

What is claimed:

1. A method for managing a plurality of instances of a Wide Area Network (WAN) optimizer executing on an intermediary device, the method comprising:
 - 5 (a) establishing, on a device intermediary to a plurality of clients and a plurality of servers, a plurality of instances of a Wide Area Network (WAN) optimizer to accelerate WAN communications between the plurality of clients and the plurality of servers;
 - (b) monitoring, by the device, network traffic traversing the device for each of
10 the plurality of instances of the WAN optimizer; and
 - (c) selecting, by a manager executing on the device responsive to the monitoring, a change of a load balancing scheme to load balance the plurality of instances of the WAN optimizer.
- 15 2. The method of claim 1, wherein step (a) further comprises automatically establishing, by the device, a configuration of a size of each of the plurality of instances of the WAN optimizer based on data stored from monitoring of previous execution of the plurality of instances of the WAN optimizer.
- 20 3. The method of claim 1, wherein step (a) further comprises executing, by the device, each of the plurality of instances of the WAN optimizer as a virtual machine in a virtualized environment.
4. The method of claim 1, wherein step (b) further comprises monitoring, by the device, compression history allocation, compression fragmentation and compression ratios of each of the plurality of instances of the WAN optimizer.
- 25 5. The method of claim 1, wherein step (b) further comprises monitoring, by the device, or more of the following of each of the plurality of instances of the WAN optimizer: resource utilization, number of connections, number of claims and bandwidth usage.
6. The method of claim 1, wherein step (c) further comprises determining, by the device, that a metric computed from monitoring network traffic has exceeded a threshold and
30 responsive to the determination, automatically selecting by the device a second load balancing scheme to load balance the plurality of instances of the WAN optimizer.
7. The method of claim 1, wherein step (c) further comprises automatically switching, by the device, from the load balancing scheme to the selected load balancing scheme while executing the plurality of instances of the WAN optimizer.

8. The method of claim 1, further comprising automatically changing, by the device responsive to the monitoring, the number of instances of the WAN optimizer executing on the device.
- 5 9. The method of claim 1, further comprising automatically adjusting, by the device responsive to the monitoring, a size of resource usage used by one or more of the plurality of instances of the WAN optimizer.
- 10 10. The method of claim 1, further comprises applying, by the device one or more rules to data collected from monitoring, to determine to change one or more of the following: a number of instances of the WAN optimizer, a size of one or more WAN optimizers and the load balancing scheme.
11. A system for managing a plurality of instances of a Wide Area Network (WAN) optimizer executing on an intermediary device, the system comprising:
- a device intermediary to a plurality of clients and a plurality of servers;
 - a plurality of instances of a Wide Area Network (WAN) optimizer executing on the device to accelerate WAN communications between the plurality of clients and the plurality of servers;
 - a monitor that monitors network traffic traversing the device for each of the plurality of instances of the WAN optimizer; and
 - a manager executing on the device that, responsive to the monitor, selects a change of a load balancing scheme to load balance the plurality of instances of the WAN optimizer.
12. The system of claim 11, wherein the manager automatically establishes a configuration of a size of each of the plurality of instances of the WAN optimizer based on data stored from monitoring of previous execution of the plurality of instances of the WAN optimizer.
13. The system of claim 11, wherein each of the plurality of instances of the WAN optimizer execute as a virtual machine in a virtualized environment.
14. The system of claim 11, wherein the monitor monitors compression history allocation, compression fragmentation and compression ratios of each of the plurality of instances of the WAN optimizer.
15. The system of claim 11, wherein the monitor monitors one or more of the following of each of the plurality of instances of the WAN optimizer: resource utilization, number of connections, number of claims and bandwidth usage.

16. The system of claim 11, wherein the manager determines that a metric computed from monitoring network traffic has exceeded a threshold and responsive to the determination, automatically selects a second load balancing scheme to load balance the plurality of instances of the WAN optimizer.
- 5 17. The system of claim 11, wherein the manager automatically switches from a current load balancing scheme to the selected load balancing scheme while executing the plurality of instances of the WAN optimizer.
18. The system of claim 11, wherein the manager, responsive to the monitor, automatically changes the number of instances of the WAN optimizer executing on
10 the device.
19. The system of claim 11, wherein the manager automatically adjusts, responsive to the monitor, a size of resource usage used by one or more of the plurality of instances of the WAN optimizer.
20. The system of claim 11, wherein the manager applies one or more rules to data
15 collected from monitoring, to determine to change one or more of the following: a number of instances of the WAN optimizer, a size of one or more WAN optimizers and the load balancing scheme.

