



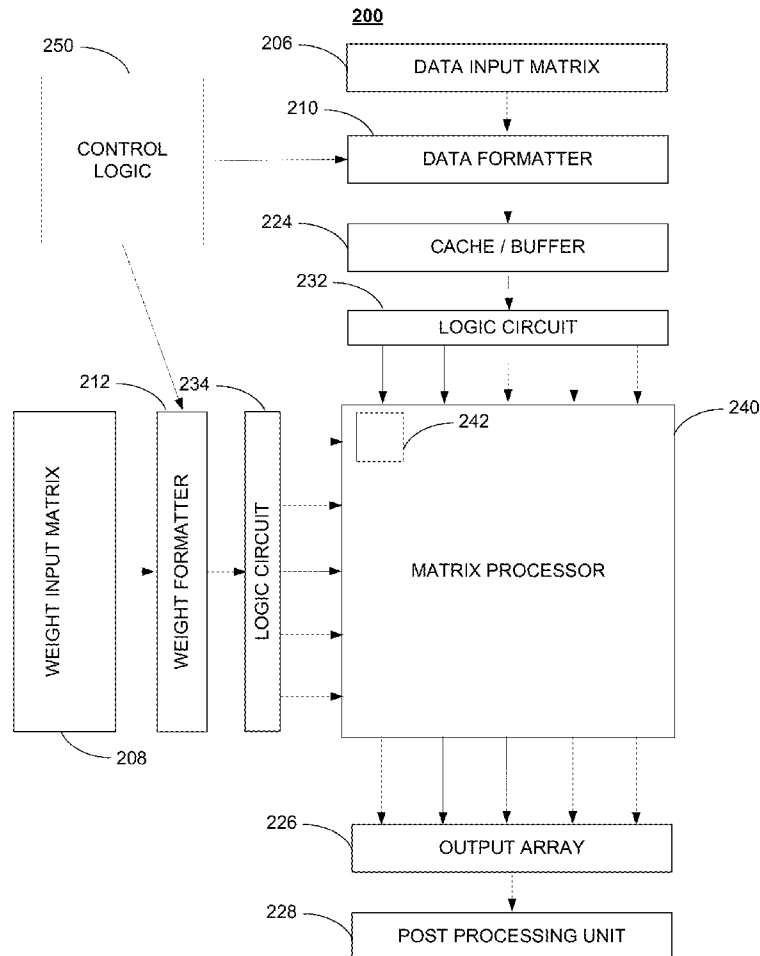
US 20190026078A1

(19) **United States**(12) **Patent Application Publication**
BANNON et al.(10) **Pub. No.: US 2019/0026078 A1**(43) **Pub. Date: Jan. 24, 2019**(54) **ACCELERATED MATHEMATICAL ENGINE****G06F 7/50** (2006.01)**G06F 7/52** (2006.01)**G06F 15/80** (2006.01)(71) Applicant: **Tesla, Inc.**, Palo Alto, CA (US)(72) Inventors: **Peter Joseph BANNON**, Woodside, CA (US); **Kevin Altair HURD**, Redwood City, CA (US); **Emil TALPES**, San Mateo, CA (US)(52) **U.S. Cl.**CPC **G06F 7/575** (2013.01); **G06T 1/20** (2013.01); **G06F 15/80** (2013.01); **G06F 7/52** (2013.01); **G06F 7/50** (2013.01)(73) Assignee: **Tesla, Inc.**, Palo Alto, CA (US)(21) Appl. No.: **15/710,433**(22) Filed: **Sep. 20, 2017****Related U.S. Application Data**

(60) Provisional application No. 62/536,399, filed on Jul. 24, 2017.

Publication Classification(51) **Int. Cl.****G06F 7/575** (2006.01)**G06T 1/20** (2006.01)(57) **ABSTRACT**

Various embodiments of the disclosure relate to an accelerated mathematical engine. In certain embodiments, the accelerated mathematical engine is applied to image processing such that convolution of an image is accelerated by using a two-dimensional matrix processor comprising sub-circuits that include an ALU, output register and shadow register. This architecture supports a clocked, two-dimensional architecture in which image data and weights are multiplied in a synchronized manner to allow a large number of mathematical operations to be performed in parallel.



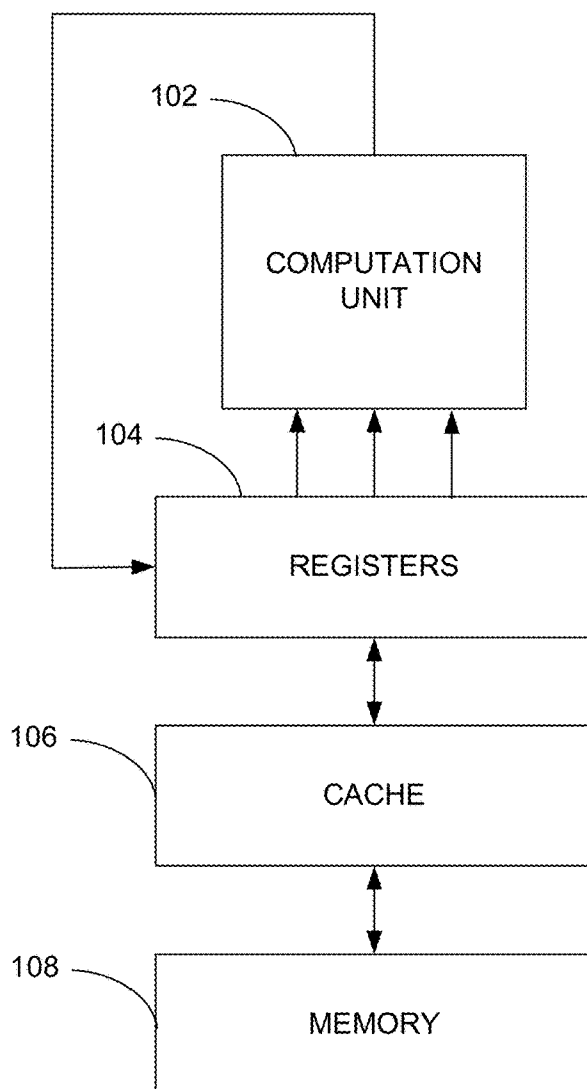
100

FIGURE 1
(PRIOR ART)

