(54) **STREAM-BASED ACCELERATOR PROCESSING OF COMPUTATIONAL GRAPHS**

STROMBASIERTE BESCHLEUNIGERVERARBEITUNG VON RECHENDIAGRAMMEN

TRAITEMENT D'ACCÉLÉRATEUR REPOSANT SUR UN FLUX DE GRAPHES DE CALCUL

(72) Inventors:
• **BARHAM, Paul Ronald**
**Mountain View, CA 94043 (US)**
• **VASUDEVAN, Vijay**
**Mountain View, CA 94043 (US)**

(74) Representative: **Suckling, Andrew Michael et al**
**Marks & Clerk LLP**
**2nd Floor, Wytham Court**
**11 West Way**
**Oxford OX2 0QL (GB)**

(56) References cited:
**US-A1- 2006 095 721      US-A1- 2006 095 722**
**US-A1- 2015 007 182**

EP 3 353 655 B1

**Description**

BACKGROUND

**[0001]** This specification relates to processing a computational graph representing a neural network by assigning a subgraph to an accelerator device, e.g., a graphical processing unit (GPU), having multiple streams and/or to use of such a processed computational graph for processing a model input.

**[0002]** Neural networks are machine learning models that employ one or more layers of models to generate an output, e.g., one or more classifications, for a received input. Some neural networks include one or more hidden layers in addition to an output layer. The output of each hidden layer is used as input to the next layer in the network, i.e., the next hidden layer or the output layer of the network. Each layer of the network generates an output from a received input in accordance with current values of a respective set of parameters for the layer.

**[0003]** In systems that exist, the operations of computational graphs can be processed by an individual device. In some implementations, the device is a GPU. The device can have a processor that performs operations, e.g., generating outputs at a layer from inputs, and stores outputs from the operations in memory. Due to the large number and size of operations generally required to generate the outputs in the computational graph, one device can take a significant amount of time to process the operations of the graph.

**[0004]** US 2006/095722 A1 proposes an apparatus for processing data under control of a program having program instructions and subgraph suggestion information identifying respective sequences of program instructions corresponding to computational subgraphs identified within said program. The apparatus comprises: a memory operable to store a program formed of separate program instructions; processing logic operable to execute respective separate program instructions from said program; and accelerator logic operable in response to reaching an execution point within said program associated with a subgraph suggestion to execute a sequence of program instructions corresponding to said subgraph suggestion as an accelerated operation instead of executing said sequence of program instructions as respective separate program instructions with said processing logic.

**[0005]** US 2015/007182 A1 proposes techniques and constructs to improve execution speed of distributed iterative computation using heterogeneous specialized resources including, for example, processors and accelerators. Iteration over an arbitrary sub-graph without loop unrolling including for algorithms with data-dependent loop termination and large iteration counts, including as a result of nested iteration, are supported in a resource-efficient manner without adding vertices to a dataflow graph to represent iteration constructs. Instead, some or all of the existing vertices within the sub-graph that is to be iterated upon based on having additional and/or modified ports and channels associated with them.

SUMMARY

**[0006]** In general, this specification describes a system or method for processing subgraphs of a computational graph using a stream-based accelerator device, e.g., a GPU.

**[0007]** In general, one innovative aspect of the subject matter described in this specification can be embodied in computer-implemented methods that include the actions of receiving a request to process a computational graph; obtaining data representing multiple subgraphs of the computational graph, the computational graph comprising a plurality of nodes and directed edges, wherein each node represents a respective operation, wherein each directed edge connects a respective first node to a respective second node that represents an operation that receives, as input, an output of an operation represented by the respective first node; assigning the multiple subgraphs to respective devices of a plurality of devices, including assigning the first subgraph to a first; determining that the first device comprises a hardware accelerator having a plurality of streams, each stream being an independent hardware queue whose operations are processed in order; in response to determining that the first device comprises a hardware accelerator having a plurality of streams, generating instructions that when executed by the first device cause the first device to: assign the operation represented by each node in the first subgraph to a corresponding stream in the plurality of streams of the hardware accelerator, including assigning two operations that do not depend on one another to different streams, and assigning to a single stream a first operation that changes the internal state of the hardware in a way that must happen before a second operation executes and assigning the second operation to be performed in the single stream after the first operation is complete; and perform the operations represented by the nodes in the subgraph in accordance with the assignment; and providing the instructions and the data to the first device. The instructions further cause the first device to determine a particular operation represented by a node has finished at a particular stream; and, in response to determining the particular operation has finished: determine a first amount of memory consumed by the particular operation that will be freed; determine, for each of a group of unassigned nodes, a respective estimated amount of memory consumed by an operation that is represented by the unassigned node; determine, from the group of unassigned nodes, a first unassigned node that represents an operation, which executes on a stream of the hardware accelerator, with the estimated amount of memory that maximizes usage of the first amount of memory; and assign an operation represented by the first unassigned node to the particular stream. A method of this aspect may be performed by

one or more computing devices, for example by one or more computing devices that comprise a computational graph system.

**[0008]** Implementations can include one or more of the following features. The request specifies identifying one or more particular outputs from one or more respective nodes in the subgraph, and the method further comprises: receiving, from the first device, the one or more particular outputs; and providing the one or more particular outputs to the client. The instructions further cause the first device to store the one or more particular outputs in memory of the first device. The operations for the subgraph comprise partial inference or training computations for a neural network. Analyzing the subgraph to identify a group of nodes in the subgraph in a chain structure; wherein the instructions cause the first device to assign the group of nodes to one stream. The assigning comprises: analyzing the subgraph to identify a first node in the subgraph has a plurality of directed edges as outputs; wherein the instructions cause the first device to assign, for each of the directed edges, a node to which the directed edge points to a unique stream of the graphical processing unit. The instructions cause the first device to determine, for each node, a respective amount of memory resources in the graphical processing unit consumed by the operation represented by the node based on the directed edges to the node, wherein the assigning is based at least on the respective amount of memory resources. The instructions cause the first device to determine a particular operation represented by a node has finished at a particular stream: in response to determining the particular operation has finished: determine at least one subsequent operation that uses the output of the particular operation as input; and reuse memory allocated for the output of the particular operation after the at least one subsequent operation has executed. Determining at least one subsequent operation that uses the output of the particular operation as input includes: determining that at least two subsequent operations, a first operation in a first stream and a second operation in a second stream, use the output of the particular operation as input; placing a first marker in a first stream that indicates when the first operation has used the particular operation as input; placing a second marker in a second stream that indicates when the second operation has used the particular operation as input; determining that both operations have used the particular operation upon indication from the first and second markers.

**[0009]** In an implementation the method further comprises: receiving a model input; and processing, by the hardware accelerator, the model input according the operations represented by the nodes in the subgraph.

**[0010]** In another aspect, the subject matter described in this specification may be embodied in methods that may include the actions of providing a machine learning model corresponding to a processed computational graph obtained by a method of the first aspect; and processing, using the machine learning model, a model input.

**[0011]** In another aspect, the subject matter described in this specification may be embodied in methods that may include the actions of executing, by a hardware accelerator, a subgraph of a processed computational graph obtained by a method of the first aspect.

**[0012]** Other implementations of these and other aspects include corresponding systems, and computer programs, configured to perform the actions of the methods, encoded on computer storage media (which may or may not be non-transitory storage media).

**[0013]** Particular embodiments of the subject matter described in this specification can be implemented so as to realize one or more of the following advantages. Operations, e.g., an operation to generate an inference from an input, of a neural network can be represented as a computational graph of nodes and directed edges. A system processes this computational graph representation to efficiently perform the operations. The system achieves this efficiency because the computational graph has multiple streams. Using multiple streams can allow logically independent operations to be reordered or executed concurrently. When the system has a goal of lowering end-to-end latency for a whole computation, the example system may reorder logically independent operations. When the system has a goal to achieve higher throughput, the example system may execute operations simultaneously. The computational graph can be more easily partitioned for parallel operations than the conventional representation. By way of illustration, subgraphs of the computational graph can be assigned to unique devices, each of which performs operations in the respective subgraph, to reduce an overall time required to perform operations of the neural network.