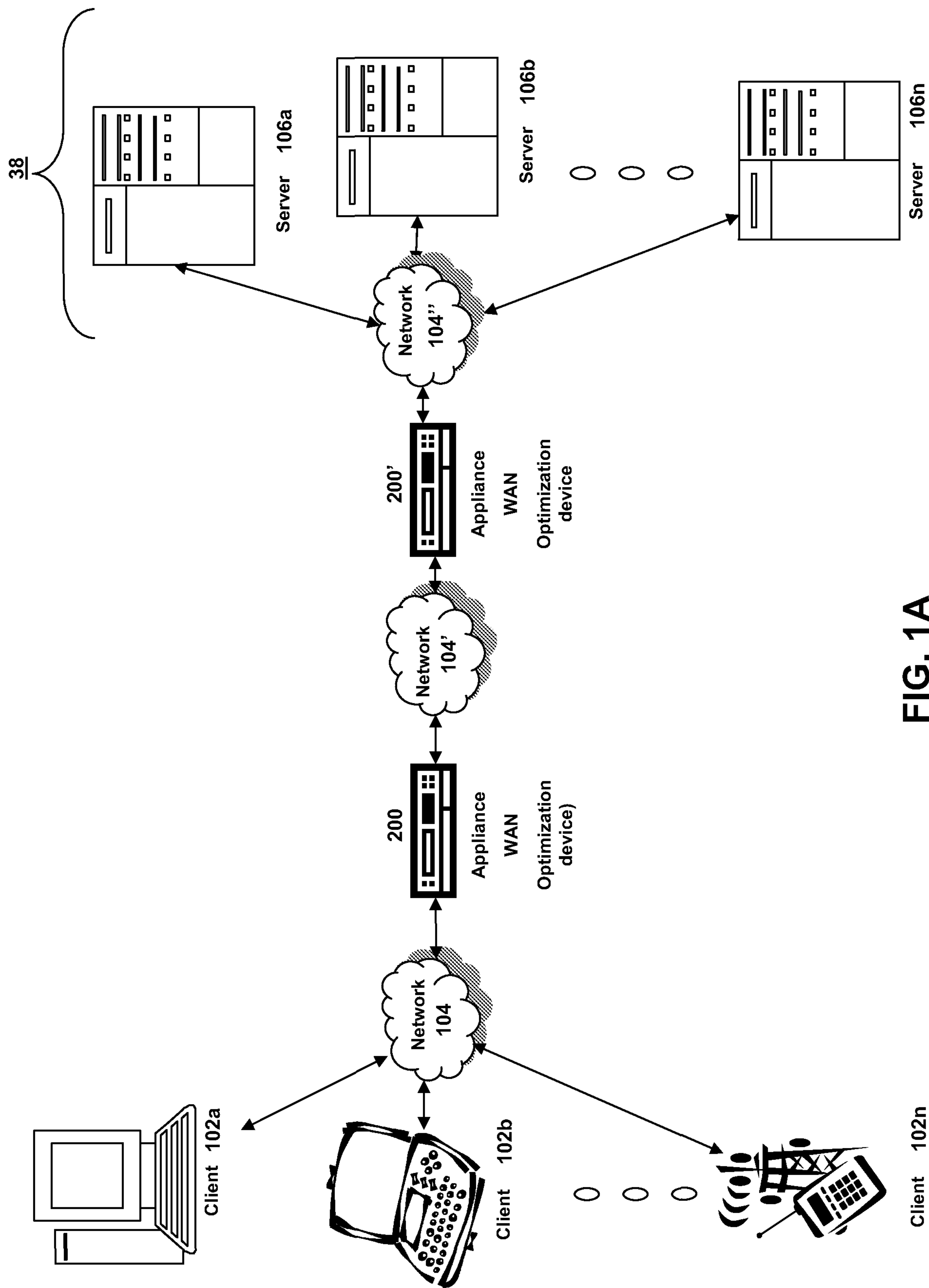


FIG. 1A

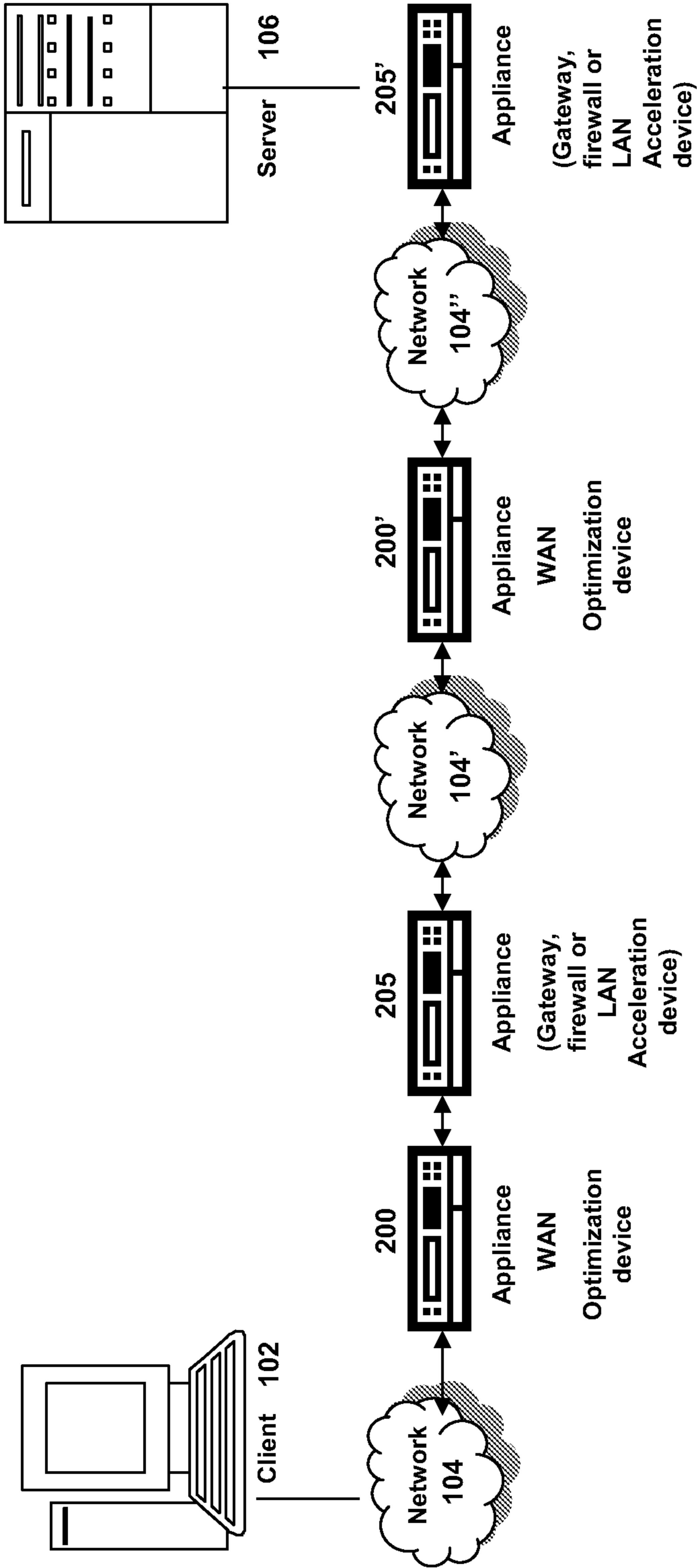


FIG. 1B