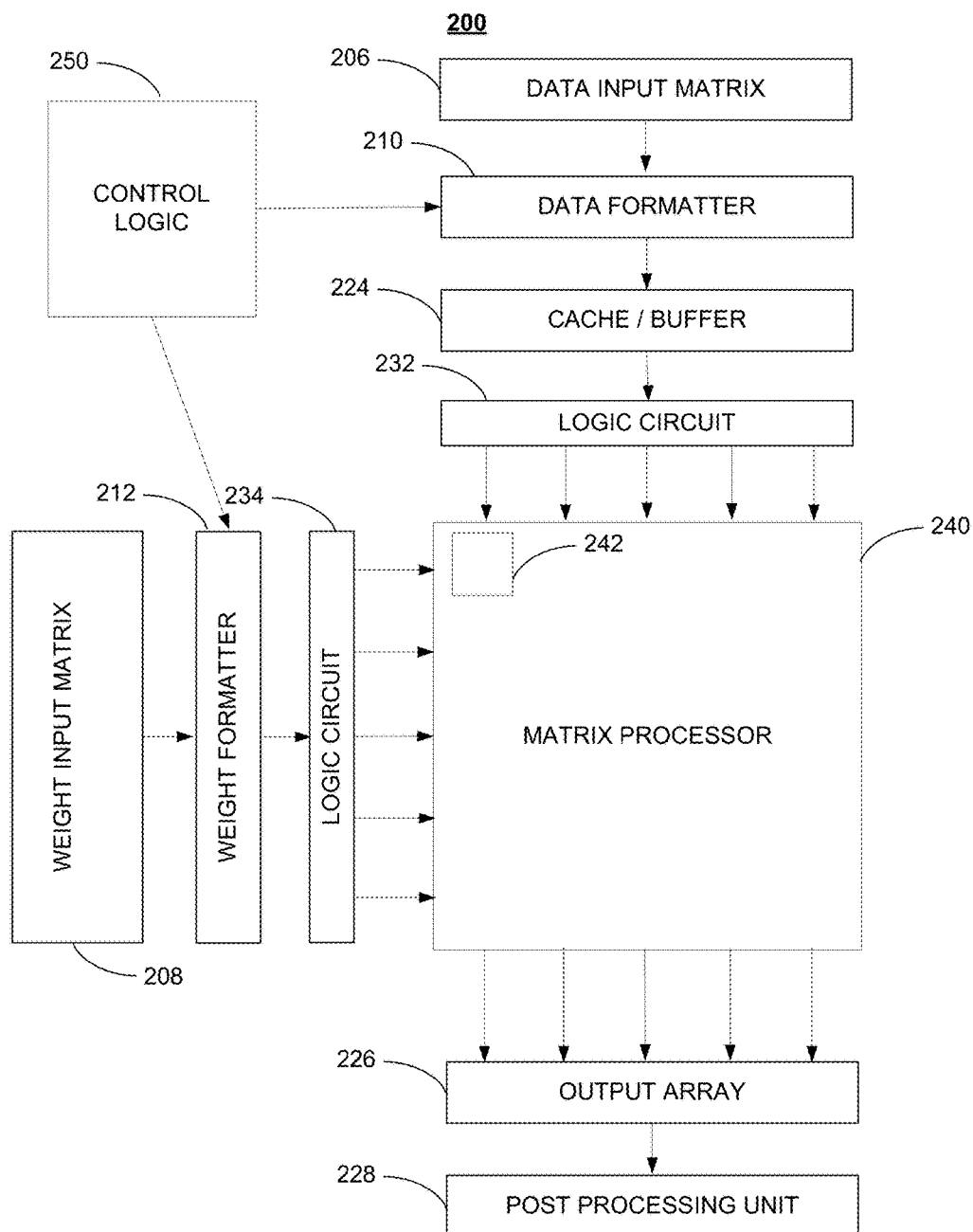


FIGURE 2

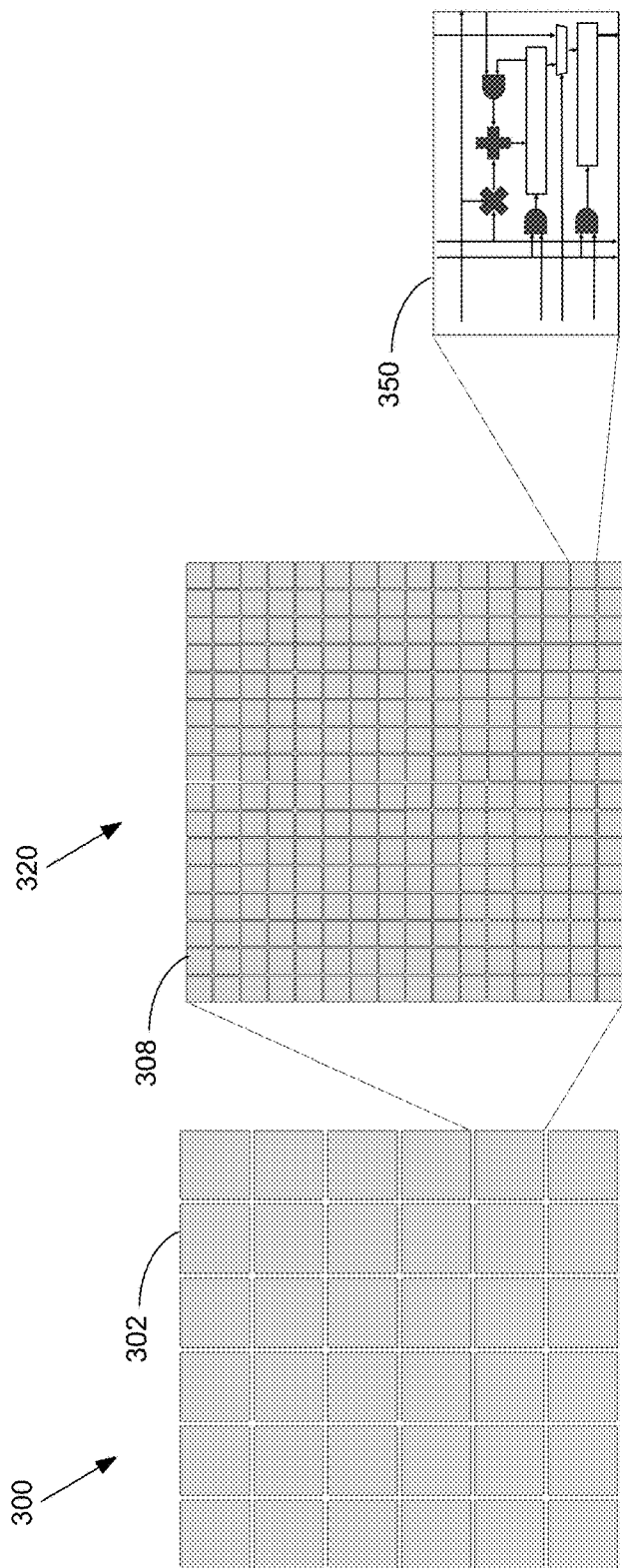


FIGURE 3

400

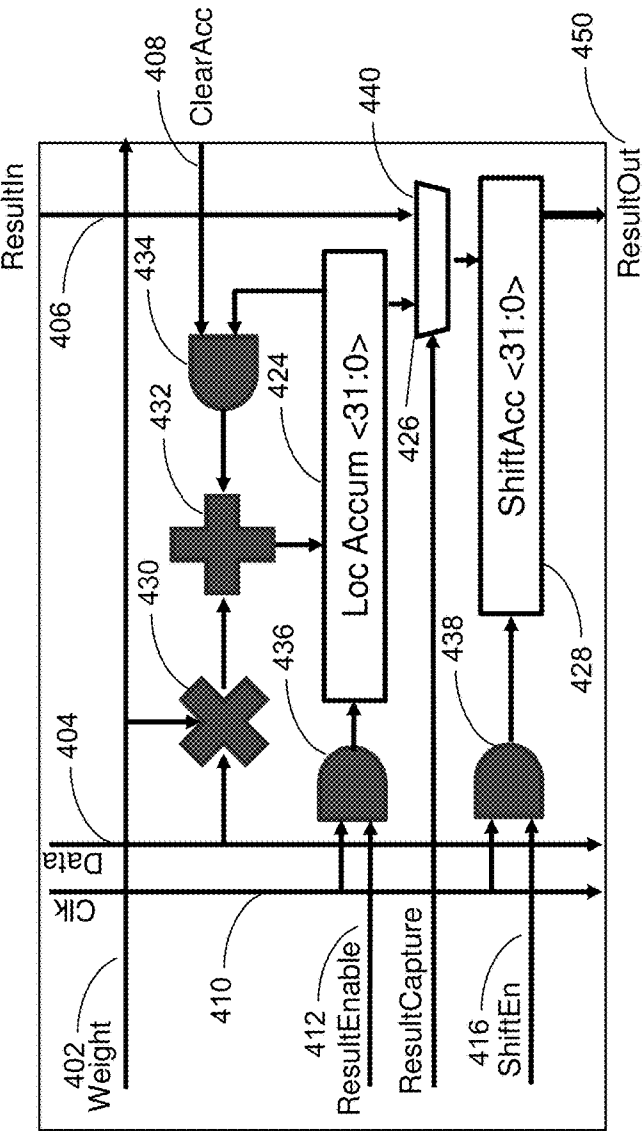


FIGURE 4

500

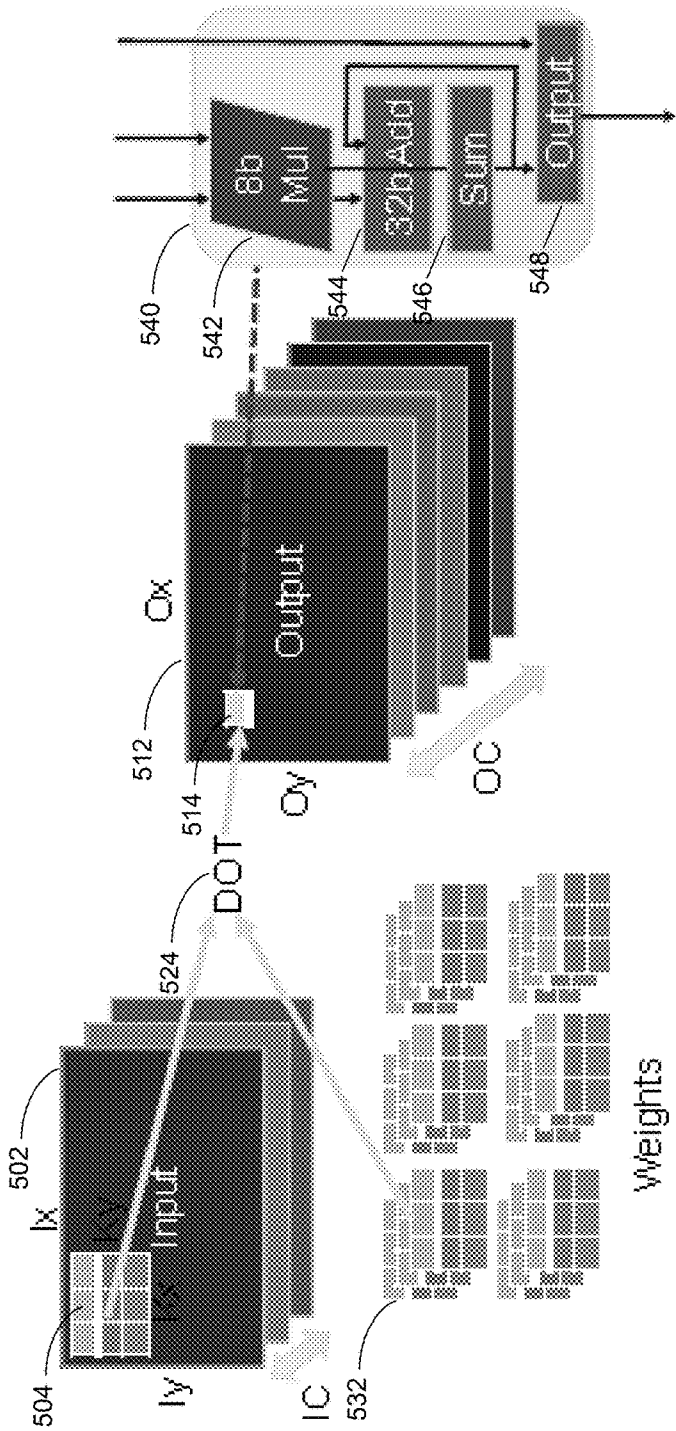


FIGURE 5

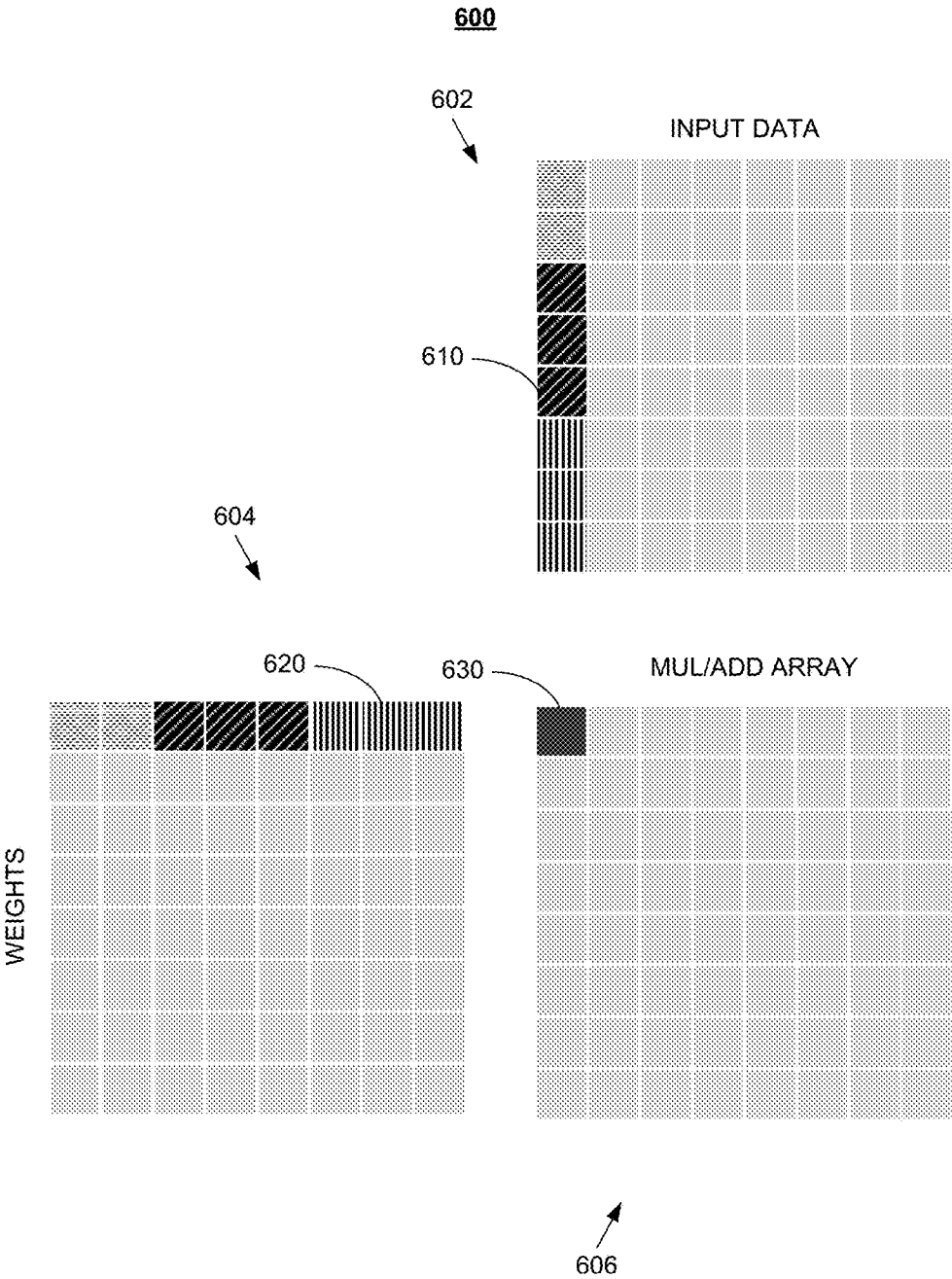


FIGURE 6

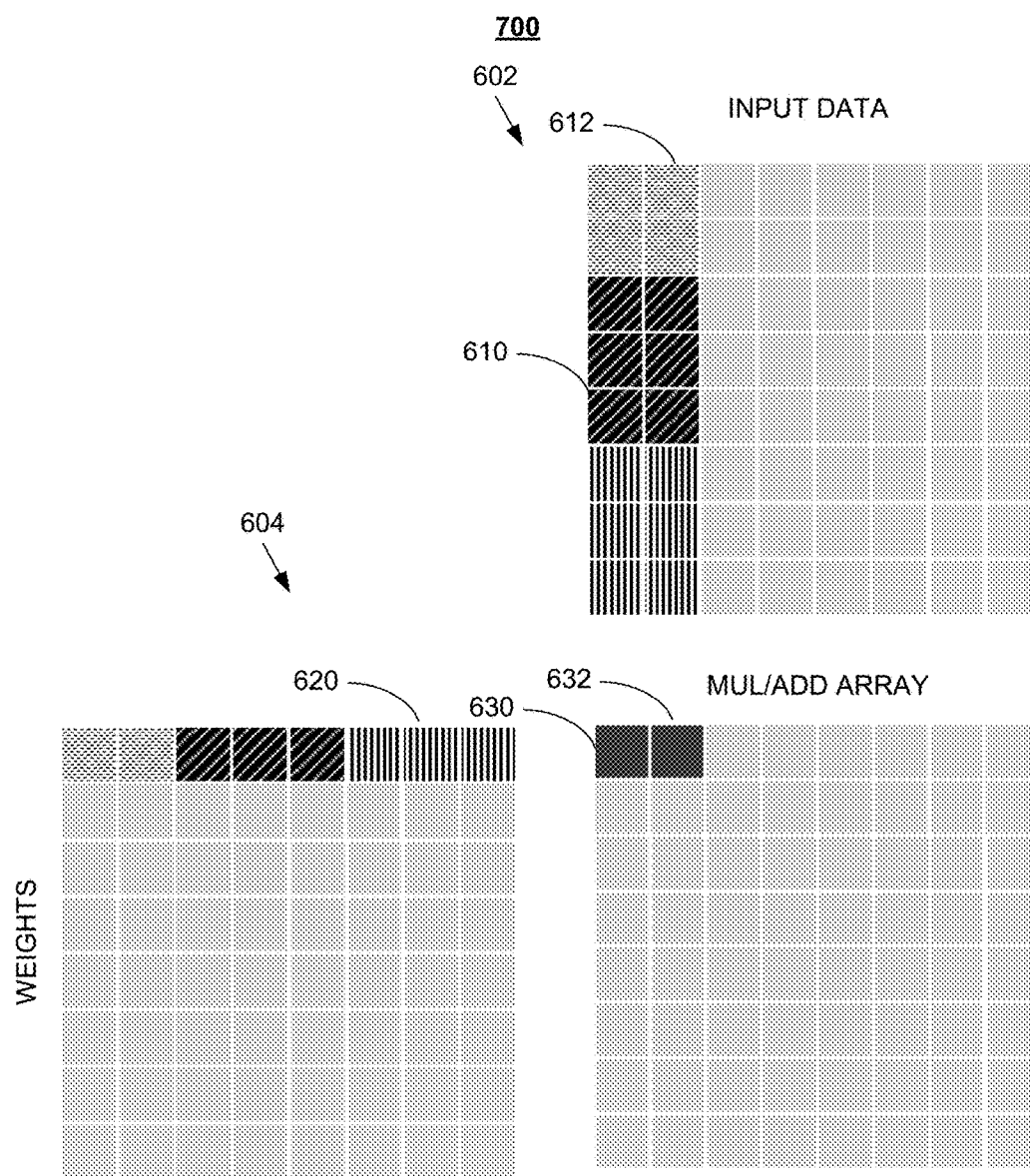


FIGURE 7

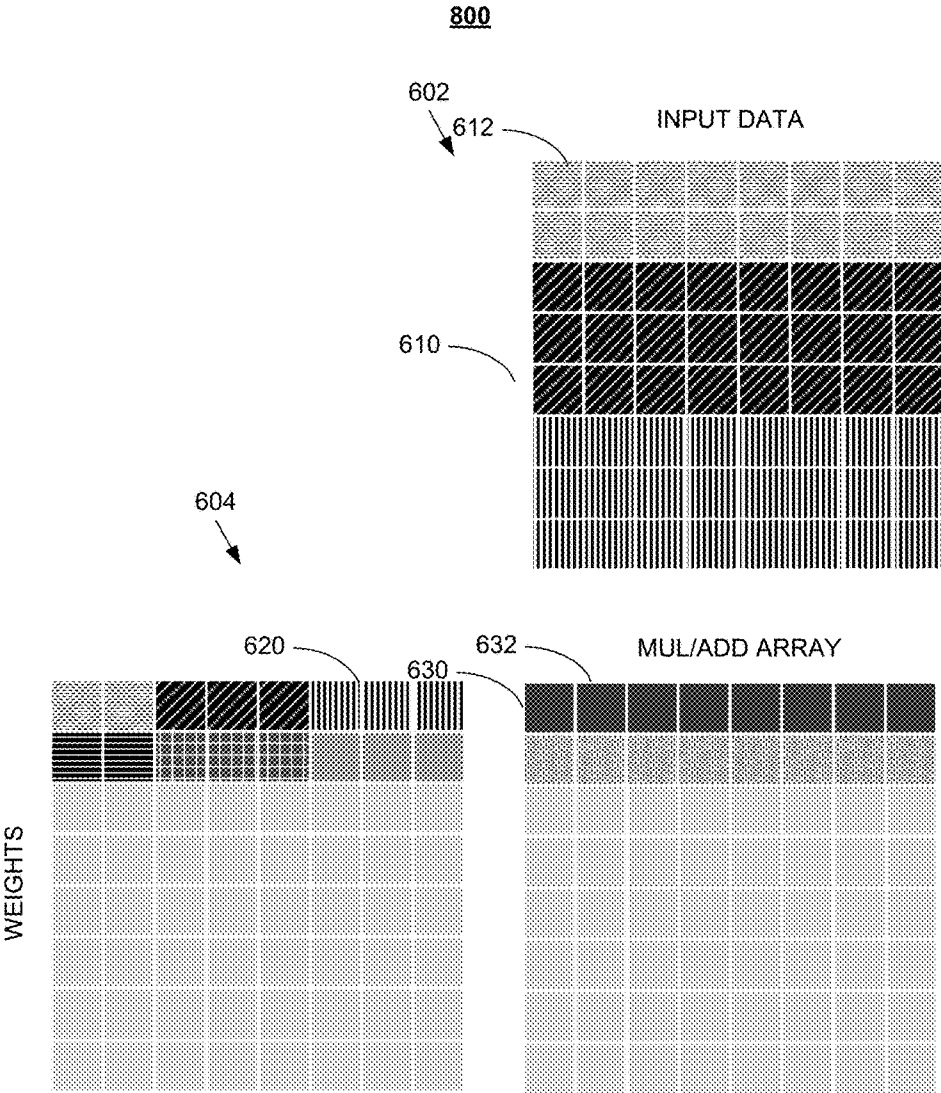


FIGURE 8

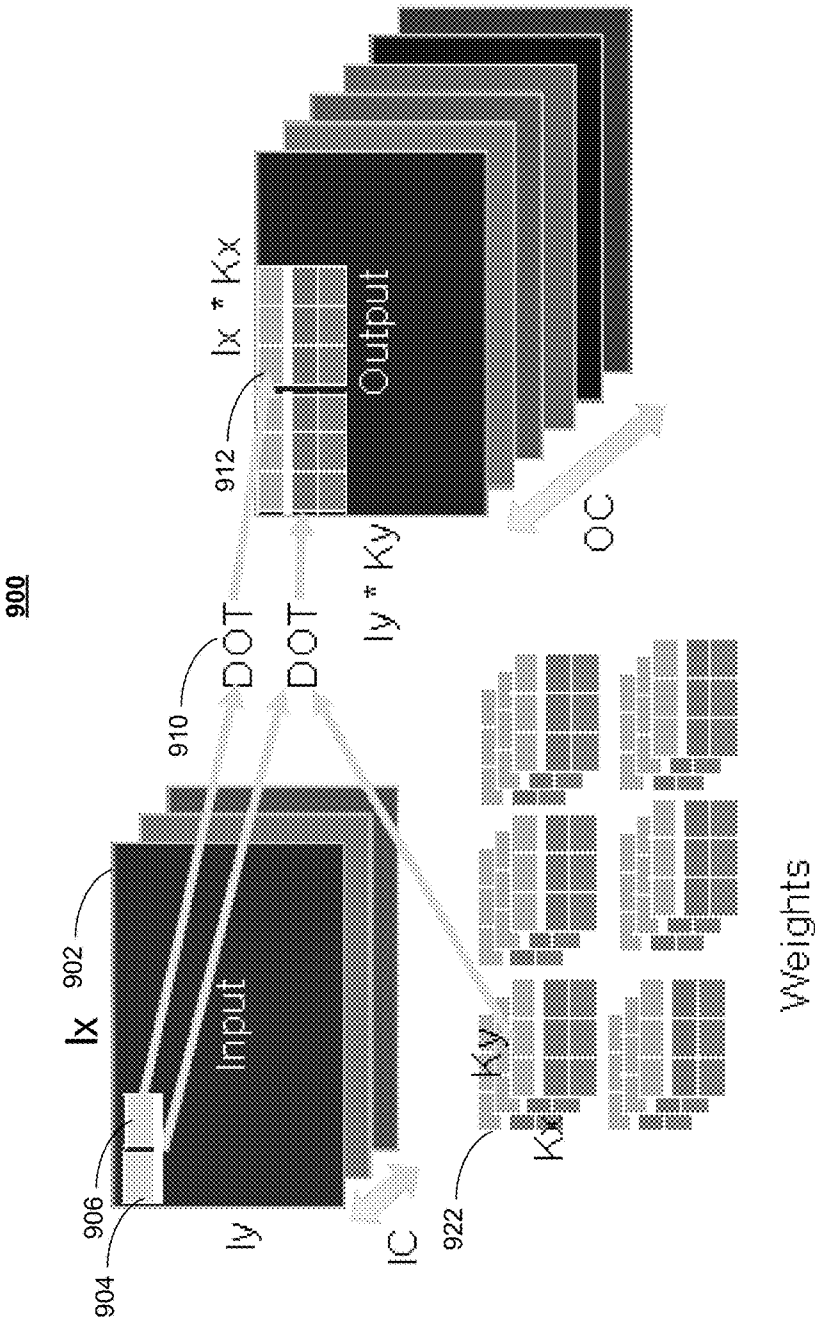


FIGURE 9

1000

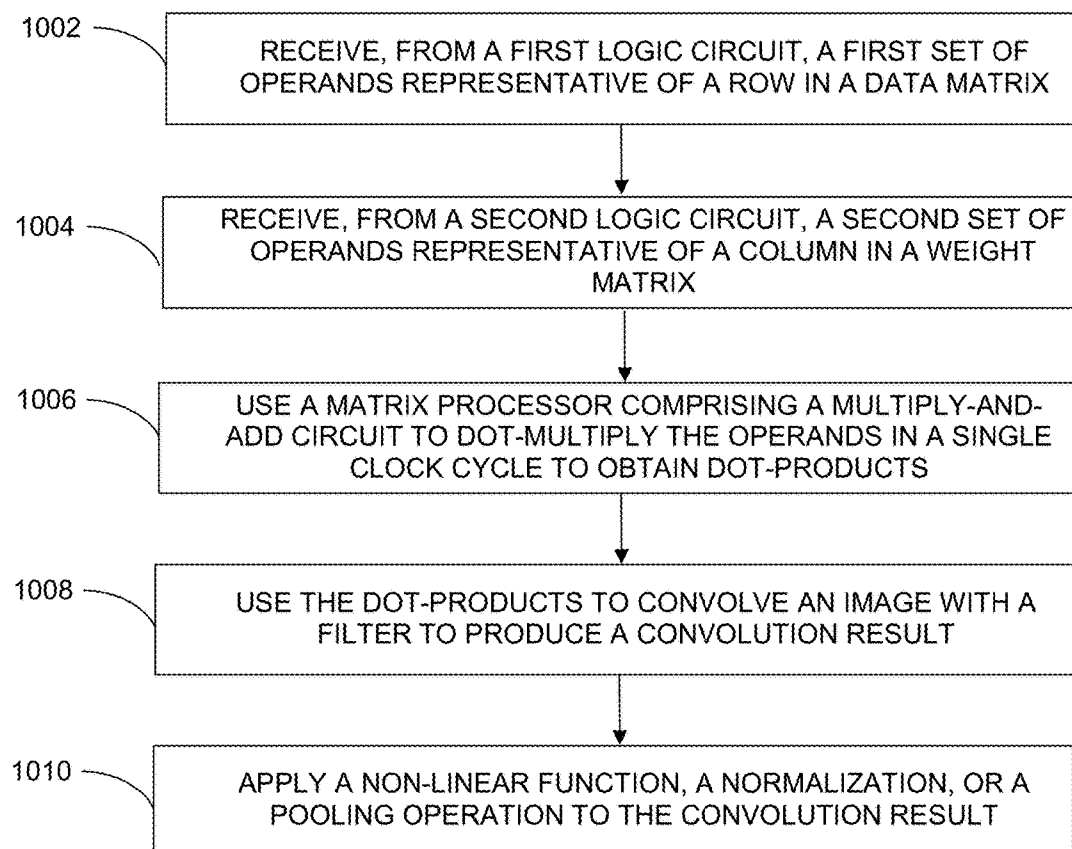


FIGURE 10

ACCELERATED MATHEMATICAL ENGINE

CROSS-REFERENCE TO RELATED APPLICATION

[0001] This application claims the priority benefit under 35 USC § 119(e) to U.S. Prov. Pat. App. Ser. No. 62/536,399 (20150-2154P (P0822-1PUS)), filed on Jul. 24, 2017, entitled “Accelerated Mathematical Engine,” and listing Peter Joseph Bannon, Kevin Altair Hurd, and Emil Talpes as inventors. The aforementioned patent document is incorporated by reference herein in its entirety and for all purposes.

BACKGROUND

A. Technical Field

[0002] The present disclosure relates to an accelerated mathematical engine for operating on large amounts of data, and more particularly, to an accelerated mathematical engine for performing complex convolution operations based on matrix multiply operations.

B. Description of the Related Art

[0003] One skilled in the art will recognize the ever-increasing demands of speed and performance on general processors and systems that are used to implement time-sensitive and complex mathematical operations. As these general systems are used to process large amounts of data and perform complex mathematical operations, the computational resources and the rate of calculations are limited by the capabilities of existing general hardware designs that perform those calculations. For example, general-purpose computing devices and processors that execute matrix operations may be unable to perform these operations in a timely manner under certain circumstances. Many conventional multipliers that perform digital signal processing operations rely on a series of software and hardware matrix manipulation steps (address generation, transpositions, bit-by-bit addition and shifting, etc.) and may represent a bottleneck within a time-sensitive system. Oftentimes, these manipulation steps require the use of a processor’s arithmetic functions to generate intermediate results at the expense of wasting computing time due to the added steps of storing and fetching intermediate results from various locations to complete an operation.

