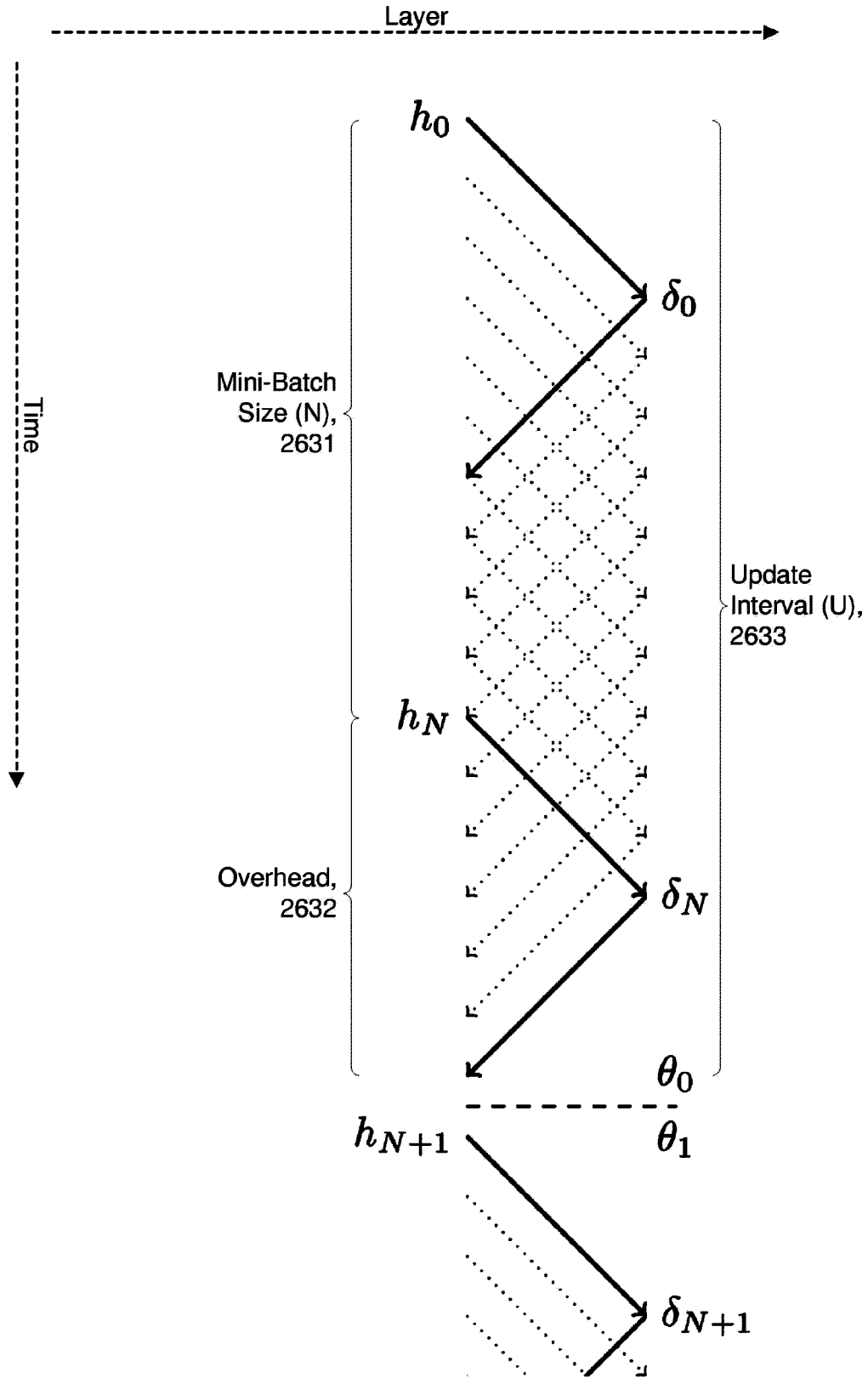


Fig. 26B

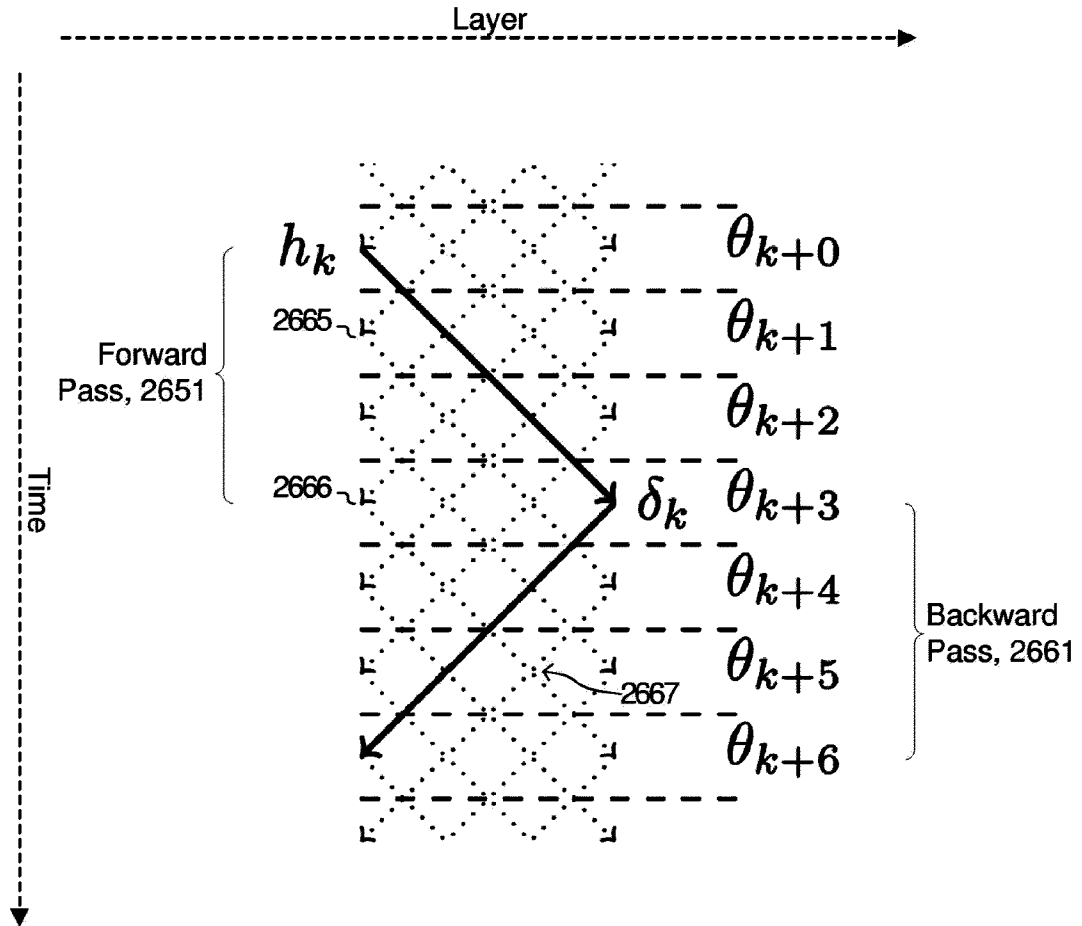