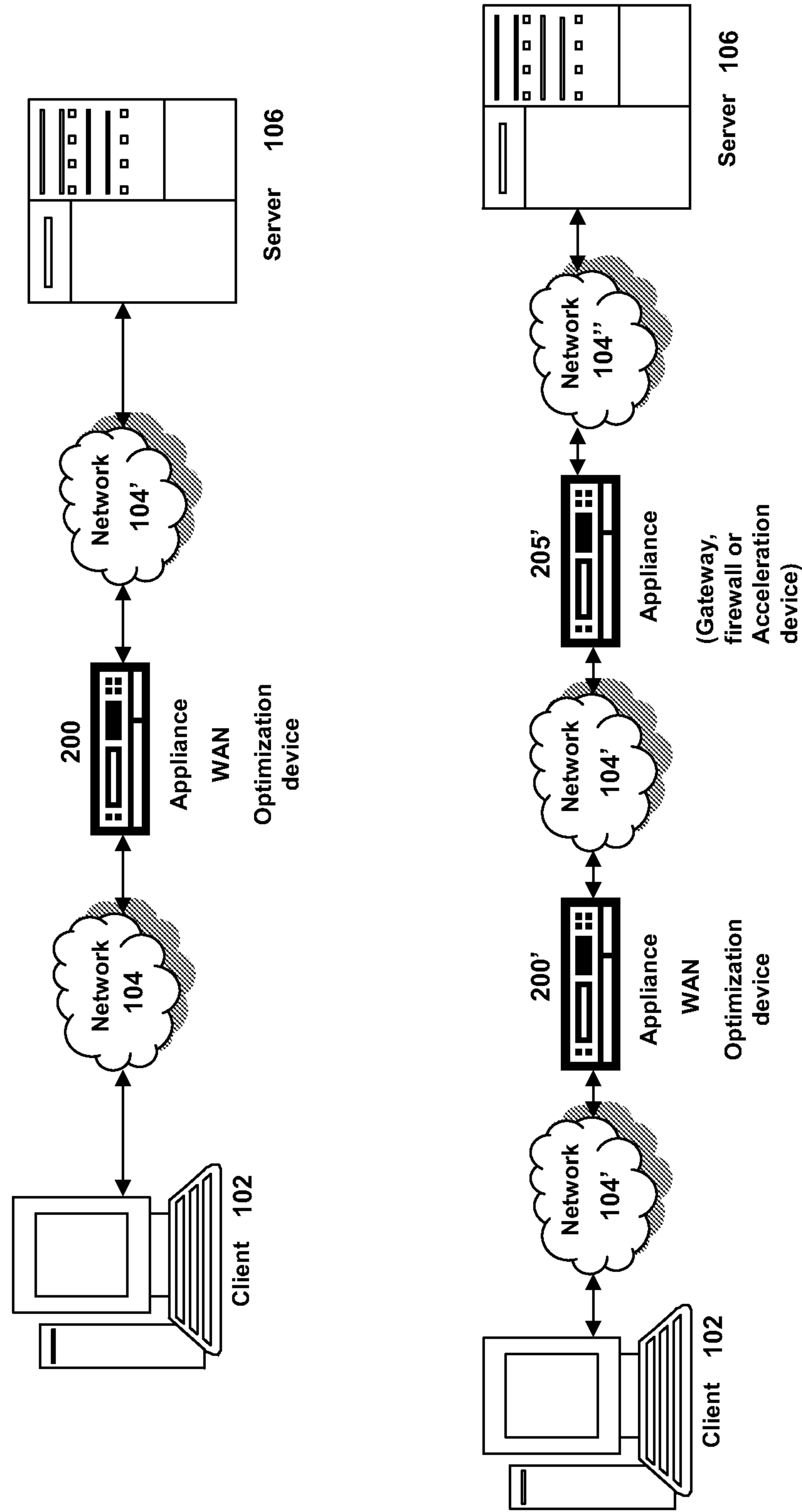


FIG. 1C

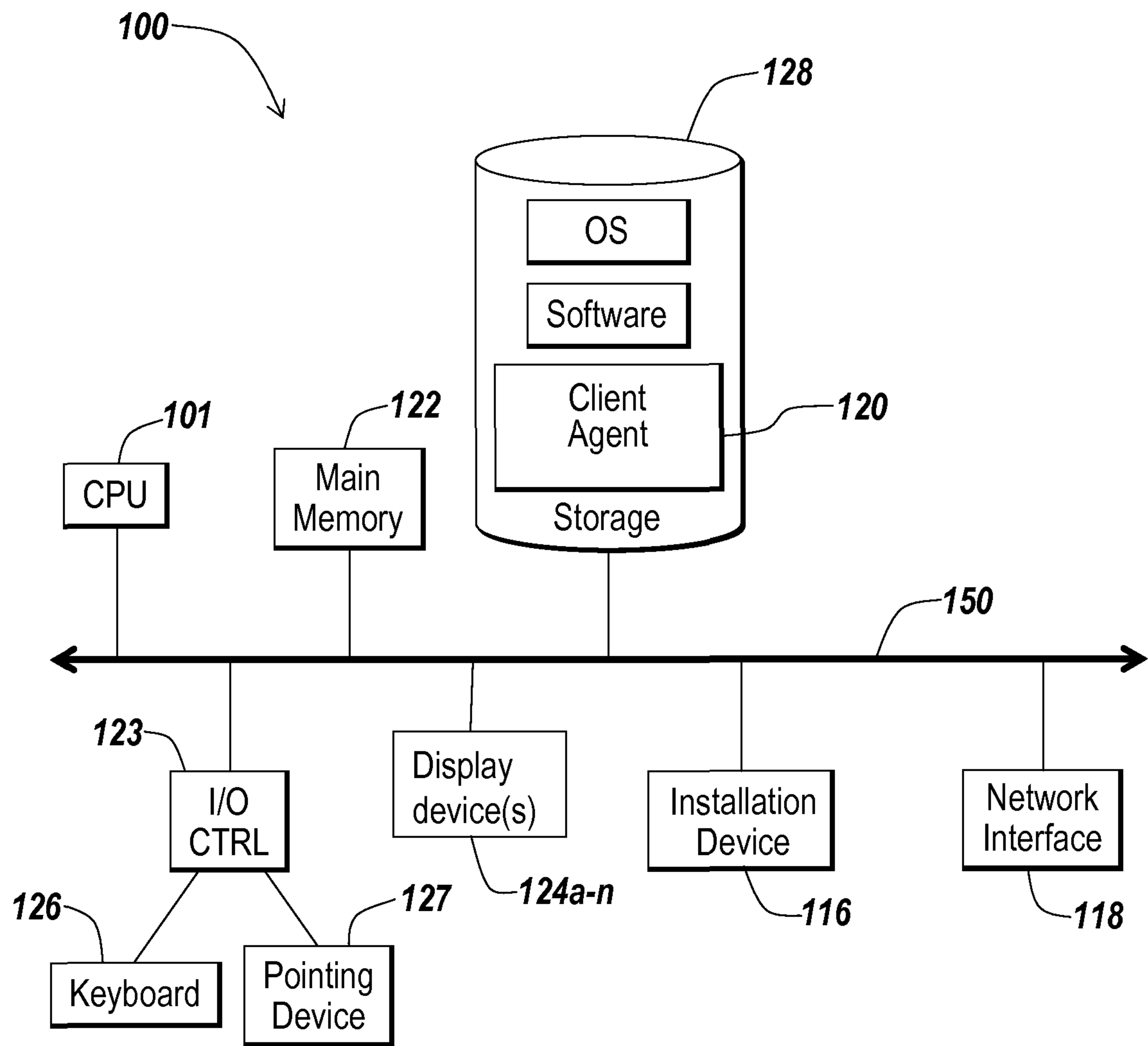
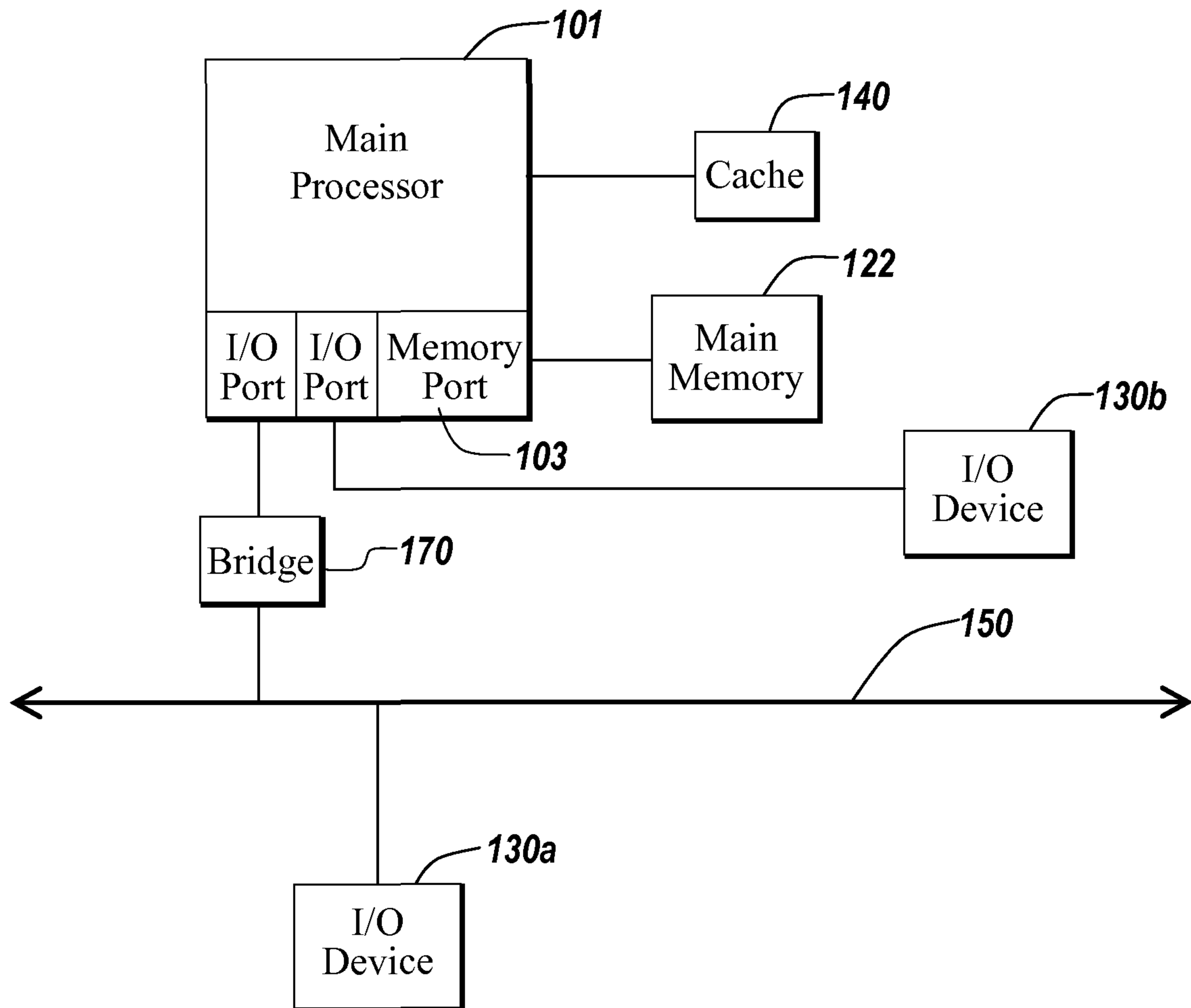