[0004] FIG. 1 shows an example of a conventional multiplier system. Multiplier system 100 is a scalar machine that comprises computation unit 102, registers 104, cache 106, and memory 108. In operation, computation unit 102 uses registers 104 and cache 106 to retrieve data stored in memory 108. Typically, computation unit 102 is a micro-processor, such as a CPU or GPU, capable of performing various computational procedures including matrix multiplication on input matrices to obtain a resultant matrix, e.g., by converting multiplications into additions and outputting the result into some internal register.

[0005] For example, a dot product that represents an output pixel of an image is typically generated by dot-multiplying individual matrix elements from two matrices to obtain partial results, which are then added to obtain the final dot product. A multiplication of individual matrix elements, i.e., a scalar multiplication, is typically performed on individual data elements by breaking up the dot multiplication into a series of individual sub-operations. As a result, partial

products have to be stored and fetched from one or more of registers 104, cache 106, and memory 108 to complete a single arithmetic operation.

[0006] Computationally demanding applications, such as a convolution, oftentimes require a software function be embedded in computation unit 102 and used to convert convolution operations into alternate matrix-multiply operations. This is accomplished by rearranging and reformatting data into two matrices that then can be raw matrix-multiplied. However, there exists no mechanism to efficiently share or reuse data in scalar machine 100, such that data necessary to execute each scalar operation has to be re-stored and re-fetched from registers many times. The complexity and managerial overhead of these operations becomes significantly greater as the amount of image data subject to convolution operations increases.

[0007] The inability to reuse much of the data in scalar machine 100 coupled with the added and inefficient steps of storing and fetching intermediate results from registers 104, cache 106, and memory 108 to complete an arithmetic operation are only some of the shortcoming of existing systems, such as multiplier system 100.

[0008] Accordingly, what is needed are high-computational-throughput systems and methods that can perform matrix mathematical operations quickly and efficiently.

BRIEF DESCRIPTION OF THE DRAWINGS