Fig. 26C

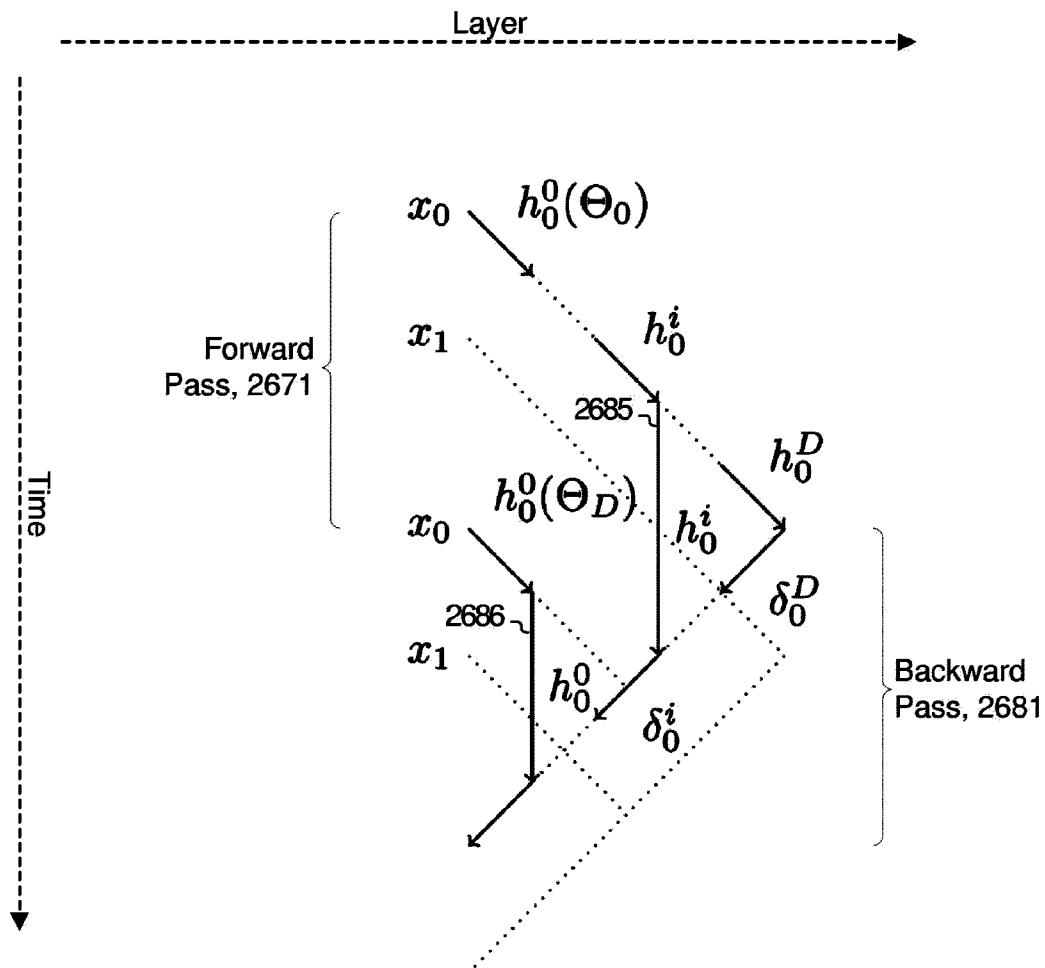


Fig. 26D

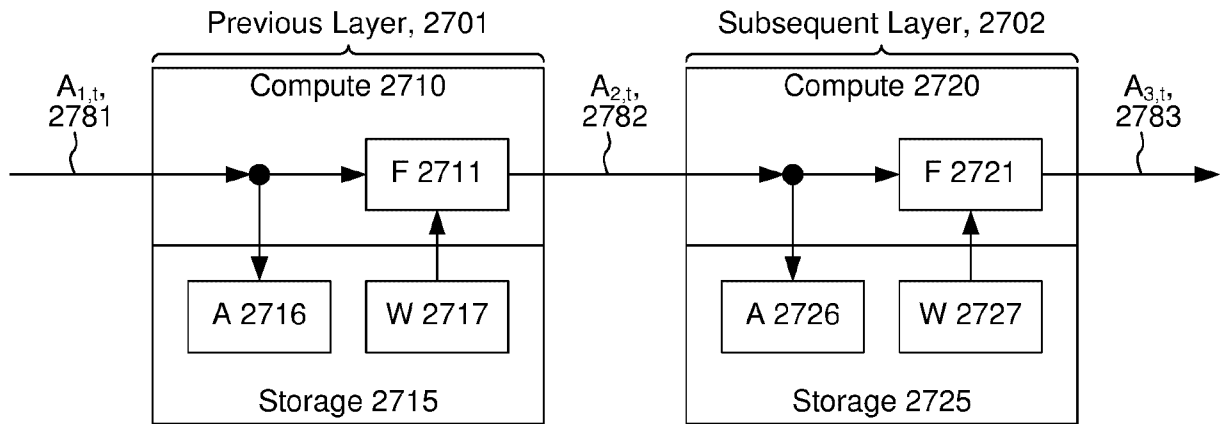