**[0014]** A device to which a subgraph is assigned can be a GPU. The subgraph can be partitioned into multiple streams of the GPU to more efficiently perform the operations of the subgraph. The details of one or more embodiments of the subject matter of this specification are set forth in the accompanying drawings and the description below. Other features, aspects, and advantages of the subject matter will become apparent from the description, the drawings, and the claims. It will be appreciated that aspects and implementations can be combined, and that features described in the context of one aspect or implementation can be implemented in the context of other aspects or implementations.

BRIEF DESCRIPTION OF THE DRAWINGS

**[0015]**

FIG. 1 illustrates an example computational graph system for distributing operations for neural networks represented as computational graphs.
FIG. 2 is a flow diagram of an example process for processing a subgraph of a computational graph using a GPU.

FIG. 3 illustrates an example subgraph of a computational graph being processed by a GPU.

FIG. 4 is a flow diagram of an example process for assigning nodes to streams.

[0016] Like reference numbers and designations in the various drawings indicate like elements.

DETAILED DESCRIPTION

[0017] This specification generally describes a computational graph system that performs operations represented by a computational graph in a distributed manner.

[0018] The computational graph includes nodes connected by directed edges. Each node in the computational graph represents an operation. An incoming edge to a node represents a flow of an input into the node, i.e., an input to the operation represented by the node. An outgoing edge from a node represents a flow of an output of the operation represented by the node to be used as an input to an operation represented by another node. Thus, a directed edge connecting a first node in the graph to a second node in the graph indicates that an output generated by the operation represented by the first node is used as an input to the operation represented by the second node.

[0019] Generally, the input and outputs flowing along directed edges in the computational graph are tensors. A tensor is a multidimensional array of numeric values or other values, e.g., strings, having a specific order that corresponds to the dimensionality of the array. For example, a scalar value is a 0th-order tensor, a vector of numeric values is a 1st-order tensor, and a matrix is a 2nd-order tensor.

[0020] In some implementations, the operations represented in the computational graph are neural network operations or operations for a different kind of machine learning model. A neural network is a machine learning model that employs one or more layers of nonlinear units to predict an output for a received input. Some neural networks are deep neural networks that include one or more hidden layers in addition to an output layer. The output of each hidden layer is used as input to another layer in the network, i.e., another hidden layer, the output layer, or both. Some layers of the network generate an output from a received input in accordance with current values of a respective set of parameters, while other layers of the network may not have parameters.

[0021] For example, the operations represented by the computational graph may be operations necessary for the neural network to compute an inference, i.e., to process an input through the layers of the neural network to generate a neural network output for the input. As another example, the operations represented by the computational graph may be operations necessary to train the neural network by performing a neural network training procedure to adjust the values of the parameters of the neural network, e.g., to determine trained values of the parameters from initial values of the parameters. In some cases, e.g., during training of the neural network, the operations represented by the computational graph can include operations performed by multiple replicas of the neural network.

[0022] By way of illustration, a neural network layer that receives an input from a previous layer can use a parameter matrix to perform a matrix multiplication between the parameter matrix and the input. In some cases, this matrix multiplication can be represented as multiple nodes in the computational graph. For example, a matrix multiplication can be divided into multiple multiplication and addition operations, and each operation can be represented by a different node in the computational graph. The operation represented by each node can generate a respective output, which flows on a directed edge to a subsequent node. After the operation represented by a final node generates a result of the matrix multiplication, the result flows, on a directed edge, to another node. The result is equivalent to an output of the neural network layer that performs the matrix multiplication.

[0023] In some other cases, the matrix multiplication is represented as one node in the graph. The operations represented by the node can receive, as inputs, an input tensor on a first directed edge and a weight tensor, e.g., a parameter matrix, on a second directed edge. In some implementations, the weight tensor is associated with the shared persistent state of the model. The node can process, e.g., perform a matrix multiplication of, the input and weight tensors to output, on a third directed edge, an output tensor, which is equivalent to an output of the neural network layer.

[0024] Other neural network operations that may be represented by nodes in the computational graph include other mathematical operations, e.g., subtraction, division, and gradient computations; array operations, e.g., concatenate, splice, split, or rank; and neural network building block operations, e.g., SoftMax, Sigmoid, rectified linear unit (ReLU), or convolutions.

[0025] Representing a neural network as a computational graph provides for a flexible and granular way to efficiently implement the neural network, especially if the operations for the neural network are distributed across multiple devices with different hardware profiles.

[0026] FIG. 1 illustrates an example computational graph system 100 for distributing operations for neural networks represented as computational graphs. The system 100 is an example of a system implemented as computer programs on one or more computers in one or more locations, in which the systems, components, and techniques described below can be implemented.

[0027] A user of a client 102 can request actions be performed on a computational graph representing a neural network. For example, a client can register a graph with the session manager, feed data input into the graph, or evaluate one or more of the outputs of a graph. The client 102 can be an application running on a computer.

[0028] As part of the request, the client 102 provides

data identifying a computational graph to the system 100 and specifies types of actions to be performed on the computational graph.

[0029]    For example, the request can identify a computational graph representing an inference for a particular neural network and can identify an input on which the inference should be performed.

[0030]    As another example, the request can identify a computational graph representing a training procedure for a particular neural network and can identify an input, such as training data, on which the training should be performed. In this example, when receiving a request to process a computation graph representing a training procedure, the system 100 can determine modified values for parameters for one or more edges of the computational graph, e.g., using conventional backpropagation or other neural network training techniques. The system 100 can store the modified parameters in memory of a device, and an executor 106 can retrieve and store, at the system 100, addresses of the modified weights. Upon further requests from the client 102 for inference, training, or other operations requiring the modified weights, the system 100 can access the modified weights using the addresses.

[0031]    In some cases, the request may specify a response that should be transmitted in response to the request. For example, for a neural network training request, the client 102 can request an indication that the requested neural network training operations have been completed and, optionally, trained values of the parameters of the neural network or an indication of a memory location from which the trained values can be accessed by the client 102. As another example, for a neural network inference request, the client 102 can request output values that represent an inference operation from one or more particular nodes of the computational graph.

[0032]    The system 100 performs the operations to generate the particular output by partitioning the operations represented by the computational graph across multiple devices 116-122. The system 100 partitions the operations to the multiple devices 116-122 over a data communication network 114, e.g., local area network (LAN) or wide area network (WAN). The devices 116-122 perform the operations and, if applicable, return a respective output or indication to the system 100, which can return the requested output or indication to the client 102.

[0033]    Any devices performing neural network operations, e.g., devices 116-122, can include a memory, e.g., a random access memory (RAM), for storing instructions and data and a processor for executing stored instructions. Generally, each device is a hardware resource that performs operations independent of other devices. For example, each device can have its own processing unit. The devices can be Graphical Processing Units (GPUs),Central Processing Units (CPUs), or other accelerators. By way of illustration, one machine can host one or more devices, e.g., multiple CPUs and GPUs.

[0034]    Each device can also have a respective computational capability. That is, devices can have different amount of memories, processing speed, or other architectural characteristics. Thus, some devices can perform operations that other devices cannot. For example, some operations require a certain amount of memory that only particular devices have, or some devices are configured to only perform a particular type of operation, e.g., inference operations.

[0035]    A session manager 104 in the system 100 receives a request from the client 102 to start a session during which operations of the computational graph are performed. The session manager 104 manages the set of devices, e.g., devices 116-122, that can perform operations of the computational graph, and can provide a placer 108 with the set of devices that are available to perform operations.

[0036]    The placer 108 determines, for each operation to be performed in the computational graph, a respective target device, e.g., device 116, that performs the operation, and in some implementations, a time for the respective target device to perform the operation. The placer 108 performs optimal device assignment by knowing how long an operation will take on each available device given the size of the input data. The placer 108 obtains the estimate of processing time using measurements or predictive performance models. Some operations can be performed in parallel while other operations require prior operations in the computational graph to be completed, e.g., the other operations process, as inputs, outputs of the prior operations.

[0037]    After the devices perform the operations allocated by the placer 108 to generate outputs, the executor 106 can retrieve the outputs. The executor 106 can generate an appropriate response to the request, e.g., an output or an indication that the processing has been completed. Then, the executor 106 can return the response to the client 102. Although FIG. 1 illustrates one executor 106, in one implementation, there is an executor per device. This executor issues operations to the device when they become runnable (i.e. all of their inputs have been computed). This implementation also has a graph manager that partitions a graph to run on multiple devices by invoking the placer 108 and creates the necessary executors.