[0009] References will be made to embodiments of the invention, examples of which may be illustrated in the accompanying figures. These figures are intended to be illustrative, not limiting. Although the invention is generally described in the context of these embodiments, it should be understood that it is not intended to limit the scope of the invention to these particular embodiments. Items in the figures may be not to scale.

[0010] FIG. 1 shows an example of a conventional multiplier system.

[0011] FIG. 2 illustrates and exemplary matrix processor architecture for performing arithmetic operations according to various embodiments of the present disclosure.

[0012] FIG. 3 illustrates details of an exemplary configuration of the matrix processor architecture shown in FIG. 2.

[0013] FIG. 4 illustrates an exemplary multiply-and-add circuit implementation of the logic circuit shown in FIG. 3.

[0014] FIG. 5 illustrates an exemplary convolution operation according to various embodiments of the present disclosure.

[0015] FIG. 6 through FIG. 8 illustrate details of an exemplary convolution operation according to various embodiments of the present disclosure.

[0016] FIG. 9 illustrates an exemplary deconvolution operation according to various embodiments of the present disclosure.

[0017] FIG. 10 illustrates a process for performing arithmetic operations to make convolutional neural networks faster, according to various embodiments of the present disclosure.

DETAILED DESCRIPTION OF EMBODIMENTS

[0018] In the following description, for purposes of explanation, specific details are set forth in order to provide an understanding of the invention. It will be apparent, however, to one skilled in the art that the invention can be practiced

without these details. Furthermore, one skilled in the art will recognize that embodiments of the present invention, described below, may be implemented in a variety of ways, such as a process, an apparatus, a system, a device, or a method on a tangible computer-readable medium.

[0019] Components, or modules, shown in diagrams are illustrative of exemplary embodiments of the invention and are meant to avoid obscuring the invention. It shall also be understood that throughout this discussion that components may be described as separate functional units, which may comprise sub-units, but those skilled in the art will recognize that various components, or portions thereof, may be divided into separate components or may be integrated together, including integrated within a single system or component. It should be noted