Fig. 27A

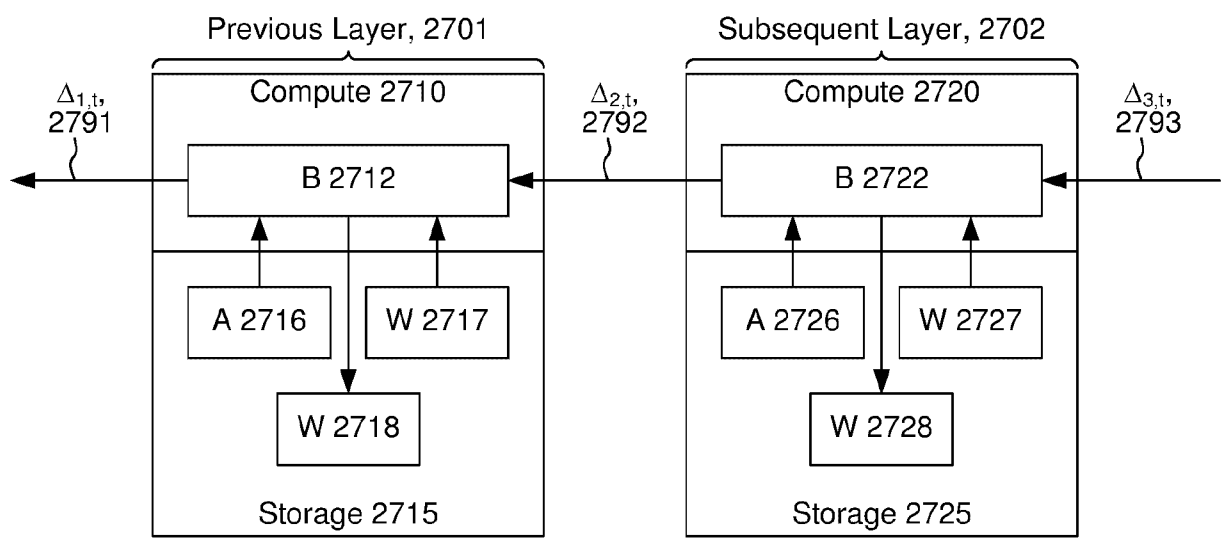


Fig. 27B