Fig. 1D

*Fig. 1E*

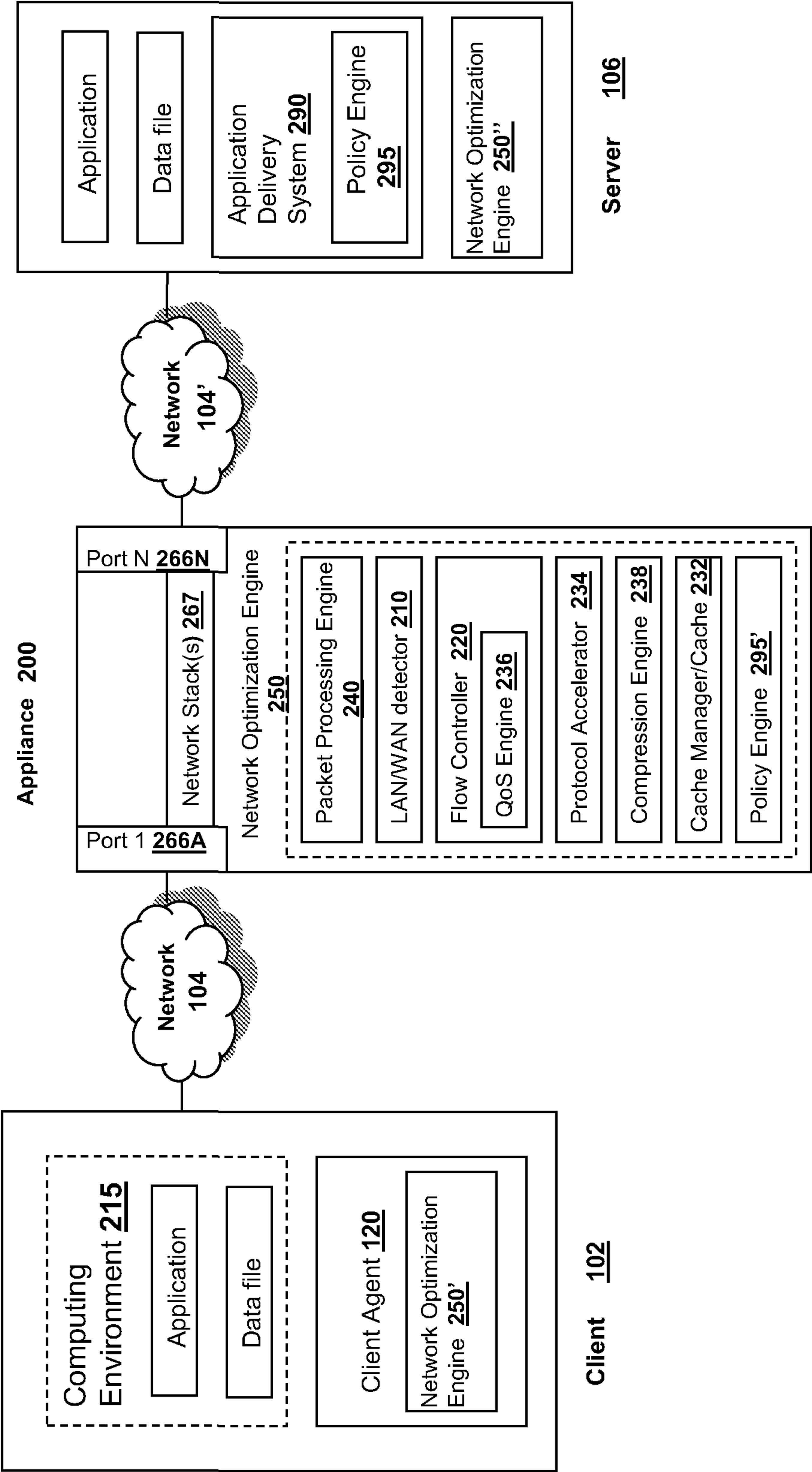


FIG. 2A

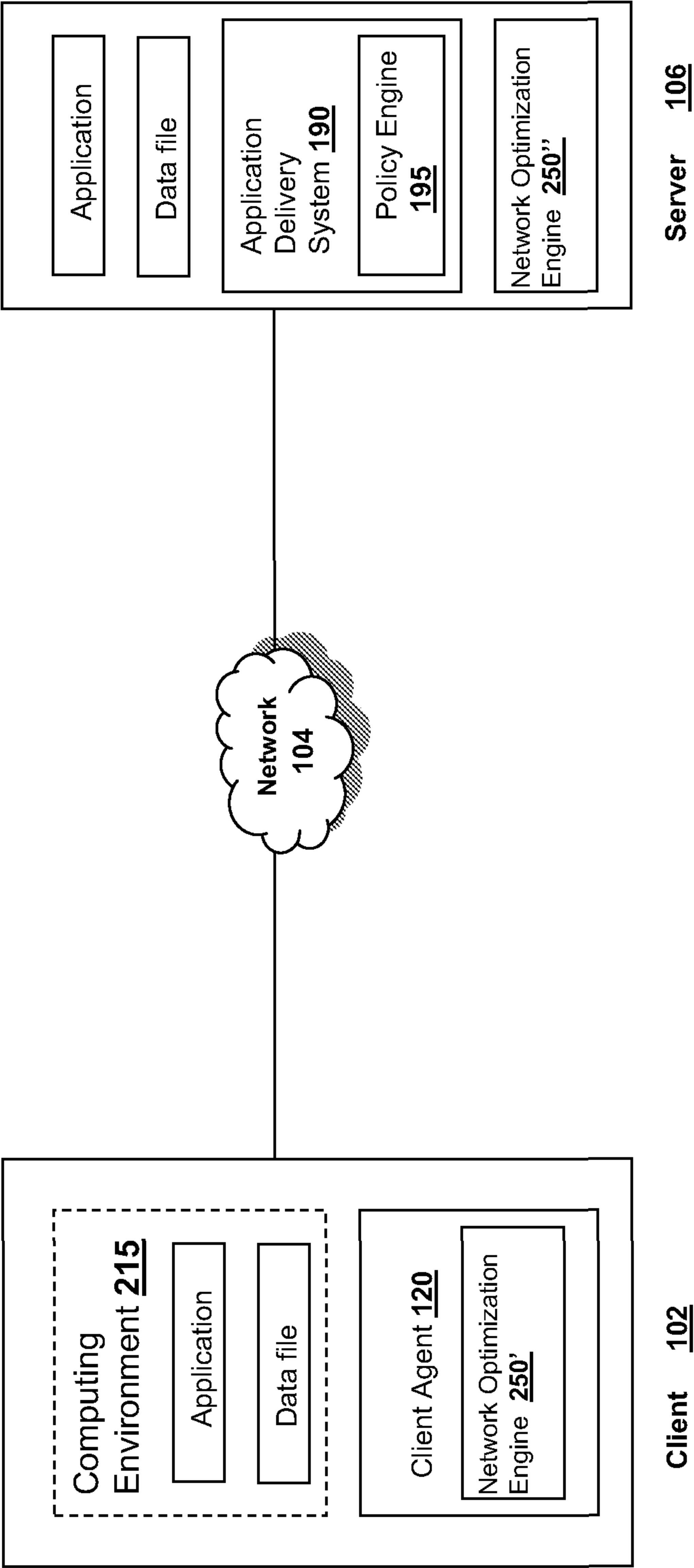


FIG. 2B

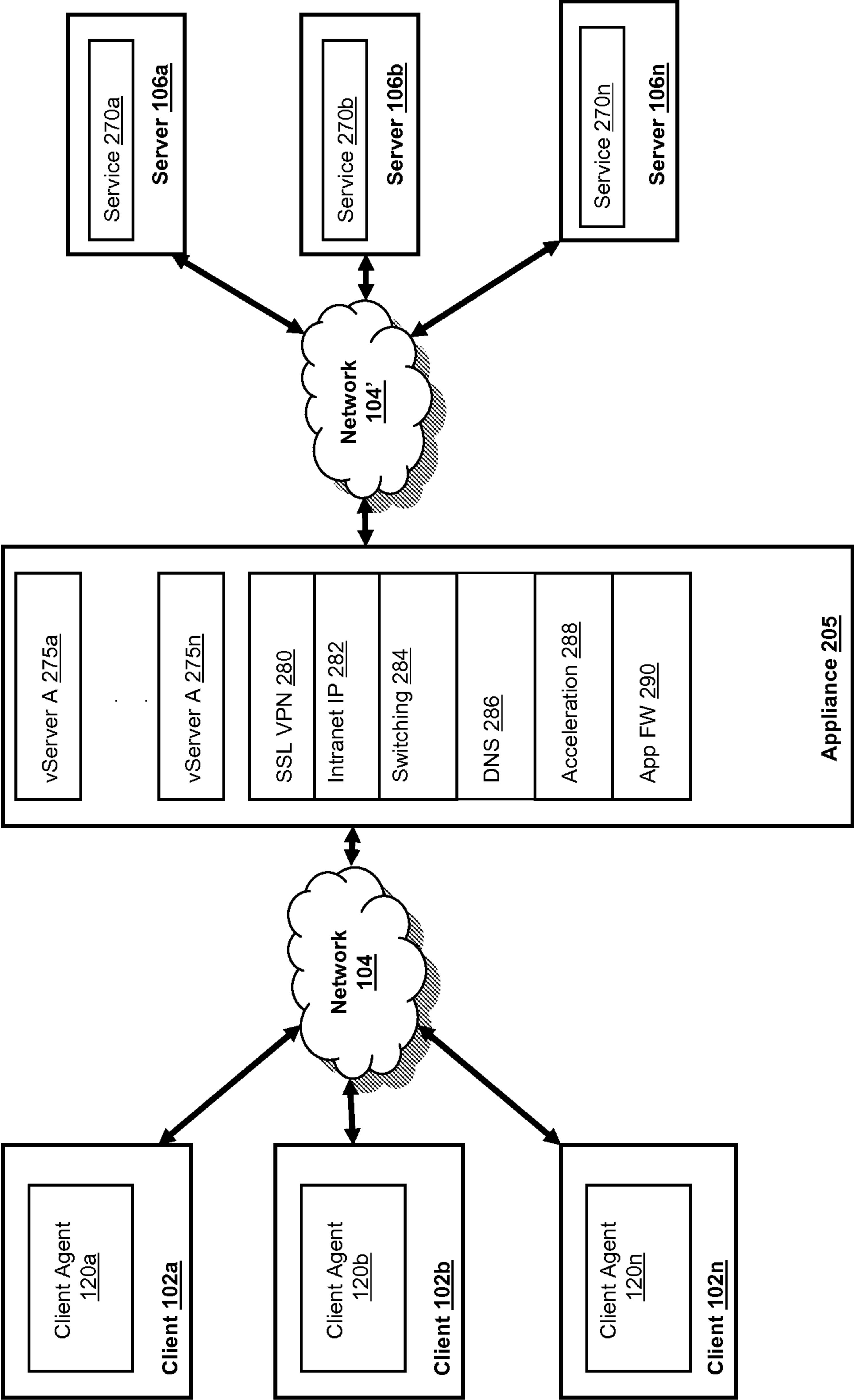
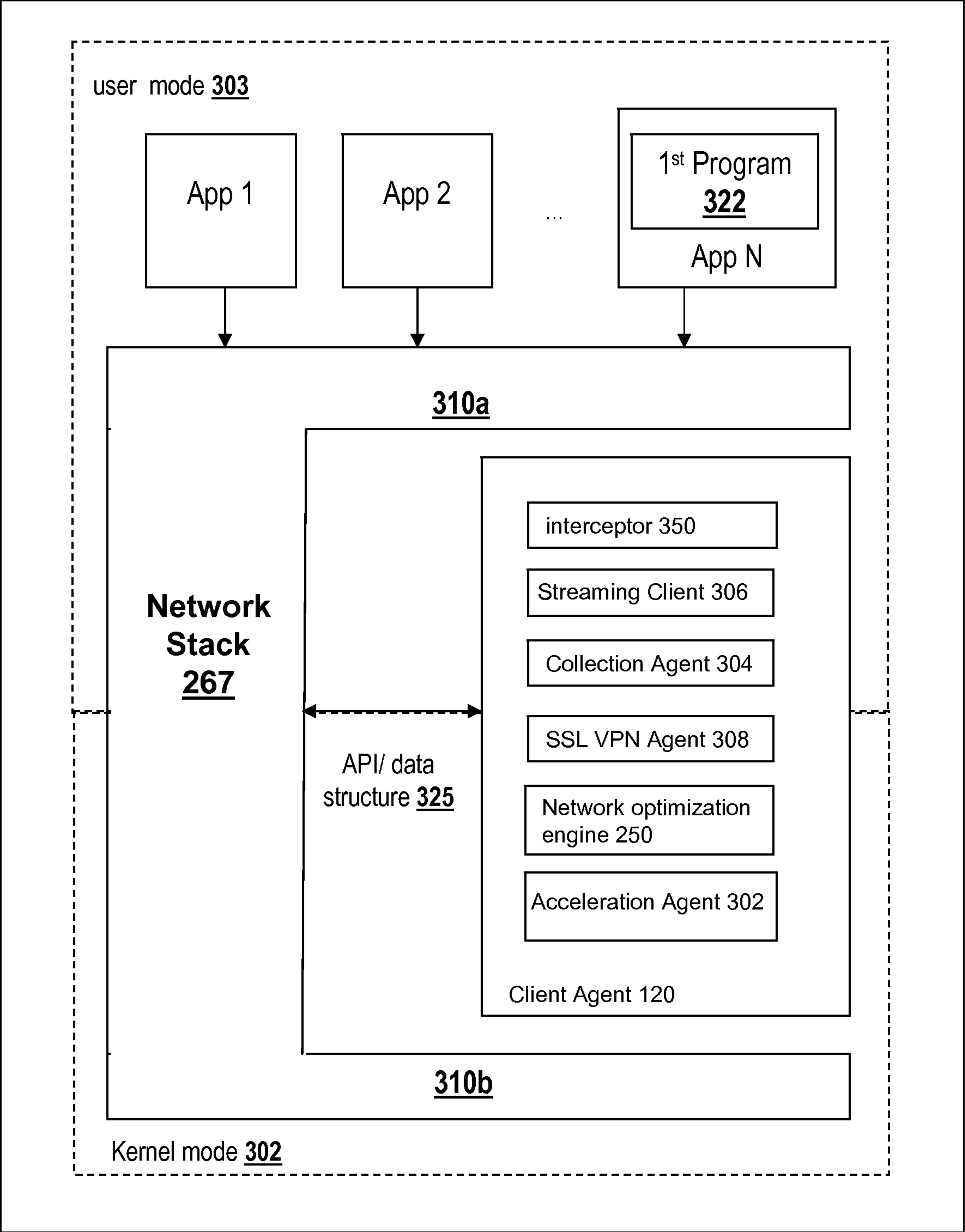