that functions or operations discussed herein may be implemented as components. Components may be implemented in software, hardware, or a combination thereof. Many components are formed through interconnection of many subcomponents. Subcomponents may be selected that are logically different in operation from what is shown herein, where these logically different subcomponents can be combined in the aggregate with other subcomponents provide similar or identical functionality at the aggregated component level to that described herein (e.g., active high signals can be active low, AND gates replaced with inverted-input NOR gates, etc).

[0020] Furthermore, connections between components or systems within the figures are not intended to be limited to direct connections. Rather, data between these components may be modified, re-formatted, or otherwise changed by intermediary components. Also, additional or fewer connections may be used. It shall also be noted that the terms “coupled,” “connected,” or “communicatively coupled” shall be understood to include direct connections, indirect connections through one or more intermediary devices, and wireless connections.

[0021] Reference in the specification to “one embodiment,” “preferred embodiment,” “an embodiment,” or “embodiments” means that a particular feature, structure, characteristic, or function described in connection with the embodiment is included in at least one embodiment of the invention and may be in more than one embodiment. Also, the appearances of the above-noted phrases in various places in the specification are not necessarily all referring to the same embodiment or embodiments.

[0022] The use of certain terms in various places in the specification is for illustration and should not be construed as limiting. A service, function, or resource is not limited to a single service, function, or resource; usage of these terms may refer to a grouping of related services, functions, or resources, which may be distributed or aggregated.

[0023] The terms “include,” “including,” “comprise,” and “comprising” shall be understood to be open terms and any lists that follow are examples and not meant to be limited to the listed items and may include subsets or supersets of the items along with additional items. Any headings used herein are for organizational purposes only and shall not be used to limit the scope of the description or any claims. Each document mentioned in this patent document is incorporated by reference herein in its entirety.