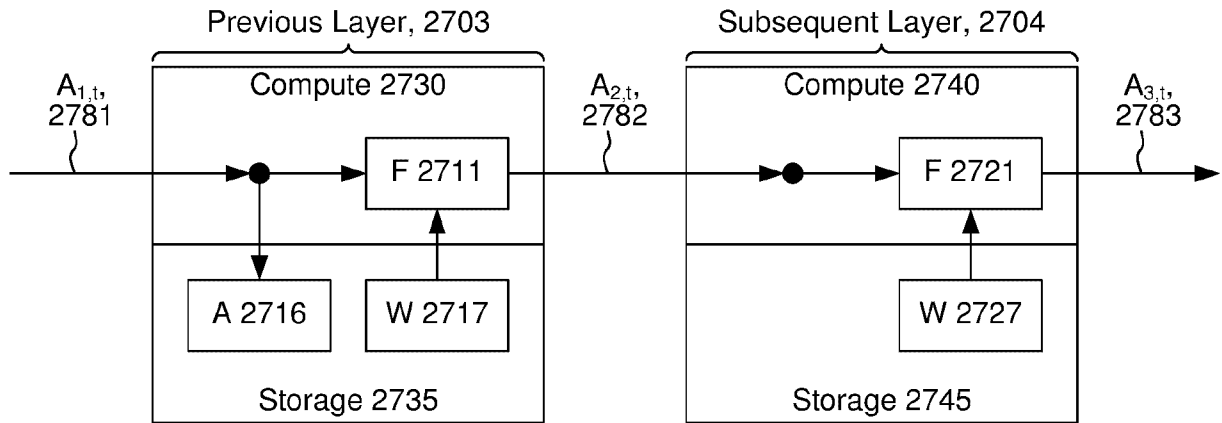


Fig. 27C

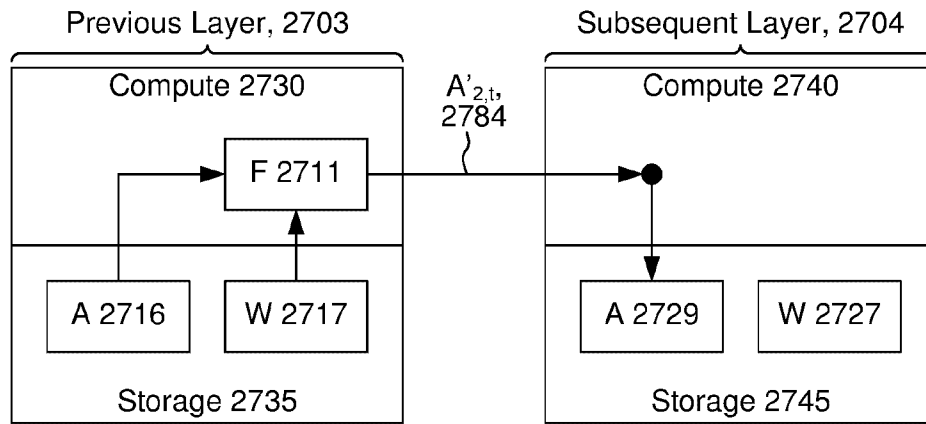


Fig. 27D