FIG. 2C

Client 102



100

Fig. 3

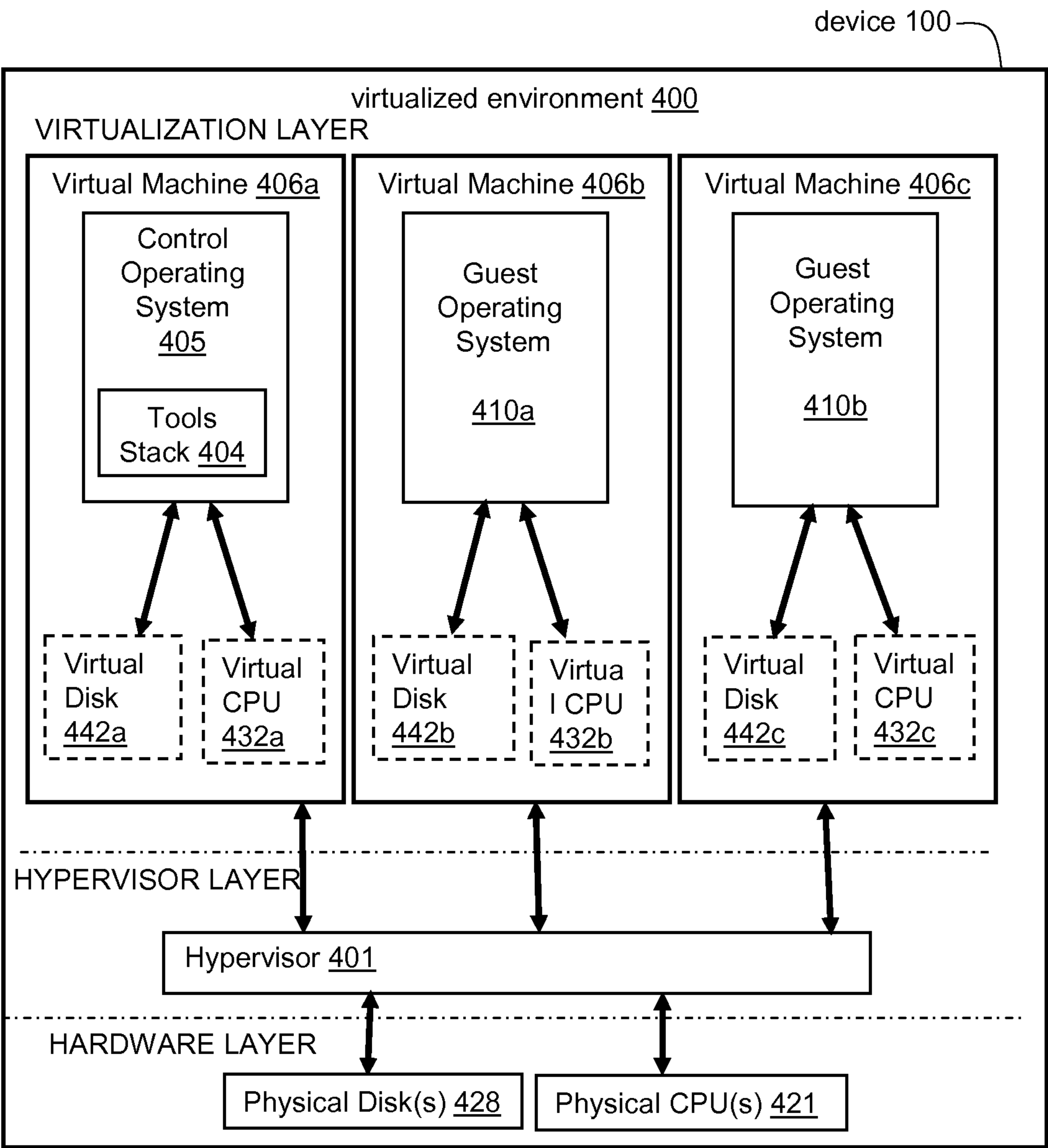
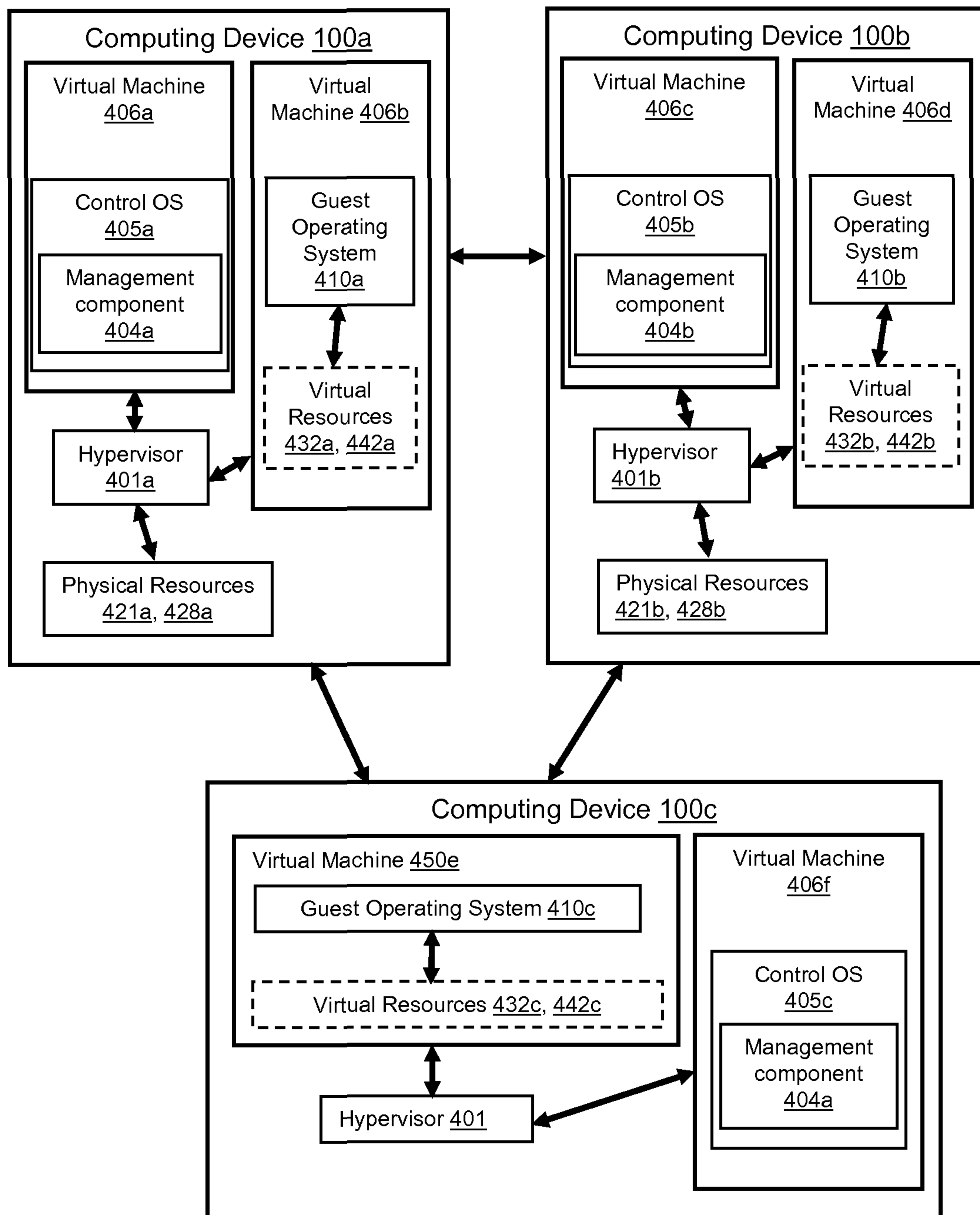
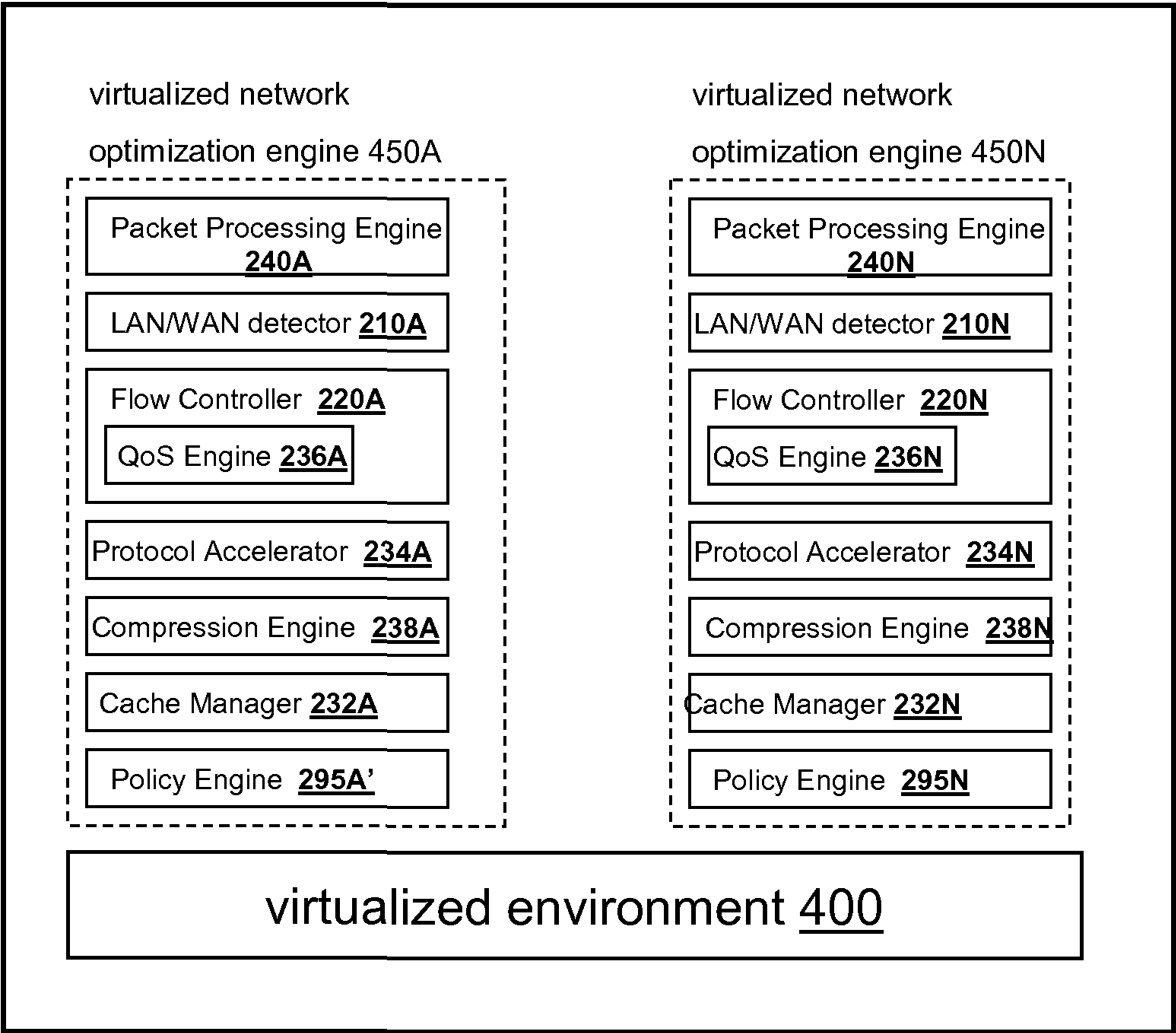