[0024] Furthermore, one skilled in the art shall recognize that: (1) certain steps may optionally be performed; (2) steps may not be limited to the specific order set forth herein; (3)

certain steps may be performed in different orders; and (4) certain steps may be done concurrently.

[0025] Although embodiments herein are discussed mainly in the context of convolutions, one of skill in the art will appreciate that a deconvolution and other matrix operations can also be structured as a matrix-matrix type multiply operation and, thus, the principles of the present invention are equally applicable to deconvolutions. Furthermore, other types of mathematical operations may be implemented in accordance with various embodiments of this disclosure.

[0026] FIG. 2 illustrates an exemplary matrix processor architecture for performing arithmetic operations according to various embodiments of the present disclosure. System 200 comprises logic circuit 232 234, cache/buffer 224, data formatter 210, weight formatter 212, data input matrix 206, weight input matrix 208, matrix processor 240, output array 226, post processing units 228, and control logic 250. Matrix processor 240 comprises a plurality of sub-circuits 242 which contain Arithmetic Logic Units (ALUs), registers and, in some embodiments, encoders (such as booth encoders). Logic circuit 232 may be a circuit that represents N input operators and data registers. Logic circuit 234 may be circuitry that inputs M weight operands into matrix processor 240. Logic circuit 232 may be circuitry that input image data operands into matrix processor 240. Weight input matrix 208 and data input matrix 206 may be stored in various types of memory including SRAM devices. One skilled in the art will recognize that various types of operands may be input into the matrix processor 240.

[0027] In operation according to certain embodiments, system 200 accelerates convolution operations by reducing redundant operations within the systems and implementing hardware specific logic to perform certain mathematical operations across a large set of data and weights. This acceleration is a direct result of methods (and corresponding hardware components) that retrieve and input image data and weights to the matrix processor 240 as well as timing mathematical operations within the matrix processor 240 on a large scale.

[0028] In embodiments, formatters 210 212, which in example in FIG. 2 are implemented as in-line formatters. In certain embodiments, formatters 210 212 are discrete components and in other embodiments the formatters 210 212 are integrated together and/or with one or more other components. Each is implemented in hardware and converts a matrix to a vector on operands to be operated upon within the matrix processor 240. In other embodiments, formatters 210 212 are implemented in software, although this typically produces a loss in speed. Data formatter 210 converts two-dimensional or three-dimensional (e.g., a 3×3×3 cube) data comprising data input matrix 206 into a single vector or string that may be represented by a row or column, thereby, linearizing or vectorizing data input matrix 206. In detail, formatter 210 receives data input matrix 206 and prepares input data to be processed by matrix processor 240. In embodiments, this is accomplished by mapping parameters of the data input matrix 206 into a suitable format according to the hardware requirements of matrix processor 240 such that matrix processor 240 can efficiently perform a matrix multiply as part of a convolution calculation when generating output pixels.

[0029] As an example, assuming matrix processor 240 comprises 96 rows and 96 columns, data mapped into a 96×96 format would cause matrix processor 240 to be

utilized to its full computational capacity and, thus, provide a preferred efficiency. In that case, formatter **210** should produce an output that is 96-columns wide. Similarly, formatter **212** should produce an output that is 96-rows wide based on the weight input matrix **208**.

[0030] In embodiments, formatter **210** uses a number of multiplexers or switches to fetch some or all of data input matrix **206** and choose different elements therefrom in order to produce data that is then lined up according to the columns of matrix processor **240**. In embodiments, the selection ensures that the appropriate data from data input matrix **206** is passed to each of the columns at defined clock cycles. In embodiments, if weights are static, they may be pre-formatted offline, stored in memory, fetched only once, and fed directly into matrix processor **240** in a modified, vectorized format without the use of formatter **212**. In other embodiments, weights may be dynamically adjusted and fed into matrix processor **240** in accordance with various formatting and fetching operations. In embodiments, matrix processor **240** allows for column and row inputs of varying sizes. That is, matrix processor **240** is designed to compute $N \times M$ computations of arbitrary size.

[0031] In other embodiments, if the number of columns of the matrix processor **240** is limited (for example to N columns) such that the number of columns in the data input matrix **206** (for example X) is greater than the number of columns of the matrix processor **240** (i.e., $X > N$), then the control logic **250** may split the data input matrix **206** into multiple submatrices with each submatrix computed by a matrix processor **240**. In such instances, each matrix processor **240** may be running in a different thread. For example, if data input matrix **206** consists of 192×96 data points, and the matrix processor has 96 columns and 96 rows (i.e., 96×96 computations may occur in one clock cycle), the control logic **250** may split the data input matrix **206** into two submatrices (such as the left half of the data input matrix **206** and the right half of the data input matrix **206**). Each submatrix will consist of 96×96 data points. Each separately threaded matrix processor **240** can compute the output channels for the submatrix sent to it with results placed into the final output array **260**, which must be large enough to hold the values from all channels (that is 192 values). More generally, data input matrix **206** may be split into any number of submatrices and sent to different matrix processors **240**, each running in a separate thread. As with the output array **226**, the data input matrix **206**, data formatter **210**, cache/buffer **224**, logic circuit **232**, and post processing unit **228** must similarly be able to accommodate the larger data.

[0032] In alternative embodiments, a CNN may be computed between multiple matrix processors **240** by having control logic **250** splitting the computations along the inner product. The segments of the inner product are computed, each in a different matrix processor **240**, and then the input products added together to compute the output vector, which is then stored in output array **260**.

[0033] Unlike common software implementations of formatting functions that are performed by a CPU or GPU to convert a convolution operation into a matrix-multiply by rearranging data to an alternate format that is suitable for a fast matrix multiplication, various hardware implementations of the present disclosure re-format data on the fly and make it available for execution, e.g., 96 pieces of data every cycle, in effect, allowing a very large number of elements of

a matrix to be processed in parallel, thus efficiently mapping data to a matrix operation. In embodiments, for $2N$ fetched input data $2N^2$ compute data may be obtained in a single clock cycle. This architecture results in a meaningful improvement in processing speeds by effectively reducing the number of read or fetch operations employed in a typical processor architecture as well as providing a paralleled, efficient and synchronized process in performing a large number of mathematical operations across a plurality of data inputs.

[0034] In embodiments, to increase efficiency of matrix processor **240** that may have any arbitrary number of columns and rows, formatter **212** **214** may reformat different shapes of input matrices data into the columns and rows suitable for matrix processor **240**. In embodiments, formatting is performed dynamically to accommodate processing of matrices having different input sizes. In embodiments, the reformatted matrixes comprising input channels are fed into cache/buffer **224**.

[0035] Cache/Buffer **224** may fetch data from data input matrix **206** only $1/k$ times as various pieces of data may be reused, where k is the convolution kernel width. For example, for any given cycle, once a row is fetched, certain columns will have access to all the data in that row. In embodiments, cache/buffer **224** may be a local buffer that stores a local copy of data that may be reused by a convolution without having to re-access and read data from SRAM.

[0036] Once matrix processor **240** has completed a computation, a set of result may be shifted, e.g., from the accumulators in the bottom row of matrix processor **240**, e.g., to output flip-flops (not shown) that effectively form a shift register that receive a dot product. In embodiments, pulling or shifting results into output array **226**, e.g., one per clock cycle, from a row that corresponds to an output channel may be accomplished by a state machine (not shown). The state machine may perform additional operations on the output channel, for example, prior to sending data to SRAM and/or post processing unit **228**. The internal operation of matrix processor **240** will be described in more detail below.

[0037] In embodiments, matrix processor **240** comprises shadow registers that enable parallel processing by storing a copy of the results that are passed through matrix processor **240** to output array **226**. In embodiments, moving an operation result from output register to shadow register involves loading the next set of values into the ALUs.

[0038] Once an accumulation has completed, a convolution may commence and accumulation may start over before all of the data of a prior convolution is output to output array **226**. As a result, in every clock cycle, the data in matrix processor **240** may move down by one row, such that for each cycle the last row may be output to output array **226**. In effect, this mode of operation ensures that a new calculation may be made in each consecutive cycle without any interruptions and independent of additional processing operations, such as storing data in SRAM, etc.

[0039] Post processing unit **228** may comprise or interact with a number of devices (not shown), such as a hardware-accelerated pooling unit, a DRAM that may be part of a direct memory access ("DMA") that retrieves data from memory and stores data (e.g., weights and results) in SRAM, and the like. The devices may be partially or entirely

controlled by control logic **250**, which may also manage formatters **210** **212** and other components within system **200**.