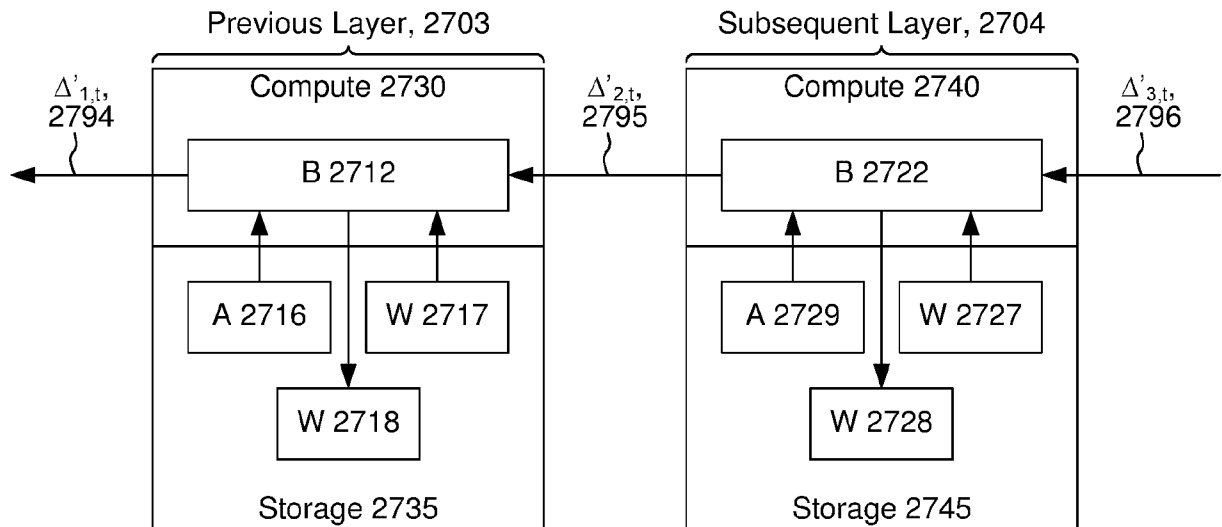


Fig. 27E

Basic Operation

matrix m x vector v = vector v'