[0038]    The session manager 104 also provides sets of operations to be performed in the computational graph to the executor 106. The executor 106 periodically retrieves runtime statistics from the devices 116-122 related to graph execution of operations. The executor 106 provides the runtime statistics to the placer 108, which can re-optimize placement and scheduling of further operations.

[0039]    FIG. 2 is a flow diagram of an example process 200 for processing a subgraph of a computational graph using a GPU. For convenience, the process 200 will be described as being performed by a system of one or more computers located in one or more locations. For example, a computational graph system, e.g., the computational

graph system 100 of FIG. 1, appropriately programmed, can perform the process 200.

**[0040]** The system receives a request from a client to process a computational graph (step 202). For example, the request can be a request to perform a neural network inference represented by the computational graph on a specified input, a request to perform neural network training operations represented by the computational graph on a specified set of training data, or a request to perform other neural network operations represented by the computational graph, as described above with reference to FIG. 1.

**[0041]** In some cases, a computational graph is sent with the request from the client. In other cases, the request identifies the computational graph and the system retrieves the data representing the identified graph from memory.

**[0042]** The system can partition the computational graph into multiple subgraphs. In some implementations, the subgraphs are specified by the client sending the request, and the system partitions the computational graph according to the specifications. In some other implementations, the system partitions the computational graph such that each subgraph requires a similar amount of resources for performing operations compared to the other subgraphs.

**[0043]** The system can assign each subgraph to an available device, e.g., using placer 108 of FIG. 1.

**[0044]** The system obtains data representing a particular subgraph of the computational graph (step 204) from the partitioned computational graph. The data can be obtained from a database or memory of the system. By way of illustration, operations of the particular subgraph represent partial inference or training computations.

**[0045]** The system determines that a device to which the subgraph is assigned is a graphical processing unit or other hardware accelerator device having multiple streams (step 206). By way of illustration, the system can assess whether the device is a GPU with multiple streams by requesting a type of the device from a resource manager that manages devices to be assigned to the computational graph. Each stream is an independent hardware queue whose operations are processed in order.

**[0046]** The system generates instructions that, when executed by the device, cause the device to perform particular operations (step 208). In particular, the instructions cause the device to assign the operation represented by each node in the subgraph to a respective stream of the device.

**[0047]** An example system may assign computations of some hardware accelerators to streams in a particular way (e.g. if one operation executes on stream A, then a later, related operation must also execute on stream A.) For example, a first operation may be stateful and execute on stream A. By executing, the first operation may change the internal state of the hardware in a way that must happen before a second operation executes. The second operation may then execute on stream A after

the first operation is complete.

**[0048]** In some implementations, two internal hardware resources cannot be used simultaneously and therefore need to be serialized.

**[0049]** Generally, the device assigns operations that do not depend on each other to different streams. By assigning operations that do not depend on each other to different streams, the hardware does not need to know how long an operation will take and can choose from a number of available operations to execute the first one that is ready to execute without expensive host intervention.

**[0050]** The instructions also cause the device to perform the operations represented by the nodes in the subgraph in accordance with the assignment. When operations are assigned to a particular stream, the operations are queued. The device can perform operations in a first-in-first-out (FIFO) manner. Thus, if the device only has one stream, the operations assigned to the device are performed serially. If the device has multiple streams, the operations in different streams can be performed in parallel and reordered with respect to each other, while the operations within a given stream are performed serially. Performing operations using multiple streams decreases a total time to perform the operations of the subgraph. This is described further below with reference to FIGS. 3 and 4.

**[0051]** The system provides the instructions and the data to the device (step 210). In some implementations, the system sends the device a request to start the operations. The device receives the request and in response, executes the instructions received from the system. For example, the device may receive a model input, and processing the model input according to the operations represented by the nodes in the subgraph.

**[0052]** FIG. 3 illustrates an example subgraph 316 of a computational graph being processed by an Accelerator 302. The subgraph 316 has nodes 308-314, each of which represent an operation to be performed by the Accelerator 302. A computational graph system, e.g., the system 100 of FIG. 1, assigned the subgraph 316 to the Accelerator 302.

**[0053]** The Accelerator 302 has two streams 304 and 306. The streams share utilization of the Accelerator 302. In GPU, streams may be symmetric, meaning that all operations can be performed on any stream. This symmetry may not be available of all accelerator devices. For example, on specific accelerator devices certain streams must be used to perform operations that copy data between host and device memory.

**[0054]** The computational graph system can analyze the subgraph 316 to determine how the subgraph 316 is assigned to the multiple streams 304 and 306. In some implementations, the system generates instructions that causes the Accelerator 302 to assign the nodes of the subgraph 316 in a way that minimizes the number of times a directed edge connects to different streams. There may be a performance cost to enforcing depend-

encies between streams. Ordering instructions has some overhead cost. Every ordering dependency reduces the number of possible execution orderings available to the device, reducing scheduling flexibility. Each time a directed edge from a first stream connects to a second stream, the second stream waits for the operation with the directed edge from the first stream to the second stream to complete processing. Waiting can cause the second stream to remain idle, which causes the GPU to be inefficiently utilized.

[0055] In some implementations, the system generates instructions that causes the Accelerator 302 to assign the nodes of the subgraph 316 based on characteristics of the Accelerator 302. For example, the Accelerator 302 has a fixed number of streams, i.e., streams 304 and 306. The system can assign the nodes so each stream will be similarly utilized by the Accelerator 302. For accelerators that are GPUs, all streams share a single large pool of threads.

[0056] Some streams also perform particular operations that other streams do not. For example, stream 306 can perform direct memory access (DMA) operations while stream 304 does not. Thus, the system can analyze each node to determine a type of operation represented by the node, and the system can assign the node to a stream that is able to perform the type of operation. In GPUs, the main congested resources are DMA engines that copy data between hosts and device memory. DMA engines can be used by any stream. If one stream is executing a DMA operation, the stream cannot simultaneously execute a computation. An example system therefore ensures that at least one other stream has some compute work to execute at the same time. The system can analyze the subgraph to identify, and thus, generate instructions that causes a software module or driver that manages assigning operations to assign nodes by following two general rules. First, the system tries to assign nodes arranged in a chain structure to the same stream. Nodes in a chain structure are nodes that are connected to each other by following one directed edge from node to node. Thus, a node in the chain must wait for operations at previous nodes in the chain to finish computing before computing its own operation. Assigning chains of nodes is not always possible since branching and merging occur in the graph, e.g. from shared input variables or common subexpressions.

[0057] Second, the system can choose to generate instructions that cause the Accelerator 302 to assign multiple nodes that each receive input from one node to unique streams. That is, if a first node has multiple outputs to multiple different nodes, the system assigns each of the different nodes to a unique stream. Each of the different nodes do not have data dependence on any of the other different nodes, and therefore, improve efficiency when operating on disjoint streams.

[0058] By way of illustration, the Accelerator 302 receives the subgraph 316. The instructions received by the system cause the Accelerator 302 to assign the initial node 308 to a first stream 306. The initial node 308 has two outputs - one directed edge to node 310 and one directed edge to node 314. Therefore, using the second rule, the instructions cause the Accelerator 302 to assigns nodes 310 and 314 to different streams. Node 312 also only receives, as input, an output of the node 310. Therefore, using the first rule, the system assigns node 312 to the same stream, i.e., stream 304, as the node 310.