[0040] Not shown in FIG. 2 are auxiliary devices that perform management functions, such as a sequencer that generates addresses for reading the data, writes the results, and keeps track of where system **200** is in the convolution in order to calculate from where to get and how to execute the data that will be used in a subsequent step of the convolution.

[0041] In certain embodiments, weight input matrix **208** is physically split and drives weights from two different sides of matrix processor **240**, such that the two-dimensional array is split into two regions (e.g., a left-hand side and a right-hand side) that each receive a portion of the data in weight input matrix **208**. Such an implementation reduces data latency by taking advantage of the fact that weights are known. In embodiments, in order to reduce peak power consumption, the timing of operations may be chosen such that multiplications of weight and data are spread out over a certain number of cycles. This efficient timing of operations results in a reduction of energy consuming steps including a decrease in the number of read operations performed by the matrix processor and improving the efficiency of data movement within the matrix (e.g., between sub-circuits).

[0042] In embodiments, a state machine (not shown) that is configured to identify redundant data may be employed. Identified redundant data may be reused across columns, such that the data does not need to be re-fetched. The state machine may be configured to determine how and where to shift data that is to be executed, e.g., based on inputs related to image size, filter size, stride, number of channels, and similar parameters.

[0043] In embodiments, a booth encoder is shared across a number of elements in the multiplication architecture of matrix processor **240**. The booth encoder may be any booth encoder known in the art and may be used to multiply two numbers and encode one of the two numbers, e.g., from an 8-bit value to a 12-bit or any other value that makes multiplication operations easier on the multiplier logic and, thus, faster. In embodiments, the booth encoder may be applied in parallel across an entire row so as to share the same encoded, alternate weight value across all columns. By loading an operand across all columns, a multiplication may be performed in a single clock cycle across an entire row. The cost for leveraging re-encoding to share the same data (e.g., weights) across for N computational elements is thus paid only once for each column (or row). In comparison, in existing computing architectures, every single scalar would require a booth encoder for every single multiplication operation.

[0044] FIG. 3 illustrates details of an exemplary configuration of the matrix processor architecture shown in FIG. 2. In embodiments, matrix processor **300** may accommodate a predetermined vector length on each axis. As depicted in FIG. 3, matrix processor **300** may comprise an array of 6×6 tiles **302** that are arranged in a matrix format. Each tile **302** may comprise a matrix **320** that, in turn, comprises sub-circuits **350**. As discussed in detail below with reference to FIG. 4, each sub-circuit circuit **350** may be a cell capable of performing arithmetic operations. In embodiments, sub-circuit circuit **350** performs simultaneously multiplication, accumulation, and shift operations.

[0045] In embodiments, arithmetic operations are parallelized by utilizing multiple rows and columns of matrix processor **300** to generate an NxN tile output. For example, a given row size of 96 and a corresponding column size of 96 facilitate an output of 2*9216 mathematical calculations. In other embodiments, the number of rows and columns may be different. That is, there may be N rows and M columns and an NxM tile output may be generated. For example, for a row size of 96 and a corresponding column size of 192, an output of 2*18,432 calculations is generated in a single clock cycle.

[0046] FIG. 4 illustrates an exemplary multiply-and-add circuit implementation of the sub-circuit shown in FIG. 3. As depicted in FIG. 4, multiply-and-add circuit **400** comprises multiplier **430**, adder **432**, logic **434** **436** **438**, accumulator **424**, shadow register **428**, and output register **440**. In embodiments, accumulator **424** may be implemented as an accumulation register.

[0047] In embodiments, accumulator **424** may comprise a set of ALUs that comprise registers and shadow register **428** that may be configured to receive the outputs of the ALUs.

[0048] In operation, multiplier **430** receives and multiplies weights **402** and data **404** to generate products therefrom. Each product may be provided to adder **432** that, in response to receiving the product from multiplier **430**, adds the product to the current value of the accumulator **424**.

[0049] In embodiments, accumulator **424** generates an accumulated value that is stored, e.g., in output register **440**. The accumulated value is the result of a convolution and, as mentioned with reference to FIG. 2, may correspond to the dot product of two formatted matrices.

[0050] In embodiments, a copy of the result in output register **440** may be provided to shadow register **428**, which may output result **450**, such that accumulator **424** can be accessed again to commence new calculations. In embodiments, multiply-and-add circuit **400** in FIG. 4 may perform a multiplication, an addition operation, and a shift operation at the same time, i.e., within a single cycle, thereby doubling the total number of operations that occur each cycle.

[0051] In embodiments, ClearAcc signal **408** clears the contents of accumulator **424**, e.g., when multiplier **430** performs a multiply operation, such that accumulation operations can start over. In embodiments, ResultEnable signal **412** is activated in response to a determination that data **404** is valid. It is understood that accumulator **424** may accumulate and save data, accumulate and clear data, or just clear data.

[0052] In embodiments, results are moved from output register **440** to shadow register **428** in a single clock cycle, i.e., without the need of intermediate execute and save operations.

[0053] FIG. 5 illustrates an exemplary convolution operation according to various embodiments of the present disclosure. Convolution **500** comprises input channels IC of input image **502**, weights **532**, dot product **514**, output channels OC, and accumulator **540**.

[0054] In embodiments, convolution operation **500** applies individual filters (i.e., weights) **532** to input image **502**, e.g., to detect small features within input image **502**. By analyzing a sequence of different features in a different order, macro features may then be identified in input image **502**. In other embodiments, input **502** is non-image data. For example, input **502** may be non-image sensor data, such as ultrasonic, radar, LIDAR, or other sensor data. Input **502**

may also be general mathematical computations or any other types of data known to one of skill in the art.

[0055] Convolution 500 may use a different set of weights 532 for each input channel IC, as each input channel IC may contain a different set of information, and each weight matrix 532 may be designed to help identify a different feature. In embodiments, convolution 500 multiplies a rectangular input matrix 504 with a rectangular weight matrix 532 to obtain partial dot products. The partial dot products may then summed by adder 546 in order to generate an accumulated dot product 514 (i.e., an integer) that represents an output pixel 514 in the output image.

[0056] In embodiments, each pixel in output channel OC is generated by multiplier 542 and adder 544. In embodiments, the value of the partial dot products correspond to the application of weight matrix 532 in its entirety to area 504 of the input image 502. In other words, each weight 532 is dot multiplied by multiplier 542 with area 504 to produce a partial dot product, then the partial dot products are accumulated in accumulator 540 to generate an accumulated output that represents the convolution.

[0057] One or more input channels IC, e.g., one for each color (e.g., RGB) may be used. For example, each convolution may use weights 532 that represent three different matrices, one for each color. Each output channel OC 512 may be generated using a different filter or weight 532 that represents a different feature in input data 502. The number of output channels may depend on the number of features. The number of convolutions is equal to the number of output channels OC times the number of input channels IC, and each convolution may have N convolutions for each input channel IC. One skilled in the art will recognize that the number and type of input channels may vary and may include color and/or clear inputs.

[0058] As depicted in FIG. 5, input matrix 504 is a $K \times K \times Y$ (i.e., 3×3) matrix that may be combined with a 3×3 weight matrix 532 across 3 input channels, i.e., $3 \times 3 \times IC$, such that the depths match and produce a single element, dot product 514, in the output plane. Each dot product 514 in output channel 512 is the result of a dot multiplication.

[0059] FIG. 6 through FIG. 8 illustrate details of an exemplary convolution operation according to various embodiments of the present disclosure. Convolution 600 comprises input data matrix 602, weight data matrix 604, array 606, and dot product 630. In embodiments, array 606 is a matrix processor architecture as shown in FIG. 2 and FIG. 3.

[0060] Input data matrix 602 in FIG. 6 comprises column 610 that, in embodiments, may be obtained by linearizing an input matrix, such as rectangular input matrix 504 shown in FIG. 5, to obtain a vectorized form of the input matrix. Similarly, weight data matrix 604 comprises row 620 that may be a vectorized form of a weight matrix, such as rectangular weight matrix 532 in FIG. 5. As an example, a 3×3 input matrix and 3 input channels may be re-formatted into a vector that comprises $3 \times 3 \times 3 = 27$ elements from which a 27-element column 610 may be produced for use in input data matrix 602. Conversely, a 3×3 weight matrix for the same 3 input channels may be used to generate a 27-element row 620 for use in weight data matrix 604. One skilled in the art will recognize that the sizes of input matrices and number of input channels may vary across different applications.

[0061] In embodiments, the input channels and input weights drawn as rectangles in FIG. 5 are reformatted, e.g.,

by the formatter discussed with reference to FIG. 2, into a vector formats (e.g., vectors having 96 elements) that are provided to a matrix multiplier/processor (denoted as element 240 FIG. 2), such that a 96×96 element dot product operation can be performed in parallel. In detail, input data 504 and input weights 532 shown in FIG. 5 as rectangles for each input channel are reformatted into vector formats.