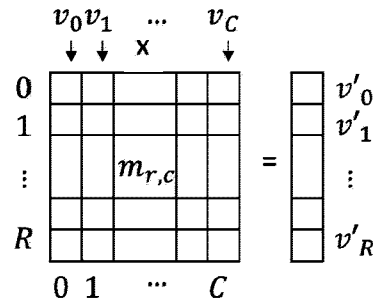


Fig. 28A

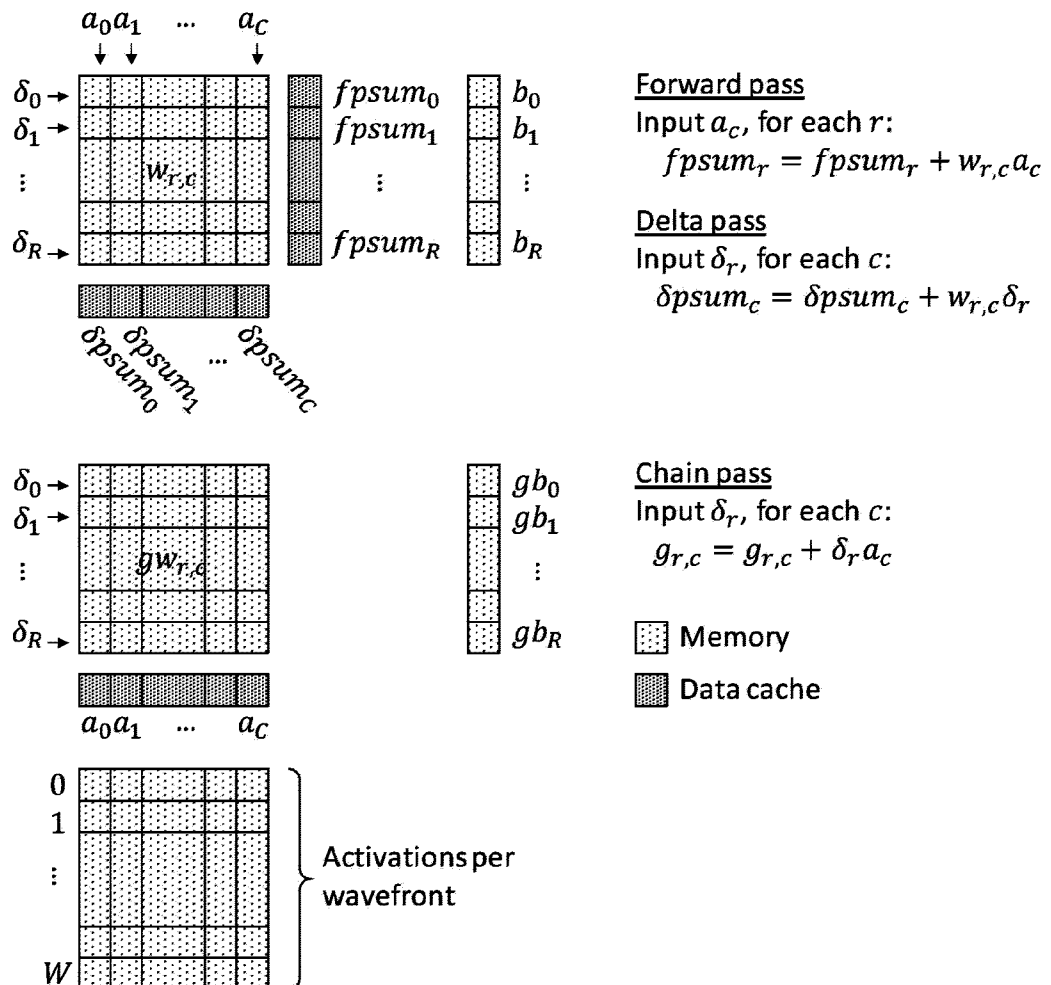


Fig. 28B

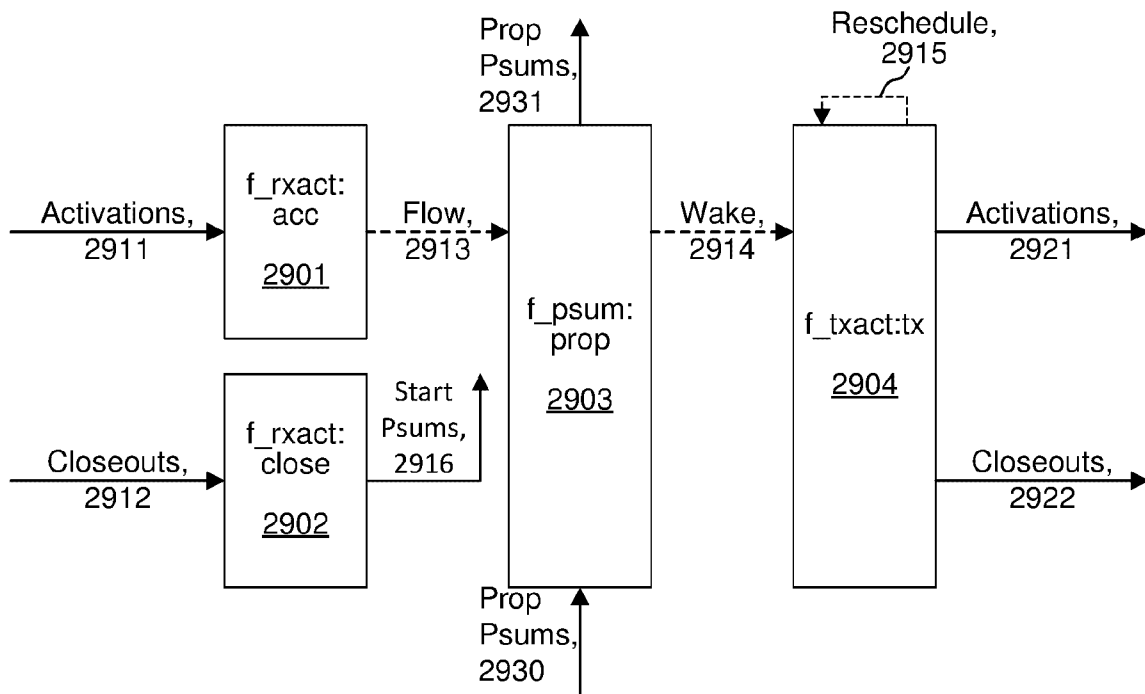


Fig. 29

Deep Learning Accelerator, 400

Forward, 401

Delta, 402