[0059] As described above, streams are hardware queues whose operations are performed in order. Thus, the order in which the Accelerator 302 assigns nodes to streams matters. The Accelerator 302 assigns nodes to streams in an order of the direction of data flow in the subgraph. That is, the Accelerator 302 identifies one or more initial nodes of the subgraph and assigns the one or more initial nodes. Then, the Accelerator 302 follows directed edges that are outputs of the one or more initial nodes to identify subsequent nodes, and the Accelerator 302 assigns the subsequent nodes to respective streams. The Accelerator 302 continues assignment of nodes until each node in the subgraph is assigned. As a result of assigning nodes in this order, operations within a given stream will also be performed in the order in which the operations were assigned, as described above. When the inputs of an operation A are produced on different streams, it is necessary to ensure that they have all been computed before operation A is executed. The execution on the stream to which operation A is assigned should be stalled until all of the inputs to operation A have been computed. The exact stalling mechanism is device specific. For GPU devices, an event can be created for each of the input streams and instructions can be added to each stream to signal the event. For each input, an instruction can also be added to the stream on which A is assigned in order for the operation to wait for the relevant event in order to execute. In cases where one or more of the inputs for operation A are computed on the same stream as operation A, dataflow dependency instructions can be safely deleted, leading to better performance. Within a given stream, operations represented by nodes assigned to the given stream that generate an output that is used as input by operations represented by one or more other nodes assigned to the given stream will have been already computed or scheduled to be computed when the Accelerator 302 performs the operations represented by the one or more other nodes.

[0060] Continuing with the illustration above, stream 304 is assigned node 310 and then assigned node 312 because data flows from the node 310 to the node 312. When executing operations in the stream, the Accelerator 302 first executes operations represented by the node 310 and then executes operations represented by the node 312.

[0061] After the final nodes, i.e., nodes 312 and 314, performs operations, the Accelerator 302 return the outputs of the nodes or an indication the operations have completed to the system. In an example system, there is a special 'send' node that copies the computation re-

sults back from the memory of the Accelerator 302 into the host memory where it can be handed to a different device by a receive node or returned to the client in a remote procedure call (RPC) response. The system can then, if necessary, return the output or the indication to the client.

**[0062]** Another implementation of assigning nodes to streams will be described further below with reference to FIG. 4.

**[0063]** FIG. 4 is a flow diagram of an example process 400 for assigning subgraphs to devices. For convenience, the process 400 will be described as being performed by a system, e.g., a GPU. For example, a GPU can receive instructions generated by a computational graph system, e.g., the computational graph system 100 of FIG. 1, that, when executed, cause the GPU to perform the process 400.

**[0064]** The system can assign a particular node to a stream based on an amount of memory resources consumed by the node or by previously assigned nodes. For example, the system can calculate a dimension of a tensor on each directed edge to and from each node of the subgraph. The dimensions of the tensors indicate a size of memory that would be consumed by a device to perform an operation. The system may need to calculate all dimensions of a tensor in order to determine the size. The system can then assign particular nodes with tensors consuming a particular size of memory to devices having the particular size of memory.

**[0065]** In particular, when the device performs the operation, the software driver or executor allocates memory to store any inputs as well as any outputs computed as a result of the operation. Because the amount of memory on the device is limited, the device frees memory when memory is no longer used.

**[0066]** By way of illustration, the system determines whether an operation represented by a node has finished at a particular stream (step 402). For example, the system can periodically poll streams to determine whether the operation in the particular stream has finished. The stream may support an action that allows the host to determine how far execution has progressed through the list of operations in the stream. In some implementations, events, or markers, can signal how far execution has progressed. When an event occurs, the event can be added to a special hardware operation queue in the stream. The host can poll this queue in order to determine which operations have occurred. Other stream implementations may only allow the host to determine when all enqueued operations are complete. Alternatively or additionally, the hardware can provide an interrupt or call-back when the stream reaches a certain point.

**[0067]** When the operation has finished, the system can determine memory used for inputs to the operation can be freed for use in other operations. The system does not free memory used for outputs of the operation because the outputs may be used in a subsequent node.

**[0068]** Thus, the system determines an amount of memory consumed that will be freed (step 404). The system can send a request to the software driver or executor to identify the size of memory that will be freed.

**[0069]** In some implementations, an example system allows the use of remote direct memory access (RDMA) network interfaces that remote machines can use to directly transfer data into the memory of a hardware accelerator at an arbitrary point in time. This memory must not be in use by any other operation running on any stream. The example system may not need to know precisely how far operations on each stream has progressed. However, the system should keep track of memory known not to be in use by any stream. This free memory can then be used for RDMA.

**[0070]** The system determines, for each of a group of unassigned nodes, a respective estimated amount of memory consumed by the unassigned node (step 406). The unassigned nodes can include nodes that receive inputs from the node whose operation has completed. The unassigned nodes can also include nodes that are independent from the node whose operation has completed but still need to be processed by the accelerator. The estimated amount of memory can be determined by evaluating dimensions of the respective tensors to the unassigned nodes, as described above.

**[0071]** The system determines, from the group of unassigned nodes, a first unassigned node that represents an operation, which when executed on a stream by the accelerator, maximizes usage of the amount of memory that will be freed (step 408). If an operation represented by an unassigned node requires more memory to execute than the amount of memory that will be free, the unassigned node will not be assigned to the stream. If a first and second operation require a respective estimated amount of memory less than or equal to the amount of memory that will be free, the system selects the operation that maximizes usage of the amount of memory that will be freed. In other words, in this case, the system determines the node representing the selected operation as the first unassigned node. An example system does not enqueue an operation on the stream until it can determine which regions of accelerator memory will be used to hold the temporary working space and outputs of the operation. In the event that memory is scarce, an example system may choose to enqueue operations that require smaller amounts of memory to execute or to preferentially enqueue operations that will consume large input tensors allowing them to be deallocated.

**[0072]** The system assigns an operation represented by the first unassigned node to the particular stream (step 410). The system can then cause the particular stream to perform the operation, and the system can continue operating as described above with reference to FIGS. 2-3.

**[0073]** Embodiments of the subject matter and the functional operations described in this specification can be implemented in digital electronic circuitry, in tangibly-embodied computer software or firmware, in computer

hardware, including the structures disclosed in this specification and their structural equivalents, or in combinations of one or more of them. Embodiments of the subject matter described in this specification can be implemented as one or more computer programs, i.e., one or more modules of computer program instructions encoded on a tangible non-transitory program carrier for execution by, or to control the operation of, data processing apparatus. Alternatively or in addition, the program instructions can be encoded on an artificially-generated propagated signal, e.g., a machine-generated electrical, optical, or electromagnetic signal, that is generated to encode information for transmission to suitable receiver apparatus for execution by a data processing apparatus. The computer storage medium can be a machine-readable storage device, a machine-readable storage substrate, a random or serial access memory device, or a combination of one or more of them. The computer storage medium is not, however, a propagated signal.

[0074]   The term "data processing apparatus" encompasses all kinds of apparatus, devices, and machines for processing data, including by way of example a programmable processor, a computer, or multiple processors or computers. The apparatus can include special purpose logic circuitry, e.g., an FPGA (field programmable gate array) or an ASIC (application-specific integrated circuit). The apparatus can also include, in addition to hardware, code that creates an execution environment for the computer program in question, e.g., code that constitutes processor firmware, a protocol stack, a database management system, an operating system, or a combination of one or more of them.

[0075]   A computer program (which may also be referred to or described as a program, software, a software application, a module, a software module, a script, or code) can be written in any form of programming language, including compiled or interpreted languages, or declarative or procedural languages, and it can be deployed in any form, including as a stand-alone program or as a module, component, subroutine, or other unit suitable for use in a computing environment. A computer program may, but need not, correspond to a file in a file system. A program can be stored in a portion of a file that holds other programs or data, e.g., one or more scripts stored in a markup language document, in a single file dedicated to the program in question, or in multiple coordinated files, e.g., files that store one or more modules, sub-programs, or portions of code. A computer program can be deployed to be executed on one computer or on multiple computers that are located at one site or distributed across multiple sites and interconnected by a communication network.

[0076]   As used in this specification, an "engine," or "software engine," refers to a software implemented input/output system that provides an output that is different from the input. An engine can be an encoded block of functionality, such as a library, a platform, a software development kit ("SDK"), or an object. Each engine can be implemented on any appropriate type of computing device, e.g., servers, mobile phones, tablet computers, notebook computers, music players, e-book readers, laptop or desktop computers, PDAs, smart phones, or other stationary or portable devices, that includes one or more processors and computer readable media. Additionally, two or more of the engines may be implemented on the same computing device, or on different computing devices.

[0077]   The processes and logic flows described in this specification can be performed by one or more programmable computers executing one or more computer programs to perform functions by operating on input data and generating output. The processes and logic flows can also be performed by, and apparatus can also be implemented as, special purpose logic circuitry, e.g., an FPGA (field programmable gate array) or an ASIC (application-specific integrated circuit).