FIG. 4A

**FIG. 4B**

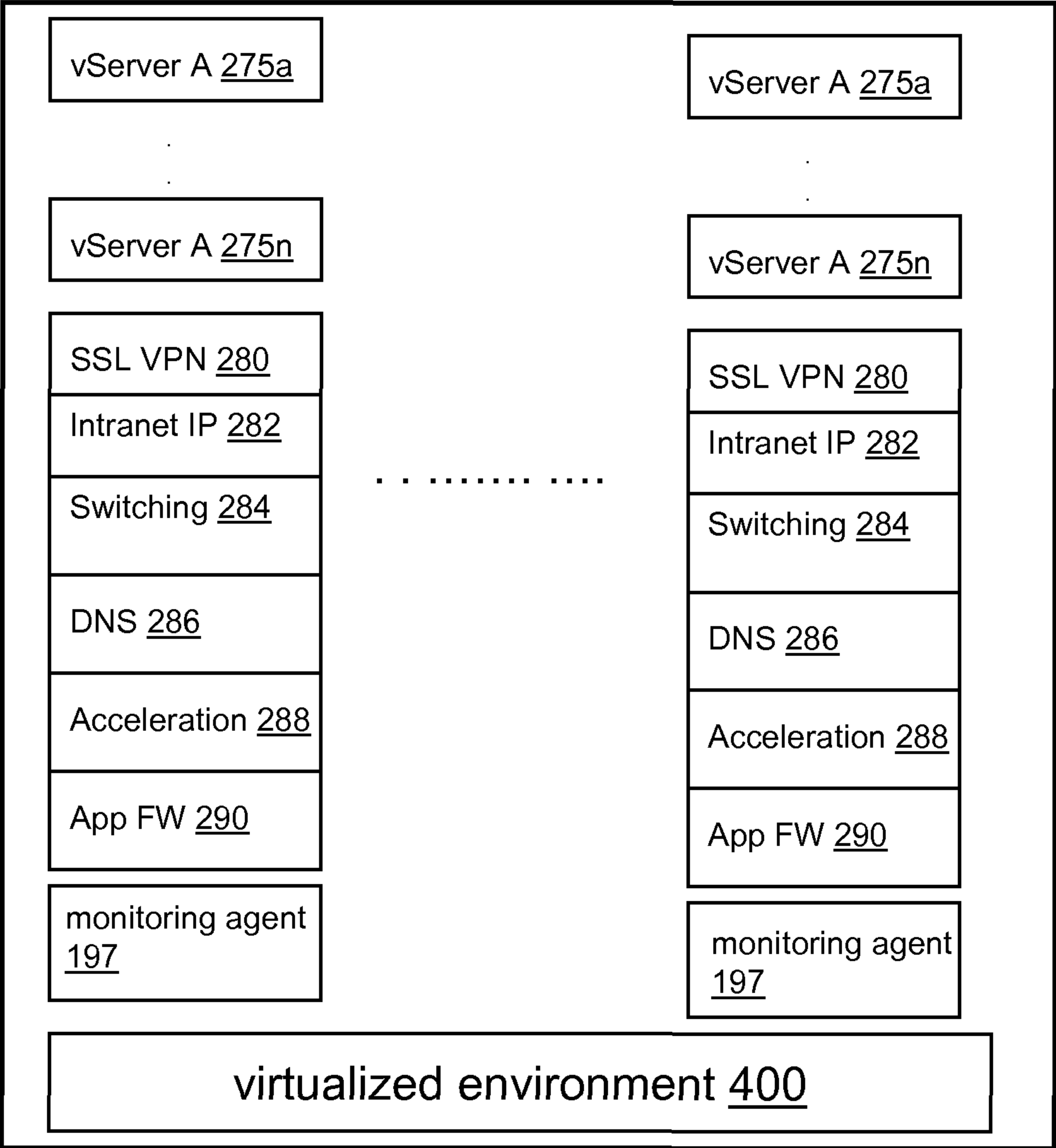
Virtualized appliance 450



device 100

FIG. 4C

virtualized application delivery controller 460



computing device 100

FIG. 4D

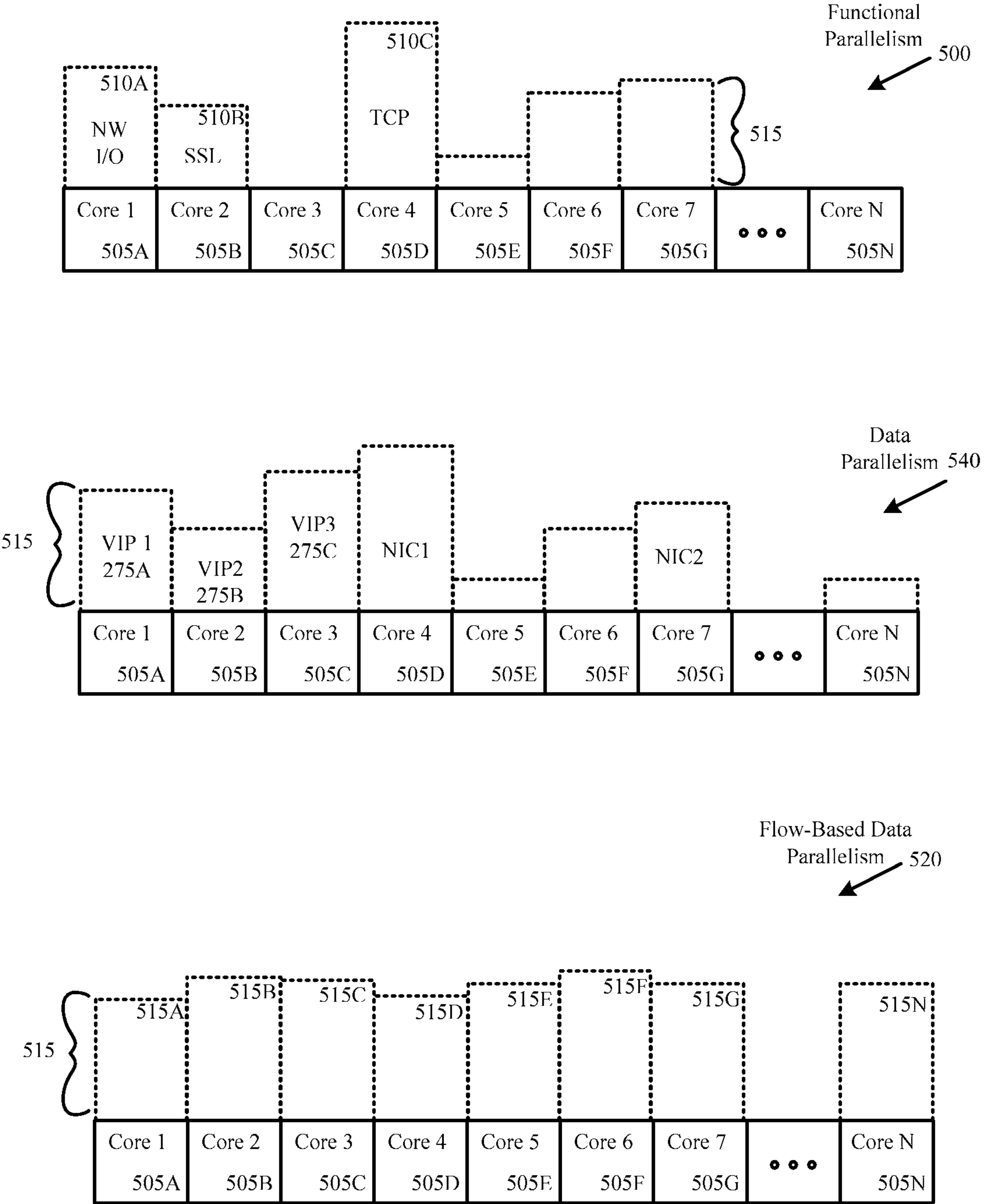


FIG. 5A

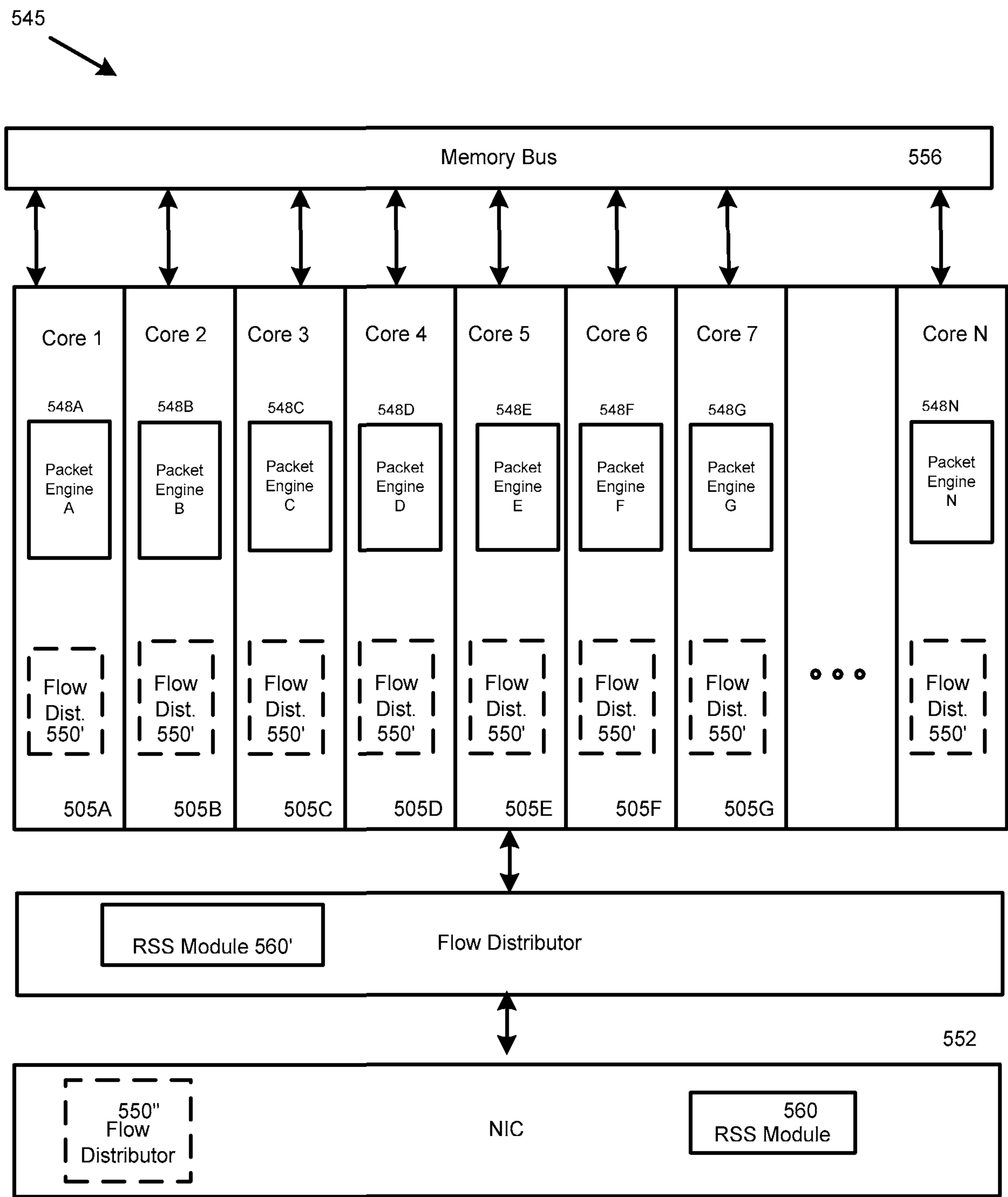


FIG. 5B

575
↘

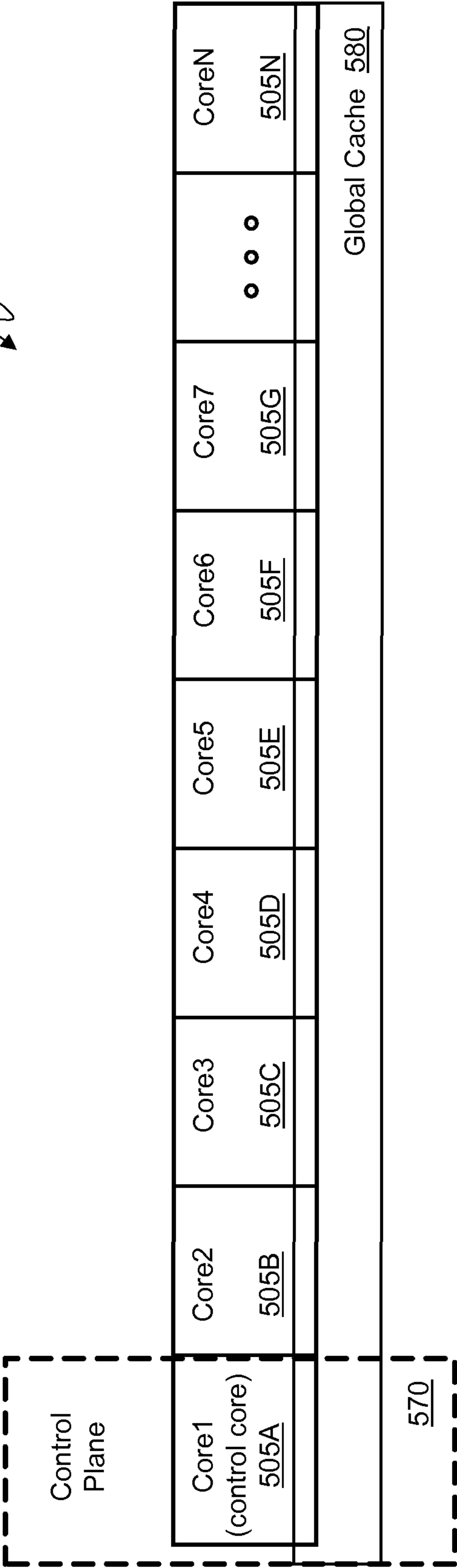


FIG. 5C

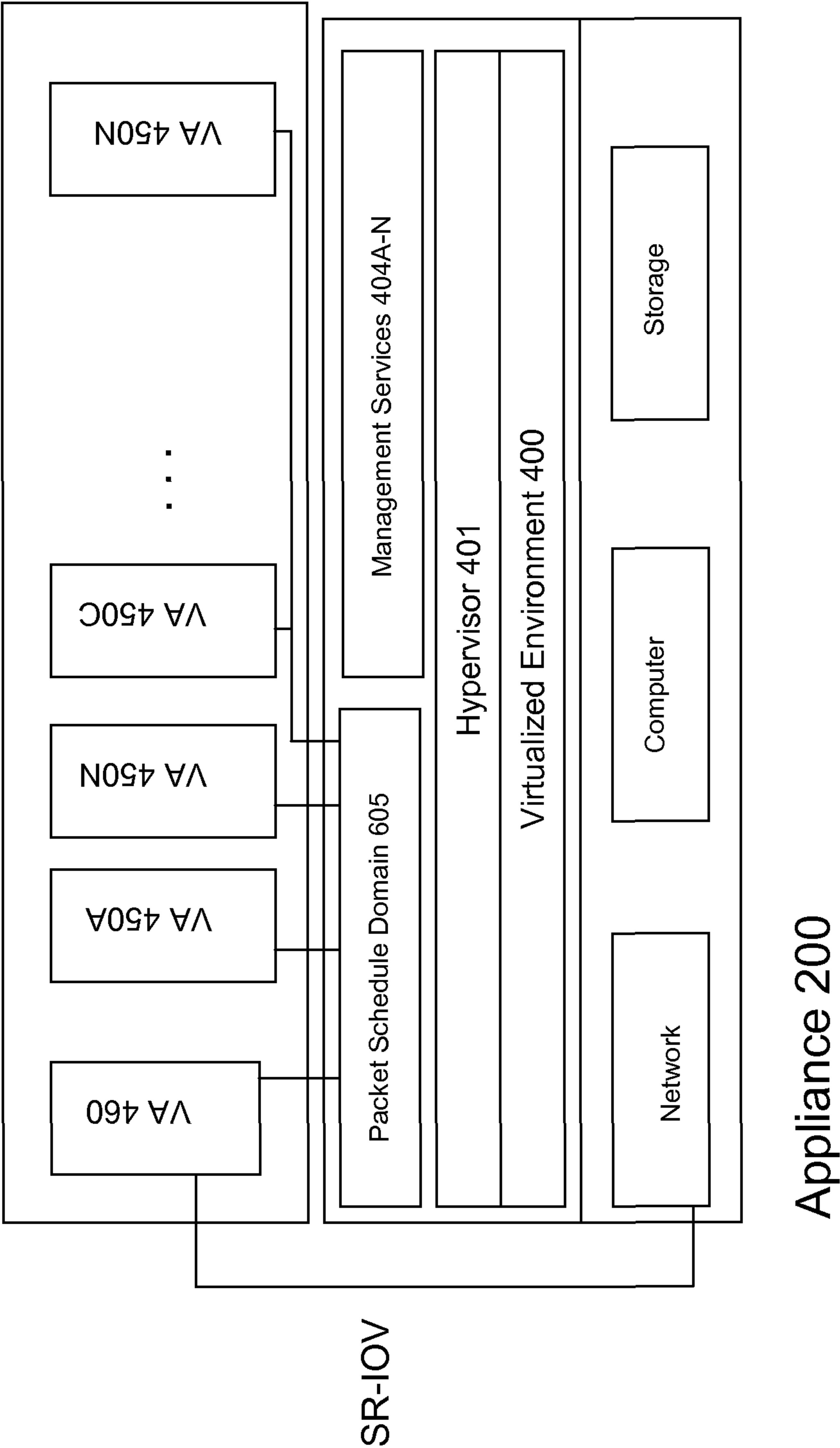


FIG. 6A

	Load Balancing Method	Balance Effectiveness	Compression History	Bandwidth Management	Simplicity	Fault Tolerant	Persistent Unit
1	Least connection (sub-optimal, easiest to configure)	↑↑	↓↓	?	↑↑	↑↑	Client
2	Static configuration (optimal WAN optimization but difficult to manage)	↑	↑↑	→	↓↓	↓↓ Single spare option is available	Branch
3	Least accumulated load using receive bandwidth + AgentID persistence (better WAN optimization characteristics)	↑	↑↑	↑↑	→	↑↑	Branch
4	Least connection + AgentID persistence (requires branches to be same size,)	↓	↑	↑↑	↑	↑	Branch
5	Source IP hash to load up one BR at a time	?	↑	↑	↑↑	↑↑	Branch

Legend:

↑↑

↑

→

↓

↓↓

?