[0062] In embodiments, the resulting vector formats, illustrated in FIG. 6 as input data 602 and input weights 604 (e.g., each having comprising 96 elements) are provided to matrix processor or matrix multiplier 240 that performs a 96×96 element dot product operation in parallel. In embodiments, in the calculation of output channels, the same output pixels are produced using the same set of input data but different set of weights (i.e., filters), such that by reading the input data once many output channels can be generated at once. As stated above, it is understood that the number of input and output channels may be arbitrarily chosen.

[0063] It is further understood that input data matrix 602, weight data matrix 604, and array 606 may have different numbers of columns and rows as those depicted in FIG. 6. In particular, the shapes of input data matrix 602 and weight data matrix 604 may be formatted such as to accommodate the columns and rows of any arbitrate configuration of array 606. In addition, in circumstances in which weight data matrix 604 is known then row 620 may be generated and stored in a vectorized format without the use of a formatter.

[0064] In embodiments, dot product 630 in FIG. 6 is generated by dot-multiplying a vector corresponding to column 610 with a vector corresponding to row 620. In embodiments, as shown in FIG. 7, the next dot product 632 may be obtained by dot-multiplying a vector corresponding to column 612 with the vector corresponding to row 620. As those of skill in the art will recognize, once all dot products in the first row of array 606 are filled, the dot product of the second row of array 606 may be calculated by dot-multiplying the elements in first column 610 of input data matrix 602 with the second row of weight data matrix 604, etc.

[0065] It is important to note that FIG. 6 through FIG. 8 merely serve illustrative purposes and that the abovementioned dot-multiplications may be simultaneously performed to generate a one-shot matrix-matrix multiply operation.

[0066] FIG. 9 illustrates an exemplary deconvolution operation according to various embodiments of the present disclosure. Deconvolution system 900 comprises input channels IC of input image 902, weights 922, dot product 904 906, and output channels OC. A person of skill in the art will recognize that, the deconvolution operation 900 is, in effect, is a mathematical transposition (approximately the inverse) of the convolution operation, for example, the convolution shown in FIG. 5. One of skill in the art will further recognize that a neural network may be used to learn deconvolution operation 900 by applying procedures similar to those used for ordinary convolutional neural networks. For purposes of brevity, a description or functions of components similar to those in FIG. 5 is not repeated here.

[0067] In embodiments, deconvolution operation 900 in FIG. 9 reassembles matrices 912 by deconstructing dot product 904 906 using weights 922. As with a convolution operation, deconvolution 900 may use a different set of weights 922 for each input channel IC. In embodiments, deconvolution 900 may be advantageously applied to an image to perform image deconvolution, for example to

improve robustness against artifacts. Other applications may include analysis and restoration of image data, and the like.

[0068] FIG. 10 illustrates a process for performing arithmetic operations to accelerate convolutional neural networks according to various embodiments of the present disclosure.

[0069] Process 1000 for performing arithmetic operations begins at step 1002 when a first set of operands that may be representative of a row in a data matrix is received from a first logic circuit. This first set of operands may be vectorized such that the operands are aligned with inputs into a matrix processor. In certain embodiments, the size of the vectorized operands is directly related to the number of inputs into a matrix processor along on axis.

[0070] At step 1004, a second set of operands that may be representative of a column in a weight matrix is received from a second logic circuit. This second set of operands may be vectorized such that the operands are aligned within corresponding inputs into the matrix processor. In certain embodiments, the size of the vectorized operands is directly related to the number of inputs into the matrix process along a different axis.

[0071] At step 1006, the first set of operands is dot-multiplied with the second set of operands to obtain one or more dot-products. In certain embodiments, this set operation across the sets of operands is performed in a single clock cycle.

[0072] At step 1008, the dot-products may be used to convolve an image with a filter to produce a convolution result.

[0073] At step 1010, the convolution result is further processed to enhance the image output. This further processing may occur using a non-linear function, a normalization operation or a pooling operation.

[0074] One skilled in the art will recognize no computing system or programming language is critical to the practice of the present invention. One skilled in the art will also recognize that a number of the elements described above may be physically and/or functionally separated into sub-modules or combined together.

[0075] It shall be noted that elements of the claims below may be arranged differently including having multiple dependencies, configurations, and combinations. For example, in embodiments, the subject matter of various claims may be combined with other claims.

[0076] It will be appreciated to those skilled in the art that the preceding examples and embodiment are exemplary and not limiting to the scope of the present invention. It is intended that all permutations, enhancements, equivalents, combinations, and improvements thereto that are apparent to those skilled in the art upon a reading of the specification and a study of the drawings are included within the true spirit and scope of the present invention.

1. A matrix processor for accelerating convolutions in a neural network, the matrix processor comprising:

- a first input circuit arranged in a first dimension of a two-dimensional array, the first input circuit being coupled to receive N operands from a first logic circuit, the N operands being formatted in accordance with a first width related to the first dimension;

- a second input circuit arranged in a second dimension of the two-dimensional array, the second input circuit coupled to receive M operands from a second logic

circuit, the M operands being formatted in accordance with a second width related to the second dimension; and

- a plurality of sub-circuits coupled to receive the N operands and the M operands, at least a subset of the plurality of sub-circuits comprising an arithmetic logic unit, an accumulator and a shadow register, the sub-circuits coupled within the two-dimensional array to perform an arithmetic operation on the N operands and the M operands.

2. The matrix processor according to claim 1 wherein the arithmetic operation is a dot product calculation associated with a convolution operation.

3. The matrix processor according to claim 2 wherein the arithmetic logic unit comprises a multiply-and-add circuit to generate the dot product.

4. The matrix processor according to claim 1 wherein the N operands represent image data and the M operands represent weight values.

5. The matrix processor according to claim 1 wherein at least some of the sub-circuits comprise an encoding element configured to encode values representing one or more of the M operands.

6. The matrix processor according to claim 5 wherein the encoding element is a booth encoder.

7. The matrix processor according to claim 1 wherein the N operands are formatted from a data input matrix.

8. The matrix processor according to claim 1 further comprising a state machine that uses at least one of a filter size and a stride to determine reusable operands within the N operands or the M operands.

9. The matrix processor according to claim 1 wherein accelerated processing speed is achieved by a reduction in read operations from a cache and accelerated data throughput via the plurality of sub-circuits.

10. A system for accelerating convolutions in a neural network, the system comprising:

- a first logic circuit that generates N operands;

- a first input circuit arranged in a first dimension of a two-dimensional array, the first input circuit being coupled to receive the N operands from the first logic circuit;

- a second logic circuit that generates M operands;

- a second input circuit arranged in a second dimension of the two-dimensional array, the second input circuit being coupled to receive M operands from the second logic circuit;

- a matrix processor comprising a plurality of sub-circuits, the plurality of sub-circuits configured to perform dot-multiplications of the N operands and the M operands to generate dot-products; and

- an output array coupled to the two-dimensional array, the two-dimensional array configured to use the dot-products to generate a result.

11. The system according to claim 10 wherein the N operands are formatted from a data input matrix into a first vector and the M operands are formatted from a weight input matrix into a second vector.

12. The system according to claim 11 wherein the first logic circuit comprises a plurality of data registers that store the N operands, the plurality of data registers having a first width corresponding to the first dimension of the two dimensional array and the second logic circuit comprises a plurality of weight registers that store the M weight oper-

ands, the plurality of weight registers having a second width corresponding to the second dimension of the two dimensional array.

13. The system according to claim **12** wherein the first width corresponds to a number of cycles that generate the result.

14. The system according to claim **12** wherein the data register and the weights register are accessed only once to fetch respective a first and second number of elements.

15. The system according to claim **10** wherein the sub-circuits comprise shadow registers configured to move data, in one or more clock cycles, to a shift register.

16. The system according to claim **10** further comprising a buffer coupled to at least one of the data input matrix and the weight input matrix, the buffer stores a copy of recently used data to enable reuse without refetching in subsequent cycles.

17. The system according to claim **10** wherein the result is an output matrix that corresponds to an application of a filter to an area of an image.

18. The system according to claim **10** further comprising a state machine that uses at least one of a filter size and a stride to identify reusable data.

19. A method for using a matrix multiplication circuit to make convolutional neural networks faster, the method comprising:

receiving, from a first logic circuit, a first set of operands representative of a row in a data matrix;

receiving, from a second logic circuit, a second set of operands representative of a column in a weight matrix; dot-multiplying the first set of operands with the second set of operands to obtain one or more dot-products; and using the dot-products to convolve an image with a filter to produce a convolution result.

20. The method according to claim **19** wherein convolving the image comprises processing the one or more dot-products by a convolution layer to generate a layer output.

21. The method according to claim **20** wherein generating the layer output comprises applying one of a non-linear function, a normalization, and a pooling to the convolution result.

* * * * *