[0078]   Computers suitable for the execution of a computer program include, by way of example, can be based on general or special purpose microprocessors or both, or any other kind of central processing unit. Generally, a central processing unit will receive instructions and data from a read-only memory or a random access memory or both. The essential elements of a computer are a central processing unit for performing or executing instructions and one or more memory devices for storing instructions and data. Generally, a computer will also include, or be operatively coupled to receive data from or transfer data to, or both, one or more mass storage devices for storing data, e.g., magnetic, magneto-optical disks, or optical disks. However, a computer need not have such devices. Moreover, a computer can be embedded in another device, e.g., a mobile telephone, a personal digital assistant (PDA), a mobile audio or video player, a game console, a Global Positioning System (GPS) receiver, or a portable storage device, e.g., a universal serial bus (USB) flash drive, to name just a few.

[0079]   Computer-readable media suitable for storing computer program instructions and data include all forms of non-volatile memory, media and memory devices, including by way of example semiconductor memory devices, e.g., EPROM, EEPROM, and flash memory devices; magnetic disks, e.g., internal hard disks or removable disks; magneto-optical disks; and CD-ROM and DVD-ROM disks. The processor and the memory can be supplemented by, or incorporated in, special purpose logic circuitry.

[0080]   To provide for interaction with a user, aspects of the subject matter described in this specification can be implemented on a computer having a display device, e.g., a CRT (cathode ray tube) monitor, an LCD (liquid crystal display) monitor, or an OLED display, for displaying information to the user, as well as input devices for providing input to the computer, e.g., a keyboard, a mouse, or a presence sensitive display or other surface. Other kinds of devices can be used to provide for interaction with a user as well; for example, feedback provided

to the user can be any form of sensory feedback, e.g., visual feedback, auditory feedback, or tactile feedback; and input from the user can be received in any form, including acoustic, speech, or tactile input. In addition, a computer can interact with a user by sending resources to and receiving resources from a device that is used by the user; for example, by sending web pages to a web browser on a user's client device in response to requests received from the web browser.

[0081] Aspects of the subject matter described in this specification can be implemented in a computing system that includes a back end component, e.g., as a data server, or that includes a middleware component, e.g., an application server, or that includes a front end component, e.g., a client computer having a graphical user interface or a Web browser through which a user can interact with an implementation of the subject matter described in this specification, or any combination of one or more such back end, middleware, or front end components. The components of the system can be interconnected by any form or medium of digital data communication, e.g., a communication network. Examples of communication networks include a local area network ("LAN") and a wide area network ("WAN"), e.g., the Internet.

[0082] The computing system can include clients and servers. A client and server are generally remote from each other and typically interact through a communication network. The relationship of client and server arises by virtue of computer programs running on the respective computers and having a client-server relationship to each other.

[0083] While this specification contains many specific implementation details, these should not be construed as limitations on the scope of any invention or of what is claimed, but rather as descriptions of features that may be specific to particular embodiments of particular inventions. Certain features that are described in this specification in the context of separate embodiments can also be implemented in combination in a single embodiment.

[0084] Similarly, while operations are depicted in the drawings in a particular order, this should not be understood as requiring that such operations be performed in the particular order shown or in sequential order, or that all illustrated operations be performed, to achieve desirable results. In certain circumstances, multitasking and parallel processing may be advantageous. Moreover, the separation of various system modules and components in the embodiments described above should not be understood as requiring such separation in all embodiments, and it should be understood that the described program components and systems can generally be integrated together in a single software product or packaged into multiple software products.

[0085] Particular embodiments of the subject matter have been described. Other embodiments are within the scope of the following claims. For example, the actions recited in the claims can be performed in a different order and still achieve desirable results.

**Claims**

1.  A computer-implemented method comprising:

    receiving (202) a request to process a computational graph;
    obtaining (204) data representing multiple subgraphs of the computational graph, the computational graph comprising a plurality of nodes and directed edges, wherein each node represents a respective operation, wherein each directed edge connects a respective first node to a respective second node that represents an operation that receives, as input, an output of an operation represented by the respective first node;
    assigning the multiple subgraphs to respective devices of a plurality of devices (116,118,120,122), including assigning a first subgraph (316) to a first device (302);
    determining (206) that the first device (302) comprises a hardware accelerator having a plurality of streams (304,306), each stream being an independent hardware queue whose operations are processed in order;
    in response to determining that the first device (302) comprises a hardware accelerator having a plurality of streams, generating (208) instructions that when executed by the first device cause the first device to:

        assign the operation represented by each node (308,310,312,314) in the first subgraph (316) to a corresponding stream in the plurality of streams (304,306) of the hardware accelerator, including assigning two operations that do not depend on one another to different streams (304,306), and assigning to a single stream (304) a first operation that changes the internal state of the hardware in a way that must happen before a second operation executes and assigning the second operation to be performed in the single stream (304) after the first operation is complete; and
        perform the operations represented by the nodes (308,310,312,314) in the first subgraph (316) in accordance with the assignment; and

    providing (210) the instructions and the data to the first device;
    wherein the instructions further cause the first device (302) to determine (402) a particular operation represented by a node has finished at a particular stream;
    in response to determining the particular operation has finished:

determine (404) a first amount of memory consumed by the particular operation that will be freed;

determine (406), for each of a group of unassigned nodes, a respective estimated amount of memory consumed by an operation that is represented by the unassigned node;

determine (408), from the group of unassigned nodes, a first unassigned node that represents an operation, which executes on a stream of the hardware accelerator, with the estimated amount of memory that maximizes usage of the first amount of memory; and

assign (410) an operation represented by the first unassigned node to the particular stream.

2. The method of claim 1, wherein the request specifies identifying one or more particular outputs from one or more respective nodes (308,310,312,314) in the first subgraph (316), further comprising: receiving, from the first device (302), the one or more particular outputs; and providing the one or more particular outputs to the client.

3. The method of claim 1 or 2, wherein the instructions further cause the first device (302) to store the one or more particular outputs in memory of the first device.

4. The method of claim 1, 2 or 3, wherein the operations for the first subgraph (316) comprise partial inference or training computations for a neural network.

5. The method of any preceding claim, further comprising:

analyzing the first subgraph (316) to identify a group of nodes in the first subgraph (316) in a chain structure;
wherein the instructions cause the first device (302) (316) to assign the group of nodes to one stream.

6. The method of any preceding claim, wherein the assigning comprises:

analyzing the subgraph to identify a first node in the subgraph having a plurality of directed edges as outputs;
wherein the instructions cause the first device to assign, for each of the directed edges, a node to which the directed edge points to a disjoint stream of the hardware accelerator.

7. The method of any preceding claim, wherein the in-

structions cause the first device to determine, for each node, a respective amount of memory resources in the hardware accelerator consumed by the operation represented by the node based on the directed edges to the node, wherein the assigning is based at least on the respective amount of memory resources.

8. The method of any preceding claim, wherein the instructions cause the first device (302) to determine a particular operation represented by a node has finished at a particular stream: in response to determining the particular operation has finished:

determine at least one subsequent operation that uses the output of the particular operation as input; and
reuse memory allocated for the output of the particular operation after the at least one subsequent operation has executed.

9. The method of claim 8, wherein determining at least one subsequent operation that uses the output of the particular operation as input includes:

determining that at least two subsequent operations, a first operation in a first stream and a second operation in a second stream, use the output of the particular operation as input;
placing a first marker in a first stream that indicates when the first operation has used the particular operation as input;
placing a second marker in a second stream that indicates when the second operation has used the particular operation as input;
determining that both operations have used the particular operation upon indication from the first and second markers.

10. The method of any preceding claim, further comprising: receiving a model input; and processing, by the hardware accelerator, the model input according to operations represented by the nodes in the subgraph.

11. A system comprising:

one or more computers; and
computer-readable medium coupled to the one or more computers and having

instructions stored thereon, which, when executed by the one or more computers, cause the one or more computers to perform operations comprising a method as defined in any one of claims 1 to 10.

12. A computer program product encoded on one or more computer storage media, the computer pro-

gram product comprising instructions that when executed by one or more computers cause the one or more computers to perform operations comprising a method as defined in any one of claims 1 to 10.