Best

Better

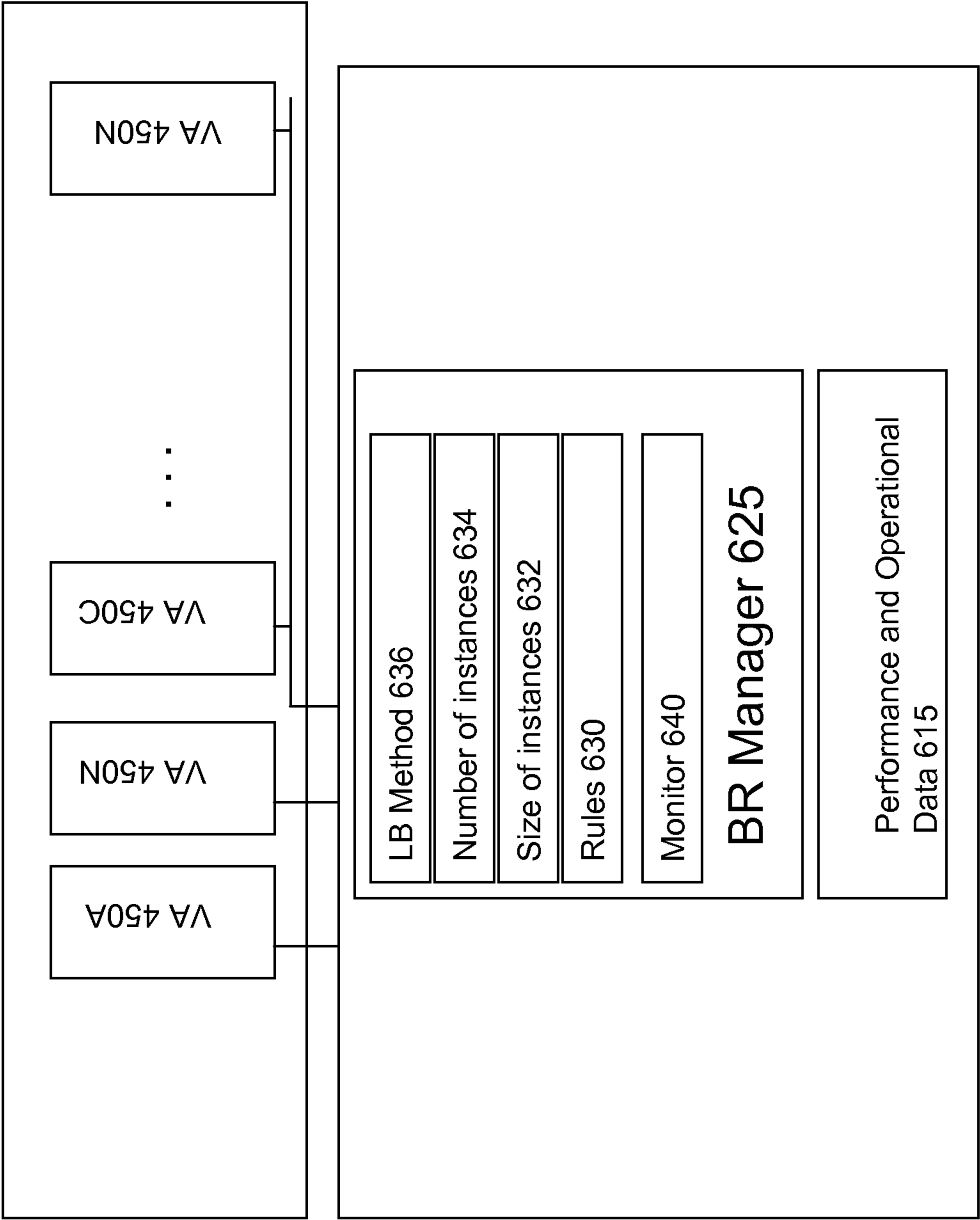
Neutral

Poor

Worst

Don't know

FIG. 6B



Appliance 200

FIG. 6C

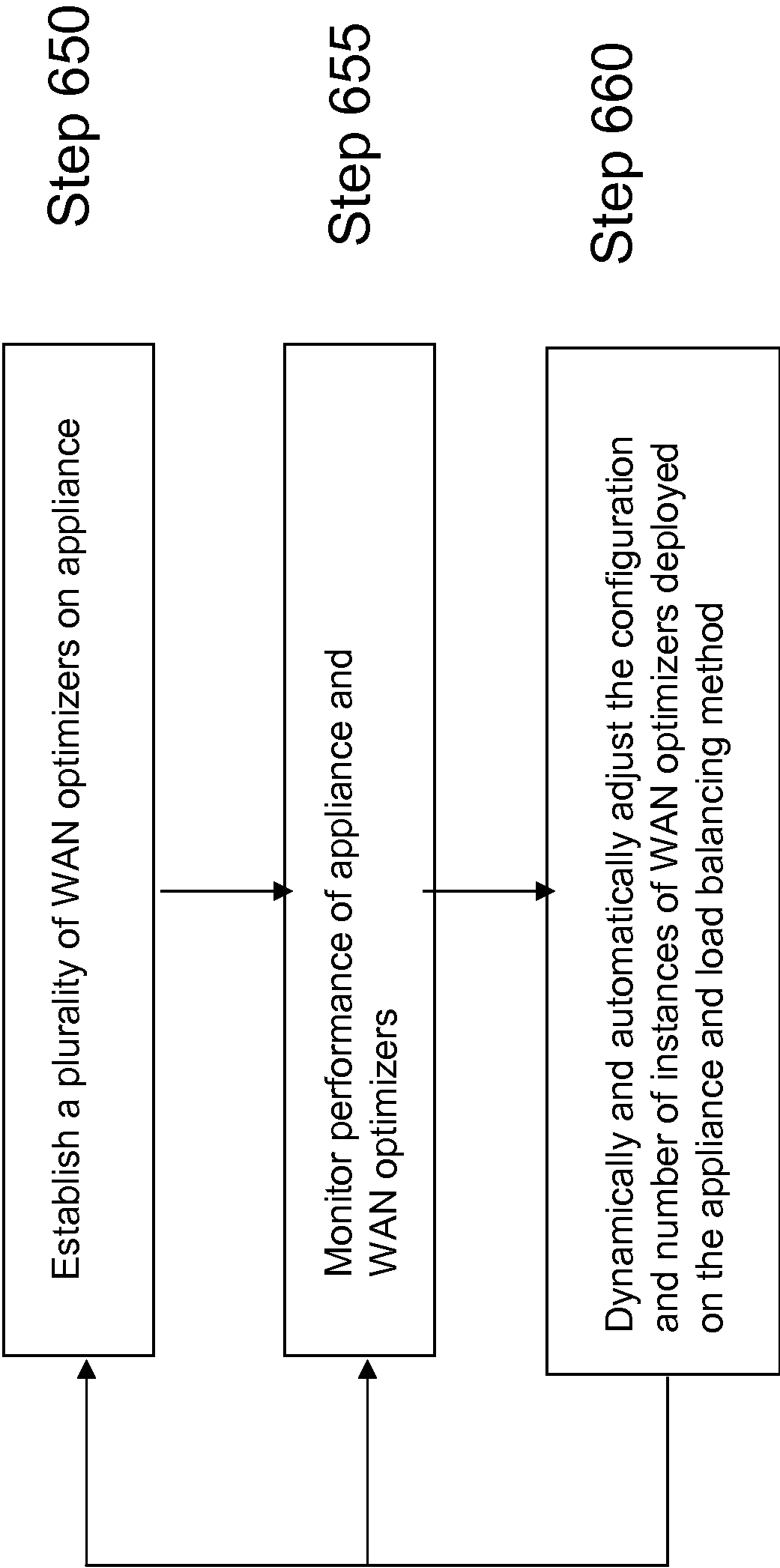


FIG. 6D

SYSTEMS AND METHODS FOR DYNAMIC ADAPTATION OF NETWORK ACCELERATORS

Related Application

This application claims the benefit of and priority to U.S. Provisional Application No.
5 61/547,493, entitled “Systems and Methods For Dynamic Adaptation of Network
Accelerators” and filed on October 14, 2011, which is incorporated herein by reference for all
purposes.

Field of the Invention

10 The present disclosure generally relates to data communication networks. In
particular, the present disclosure relates to systems and methods for the dynamic adaptation
of network accelerators on a platform.

Background of the Invention

15 Traditional network elements have been developed and deployed as discrete network
appliances and functions. As such, they have been deployed in a particular way and sized to
accommodate some sort of network traffic model. The model can be as simple as “I need
gigabit speed for this link” to more sophisticated models where traffic analysis and peak
loading are measured and placed in a model that determines network element sizing.

20

Brief Summary of the Invention

As computing in general and networking in particular moves to more virtualized
environments, there exist several problems with the prior models. As the workload become
more mobile and dynamic, traditional network engineering becomes virtually impossible and
25 new mechanisms must be used to “adapt” the network infrastructure to the offered load.

In some deployments, a load balancer may be used to load balance multiple WAN
optimizers. In further deployments, a virtualized load balancer may be used to load balance
multiple virtualized WAN optimizers. The load balancer may be configured to use or apply
any of a plurality of load balancing methods, such as but not limited to a least connection
30 method, least connection with agent id persistence, a static configuration, least accumulated
load using receive bandwidth, least accumulated load using receive bandwidth with agent id

persistence and source internet protocol (IP) hashing. Each of these methods may provide varying load balance effectiveness, compression history performance, bandwidth management and simplicity with respect to load balancing WAN optimizers. There may be no one-size-fits-all solution for all deployments. There may be a trade off between ease of
5 deployment and WAN optimization.

Systems and methods of the present solution provide a more optimal solution by dynamically and automatically reacting to changing network workload. A system that starts slowly, either by just examining traffic passively or by doing sub-optimal acceleration can learn over time, how many peer WAN optimizers are being serviced by an appliance, how
10 much traffic is coming from each peer WAN optimizers, and the type of traffic being seen. Knowledge from this learning can serve to provide a better or improved baseline for the configuration of an appliance. In some embodiments, based on resources (e.g., CPU, Memory, Disk), the system from this knowledge may determine how many WAN optimization instances should be used and of what size, and how the load should be
15 distributed across the instances of the WAN optimizer. Some example rules are as follows:

1. If a small number of peer WAN optimizers exist, a smaller number of large WAN optimizer instances should be provisioned on an appliance because compression histories will be less fragmented and compression ratios higher (better)
- 20 2. When peers WAN optimizers are of significantly different sizes, they should be distributed unevenly across the WAN optimizers instances (perhaps using WAN optimizers instances of different sizes)

Over time the system may continue to monitor the workload for changes and change the load
25 distribution as needed, for example, adding additional WAN optimizers instances because first pass compression (high CPU utilization) is being used extensively and compression history fragmentation effects would be minimal. Or traffic from a remote site has grown significantly and a larger WAN optimizers instance is required to keep compression history from fragmenting.

30 In one aspects, the present solution is directed to a method for managing a plurality of instances of a Wide Area Network (WAN) optimizer executing on an intermediary device. The method includes establishing, on a device intermediary to a plurality of clients and plurality of servers, a plurality of instances of a Wide Area Network (WAN) optimizer to accelerate WAN communications between the plurality of clients and the plurality of servers.

The method includes monitoring, by the device, network traffic traversing the device for each of the plurality of instances of the WAN optimizer and selecting, by a manager executing on the device responsive to the monitoring, a change of a load balancing scheme to load balance the plurality of instances of the WAN optimizer.

5 The method may automatically establish, by the device, a configuration of a size of each of the plurality of instances of the WAN optimizer based on data stored from monitoring of previous execution of the plurality of instances of the WAN optimizer. Each of the plurality of instances of the WAN optimizer may execute as a virtual machine in a virtualized environment.

10 In some embodiments, the method includes comprises monitoring, by the device, compression history allocation, compression fragmentation and compression ratios of each of the plurality of instances of the WAN optimizer. In some embodiments, the method includes monitoring, by the device, or more of the following of each of the plurality of instances of the WAN optimizer: resource utilization, number of connections, number of claims and
15 bandwidth usage.

 In some embodiments, the method includes determining, by the device, that a metric computed from monitoring network traffic has exceeded a threshold and responsive to the determination, automatically selecting by the device a second load balancing scheme to load balance the plurality of instances of the WAN optimizer. In some embodiments, the method
20 includes automatically switching, by the device, from the load balancing scheme to the selected load balancing scheme while executing the plurality of instances of the WAN optimizer. In some embodiments, the method includes automatically changing, by the device responsive to the monitoring, the number of instances of the WAN optimizer executing on the device. In some embodiments, the method includes automatically adjusting, by the
25 device responsive to the monitoring, a size of resource usage used by one or more of the plurality of instances of the WAN optimizer. The method may also include applying, by the device one or more rules to data collected from monitoring, to determine to change one or more of the following: a number of instances of the WAN optimizer, a size of one or more WAN optimizers and the load balancing scheme.

30 In some aspects, the present solution is directed to a system for managing a plurality of instances of a Wide Area Network (WAN) optimizer executing on an intermediary device. The system includes a device intermediary to a plurality of clients and plurality of servers and a plurality of instances of a Wide Area Network (WAN) optimizer executing on the device to accelerate WAN communications between the plurality of clients and the plurality of servers;