**Patentansprüche**

1. Computerimplementiertes Verfahren, das Folgendes umfasst:

Empfangen (202) einer Anforderung zum Verarbeiten eines Berechnungsgraphen;
Erhalten (204) von Daten, die mehrere Subgraphen des Berechnungsgraphen darstellen, wobei der Berechnungsgraph eine Vielzahl von Knoten und gerichteten Kanten umfasst, wobei jeder Knoten eine jeweilige Operation darstellt, wobei jede gerichtete Kante einen jeweiligen ersten Knoten mit einem jeweiligen zweiten Knoten verbindet, der eine Operation darstellt, die als Eingabe eine Ausgabe einer durch den jeweiligen ersten Knoten dargestellten Operation empfängt;
Zuordnen der mehreren Subgraphen zu jeweiligen Vorrichtungen aus einer Vielzahl von Vorrichtungen (116, 118, 120, 122), einschließlich des Zuordnens eines ersten Subgraphen (316) zu einer ersten Vorrichtung (302);
Feststellen (206), dass die erste Vorrichtung (302) einen Hardware-Beschleuniger mit einer Vielzahl von Strömen (304, 306) umfasst, wobei jeder Strom eine unabhängige Hardware-Warteschlange ist, deren Operationen der Reihe nach verarbeitet werden;
Erzeugen (208), als Reaktion auf die Feststellung, dass die erste Vorrichtung (302) einen Hardware-Beschleuniger mit einer Vielzahl von Strömen umfasst, von Befehlen, die bei Ausführung durch die erste Vorrichtung die erste Vorrichtung veranlassen zum:

Zuordnen der durch jeden Knoten (308, 310, 312, 314) im ersten Subgraphen (316) dargestellten Operation zu einem entsprechenden Strom in der Vielzahl von Strömen (304, 306) des Hardware-Beschleunigers, einschließlich des Zuordnens von zwei Operationen, die nicht voneinander abhängen, zu verschiedenen Strömen (304, 306) und des Zuordnens einer ersten Operation, die den internen Zustand der Hardware in einer Weise ändert, die vor der Ausführung einer zweiten Operation geschehen muss, zu einem einzelnen Strom (304) und des Zuordnens der in dem einzelnen Strom (304) auszuführenden zweiten Operation nach Abschluss der ersten Operation; und

Durchführen der durch die Knoten (308, 310, 312, 314) im ersten Subgraphen (316) dargestellten Operationen gemäß der Zuordnung; und
Bereitstellen (210) der Befehle und der Daten an die erste Vorrichtung;
wobei die Befehle ferner die erste Vorrichtung (302) veranlassen festzustellen (402), dass eine bestimmte durch einen Knoten dargestellte Operation an einem bestimmten Strom beendet ist;
als Reaktion auf die Feststellung, dass die bestimmte Operation beendet ist:

Bestimmen (404) einer ersten von der bestimmten Operation verbrauchten Speichermenge, die freigegeben wird;
Bestimmen (406), für jeden aus einer Gruppe von unzugeordneten Knoten, einer jeweiligen geschätzten Speichermenge, die durch eine durch den unzugeordneten Knoten dargestellte Operation verbraucht wird;
Bestimmen (408), aus der Gruppe von unzugeordneten Knoten, eines ersten unzugeordneten Knotens, der eine Operation darstellt, die auf einem Strom des Hardware-Beschleunigers ausgeführt wird, mit der geschätzten Speichermenge, die die Nutzung der ersten Speichermenge maximiert; und
Zuordnen (410) einer durch den ersten unzugeordneten Knoten dargestellten Operation zu dem bestimmten Strom.

2. Verfahren nach Anspruch 1, wobei die Anforderung das Identifizieren einer oder mehrerer bestimmter Ausgänge aus einem oder mehreren jeweiligen Knoten (308, 310, 312, 314) im ersten Subgraphen (316) spezifiziert, das ferner Folgendes umfasst:
Empfangen, von der ersten Vorrichtung (302), der ein oder mehreren bestimmten Ausgaben; und Bereitstellen der ein oder mehreren bestimmten Ausgaben dem Client.

3. Verfahren nach Anspruch 1 oder 2, wobei die Befehle ferner die erste Vorrichtung (302) veranlassen, die ein oder mehreren besonderen Ausgaben im Speicher der ersten Vorrichtung zu speichern.

4. Verfahren nach Anspruch 1, 2 oder 3, wobei die Operationen für den ersten Subgraphen (316) Teilinferenz- oder Trainingsberechnungen für ein neuronales Netzwerk umfassen.

5. Verfahren nach einem vorhergehenden Anspruch, das ferner Folgendes umfasst:

Analysieren des ersten Subgraphen (316), um eine Gruppe von Knoten im ersten Subgraphen (316) in einer Kettenstruktur zu identifizieren; wobei die Befehle die erste Vorrichtung (302) (316) veranlassen, die Gruppe von Knoten einem Strom zuzuordnen.

6. Verfahren nach einem vorhergehenden Anspruch, wobei die Zuordnung Folgendes umfasst:

Analysieren des Subgraphen, um einen ersten Knoten in dem Subgraphen mit einer Vielzahl von gerichteten Kanten als Ausgaben zu identifizieren; wobei die Befehle die erste Vorrichtung veranlassen, für jede der gerichteten Kanten einen Knoten, auf den die gerichtete Kante zeigt, einem disjunkten Strom des Hardware-Beschleunigers zuzuordnen.

7. Verfahren nach einem vorhergehenden Anspruch, wobei die Befehle die erste Vorrichtung veranlassen, für jeden Knoten auf der Basis der gerichteten Kanten zu dem Knoten eine jeweilige Menge an Speicherressourcen in dem Hardware-Beschleuniger zu bestimmen, die von der durch den Knoten dargestellten Operation verbraucht werden, wobei die Zuordnung mindestens auf der jeweiligen Menge an Speicherressourcen basiert.

8. Verfahren nach einem vorhergehenden Anspruch, wobei die Befehle die erste Vorrichtung (302) veranlassen festzustellen, dass eine bestimmte durch einen Knoten dargestellte Operation bei einem bestimmten Strom beendet ist: als Reaktion auf die Feststellung, dass die bestimmte Operation beendet ist:

Bestimmen mindestens einer nachfolgenden Operation, die die Ausgabe der bestimmten Operation als Eingabe verwendet; und Wiederverwenden von für die Ausgabe der bestimmten Operation zugeordneten Speichers nach der Ausführung der mindestens einen nachfolgenden Operation.

9. Verfahren nach Anspruch 8, wobei das Bestimmen mindestens einer nachfolgenden Operation, die die Ausgabe der bestimmten Operation als Eingabe verwendet, Folgendes umfasst:

Feststellen, dass mindestens zwei nachfolgende Operationen, eine erste Operation in einem ersten Strom und eine zweite Operation in einem zweiten Strom, die Ausgabe der bestimmten Operation als Eingabe verwenden; Platzieren einer ersten Markierung in einem ersten Strom, die anzeigt, wenn die erste Operation

die bestimmte Operation als Eingabe verwendet hat; Platzieren einer zweiten Markierung in einem zweiten Strom, die anzeigt, wann die zweite Operation die bestimmte Operation als Eingabe verwendet hat; Feststellen, dass beide Operationen die bestimmte Operation verwendet haben, wenn die erste und die zweite Markierung dies anzeigen.

10. Verfahren nach einem vorhergehenden Anspruch, das ferner Folgendes umfasst: Empfangen einer Modelleingabe; und Verarbeiten, durch den Hardware-Beschleuniger, der Modelleingabe entsprechend den durch die Knoten in dem Subgraphen dargestellten Operationen.

11. System, das Folgendes umfasst:

einen oder mehrere Computer; und ein computerlesbares Medium, das mit den ein oder mehreren Computern gekoppelt ist und auf dem Befehle gespeichert sind, die bei Ausführung durch die ein oder mehreren Computer die ein oder mehreren Computer veranlassen, Operationen durchzuführen, die ein Verfahren nach einem der Ansprüche 1 bis 10 umfassen.

12. Computerprogrammprodukt, das auf einem oder mehreren Computerspeichermedien codiert ist, wobei das Computerprogrammprodukt Befehle umfasst, die bei Ausführung durch einen oder mehrere Computer die ein oder mehreren Computer veranlassen, Operationen durchzuführen, die ein Verfahren nach einem der Ansprüche 1 bis 10 umfassen.

**Revendications**

1. Procédé mis en œuvre par ordinateur comprenant :

la réception (202) d'une demande de traitement d'un graphe de calcul ; l'obtention (204) de données représentant plusieurs sous-graphes du graphe de calcul, le graphe de calcul comprenant une pluralité de nœuds et d'arêtes dirigées, dans lequel chaque nœud représente une opération respective, dans lequel chaque arête dirigée connecte un premier nœud respectif à un second nœud respectif qui représente une opération qui reçoit, en tant qu'entrée, une sortie d'une opération représentée par le premier nœud respectif ; l'attribution des plusieurs sous-graphes à des dispositifs respectifs d'une pluralité de dispositifs (116, 118, 120, 122), incluant l'attribution d'un premier sous-graphe (316) à un premier dispositif (302) ;

la détermination (206) que le premier dispositif (302) comprend un accélérateur matériel ayant une pluralité de flux (304, 306), chaque flux étant une file d'attente matérielle indépendante dont les opérations sont traitées dans l'ordre ;
en réponse à la détermination que le premier dispositif (302) comprend un accélérateur matériel ayant une pluralité de flux, la génération (208) d'instructions qui, lorsqu'elles sont exécutées par le premier dispositif, amènent le premier dispositif à :

> affecter l'opération représentée par chaque nœud (308, 310, 312, 314) dans le premier sous-graphe (316) à un flux correspondant dans la pluralité de flux (304, 306) de l'accélérateur matériel, incluant l'affectation de deux opérations qui ne dépendent pas l'une de l'autre à des flux (304, 306) différents, et l'affectation à un flux unique (304) d'une première opération qui modifie l'état interne du matériel d'une manière qui doit se produire avant qu'une seconde opération ne s'exécute et l'affectation de la seconde opération à exécuter dans le flux unique (304) après que la première opération soit terminée ; et
> exécuter les opérations représentées par les nœuds (308, 310, 312, 314) dans le premier sous-graphe (316) conformément à l'affectation ; et
> fournir (210) les instructions et les données au premier dispositif ;
> dans lequel les instructions amènent en outre le premier dispositif (302) à déterminer (402) qu'une opération particulière représentée par un nœud est terminée à un flux particulier ;
> en réponse à la détermination que l'opération particulière est terminée :

>> déterminer (404) une première quantité de mémoire consommée par l'opération particulière qui sera libérée ;
>> déterminer (406), pour chacun d'un groupe de nœuds non affectés, une quantité estimée respective de mémoire consommée par une opération qui est représentée par le nœud non affecté ;
>> déterminer (408), à partir du groupe de nœuds non affectés, un premier nœud non affecté qui représente une opération, qui s'exécute sur un flux de l'accélérateur matériel, avec la quantité estimée de mémoire qui maximise l'utilisation de la première quantité de mémoire ; et

>> affecter (410) une opération représentée par le premier nœud non affecté au flux particulier.

2. Procédé selon la revendication 1, dans lequel la demande spécifie l'identification d'une ou plusieurs sorties particulières à partir d'un ou plusieurs nœuds (308, 310, 312, 314) respectifs dans le premier sous-graphe (316), comprenant en outre : la réception, à partir du premier dispositif (302), des une ou plusieurs sorties particulières ; et la fourniture des une ou plusieurs sorties particulières au client.

3. Procédé selon la revendication 1 ou 2, dans lequel les instructions amènent en outre le premier dispositif (302) à stocker les une ou plusieurs sorties particulières dans la mémoire du premier dispositif.

4. Procédé selon la revendication 1, 2 ou 3, dans lequel les opérations pour le premier sous-graphe (316) comprennent des calculs d'inférence partielle ou d'apprentissage pour un réseau neuronal.

5. Procédé selon l'une quelconque des revendications précédentes, comprenant en outre :

> l'analyse du premier sous-graphe (316) pour identifier un groupe de nœuds dans le premier sous-graphe (316) dans une structure de chaîne ;
> dans lequel les instructions amènent le premier dispositif (302) (316) à affecter le groupe de nœuds à un flux.

6. Procédé selon l'une quelconque des revendications précédentes, dans lequel l'affectation comprend :

> l'analyse du sous-graphe pour identifier un premier nœud dans le sous-graphe ayant une pluralité d'arêtes dirigées en tant que sorties ;
> dans lequel les instructions amènent le premier dispositif à affecter, pour chacune des arêtes dirigées, un nœud vers lequel l'arête dirigée pointe à un flux disjoint de l'accélérateur matériel.

7. Procédé selon l'une quelconque des revendications précédentes, dans lequel les instructions amènent le premier dispositif à déterminer, pour chaque nœud, une quantité respective de ressources de mémoire dans l'accélérateur matériel consommées par l'opération représentée par le nœud sur la base des arêtes dirigées vers le nœud, dans lequel l'attribution est basée au moins sur la quantité respective de ressources de mémoire.

8. Procédé selon l'une quelconque des revendications précédentes, dans lequel les instructions amènent

le premier dispositif (302) à déterminer qu'une opération particulière représentée par un nœud s'est terminée à un flux particulier : en réponse à la détermination que l'opération particulière s'est terminée :

déterminer au moins une opération ultérieure qui utilise la sortie de l'opération particulière en tant qu'entrée ; et
réutiliser la mémoire allouée à la sortie de l'opération particulière après l'exécution de l'au moins une opération ultérieure.

9.  Procédé selon la revendication 8, dans lequel la détermination d'au moins une opération ultérieure qui utilise la sortie de l'opération particulière en tant qu'entrée comprend :

la détermination qu'au moins deux opérations ultérieures, une première opération dans un premier flux et une seconde opération dans un second flux, utilisent la sortie de l'opération particulière en tant qu'entrée ;
le placement d'un premier marqueur dans un premier flux qui indique quand la première opération a utilisé l'opération particulière en tant qu'entrée ;
le placement d'un second marqueur dans un second flux qui indique quand la seconde opération a utilisé l'opération particulière en tant qu'entrée ;
la détermination que les deux opérations ont utilisé l'opération particulière sur indication des premier et second marqueurs.

10.  Procédé selon l'une quelconque des revendications précédentes, comprenant en outre : la réception d'une entrée de modèle ; et le traitement, par l'accélérateur matériel, de l'entrée de modèle selon des opérations représentées par les nœuds dans le sous-graphe.

11.  Système comprenant :

un ou plusieurs ordinateurs ; et
un support lisible par ordinateur couplé aux un ou plusieurs ordinateurs et comportant des instructions stockées sur celui-ci, qui, lorsqu'elles sont exécutées par les un ou plusieurs ordinateurs, amènent les un ou plusieurs ordinateurs à effectuer des opérations comprenant un procédé tel que défini selon l'une quelconque des revendications 1 à 10.

12.  Progiciel informatique codé sur un ou plusieurs supports de stockage informatique, le progiciel informatique comprenant des instructions qui, lorsqu'elles sont exécutées par un ou plusieurs ordinateurs, amènent les un ou plusieurs ordinateurs à effectuer

des opérations comprenant un procédé tel que défini selon l'une quelconque des revendications 1 à 10.

FIG. 1

```
┌─────────────────────────┐
│   Receive a request from a  │
│    client to process a      │
│   computational graph       │
│         202                 │
└─────────────────────────┘
              │
              ▼
┌─────────────────────────┐
│  Obtain data representing a │
│       subgraph of the       │
│   computational graph       │
│         204                 │
└─────────────────────────┘
              │
              ▼
┌─────────────────────────┐
│  Determine a first device is a │
│  graphical processing unit  │
│         206                 │
└─────────────────────────┘
              │
              ▼
┌─────────────────────────┐
│                             │
│  Generate instructions for the │
│        first device         │
│         208                 │
│                             │
└─────────────────────────┘
              │
              ▼
┌─────────────────────────┐
│   Provide the instructions and │
│  the data to the first device  │
│         210                 │
└─────────────────────────┘
```

FIG. 2

FIG. 3

```
┌─────────────────────────┐
│  Determine a particular  │
│ operation has finished at a │
│     particular stream     │
│           402            │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│   Determine an amount of   │
│  memory consumed that will │
│         be freed          │
│           404            │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│  Determine, for each of a  │
│  group of unassigned nodes, │
│   a respective estimated   │
│     amount of memory      │
│      consumed by the      │
│      unassigned node      │
│           406            │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│ Determine a first unassigned │
│   node with the estimated   │
│    amount of memory that   │
│ maximizes usage of the first │
│      amount of memory      │
│           408            │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│    Assign an operation    │
│   represented by the first  │
│    unassigned node to the  │
│     particular stream     │
│           410            │
└─────────────────────────┘
```
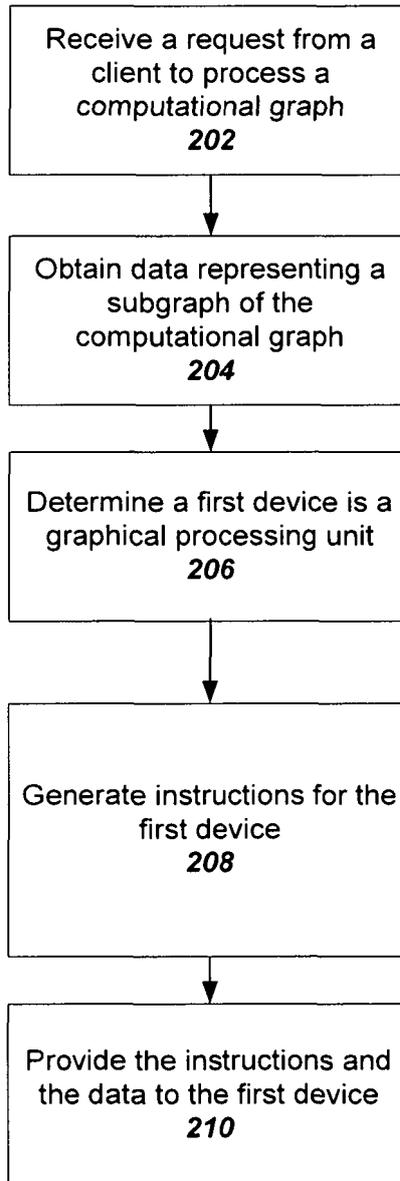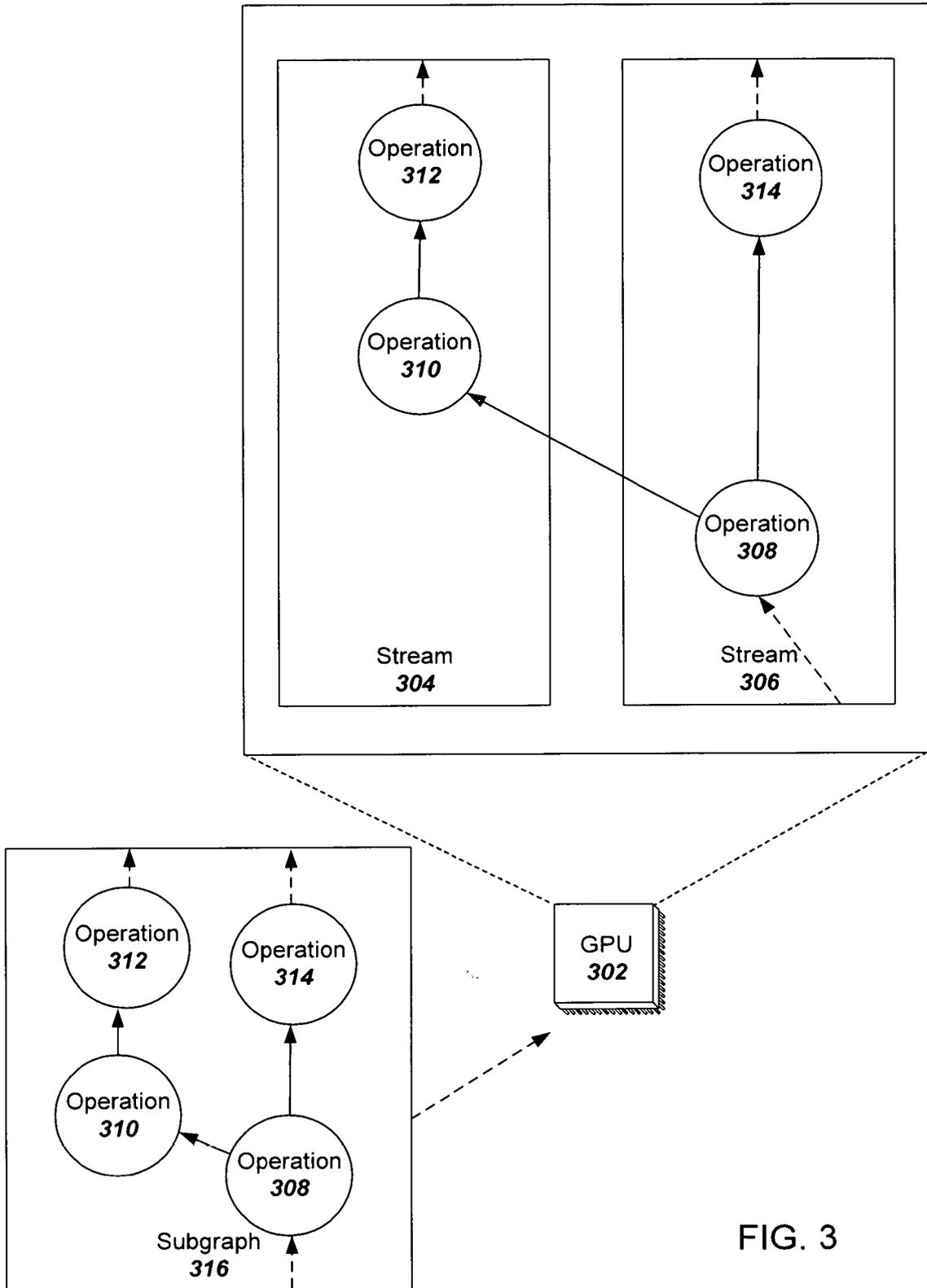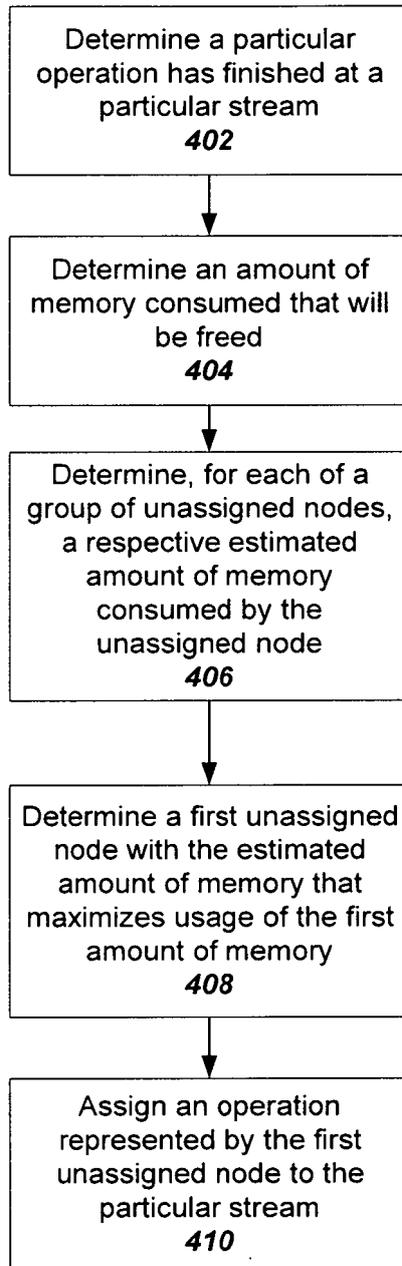
FIG. 4

## REFERENCES CITED IN THE DESCRIPTION

*This list of references cited by the applicant is for the reader's convenience only. It does not form part of the European patent document. Even though great care has been taken in compiling the references, errors or omissions cannot be excluded and the EPO disclaims all liability in this regard.*

### Patent documents cited in the description

- US 2006095722 A1 **[0004]**

- US 2015007182 A1 **[0005]**