

(19) **United States**

(12) **Patent Application Publication**  
Chou et al.

(10) **Pub. No.: US 2020/0201661 A1**

(43) **Pub. Date: Jun. 25, 2020**

(54) **GUARANTEEING AVAILABILITY OF TARGET DATA TO REMOTE INITIATORS VIA A HYBRID SOURCE/TARGET CREDIT SCHEME**

(71) Applicant: **Diamanti, Inc.**, San Jose, CA (US)

(72) Inventors: **Jeffrey Chou**, Palo Alto, CA (US);  
**Gopal Sharma**, San Jose, CA (US);  
**Amitava Guha**, San Jose, CA (US);  
**Kevin Fong**, Las Vegas, NV (US);  
**Jayasenana Sundara Ganesh**,  
Cupertino, CA (US)

(21) Appl. No.: **16/806,999**

(22) Filed: **Mar. 2, 2020**

#### Related U.S. Application Data

(63) Continuation of application No. 14/640,717, filed on Mar. 6, 2015, now abandoned.

(60) Provisional application No. 61/950,036, filed on Mar. 8, 2014, provisional application No. 62/017,257, filed on Jun. 26, 2014.

#### Publication Classification

(51) **Int. Cl.**

**G06F 9/455** (2006.01)

**H04L 29/08** (2006.01)

**G06F 3/06** (2006.01)

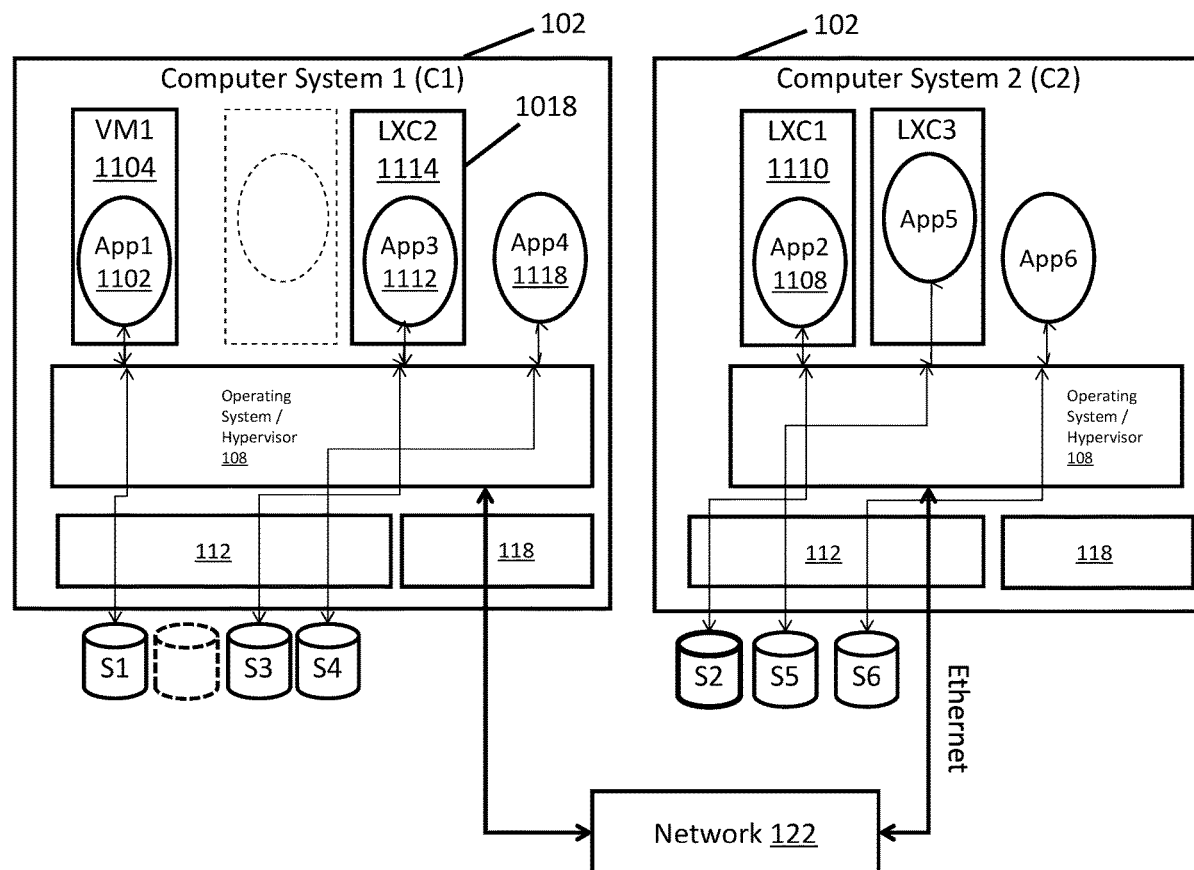
(52) **U.S. Cl.**

CPC ..... **G06F 9/45541** (2013.01); **H04L 67/1097** (2013.01); **H04L 67/1002** (2013.01); **G06F 3/0664** (2013.01); **G06F 3/061** (2013.01); **G06F 3/067** (2013.01)

(57)

#### ABSTRACT

A device includes a converged input/output controller that includes a physical target storage media controller, a physical network interface controller and a gateway between the storage media controller and the network interface controller, wherein gateway provides a direct connection for storage traffic and network traffic between the storage media controller and the network interface controller.



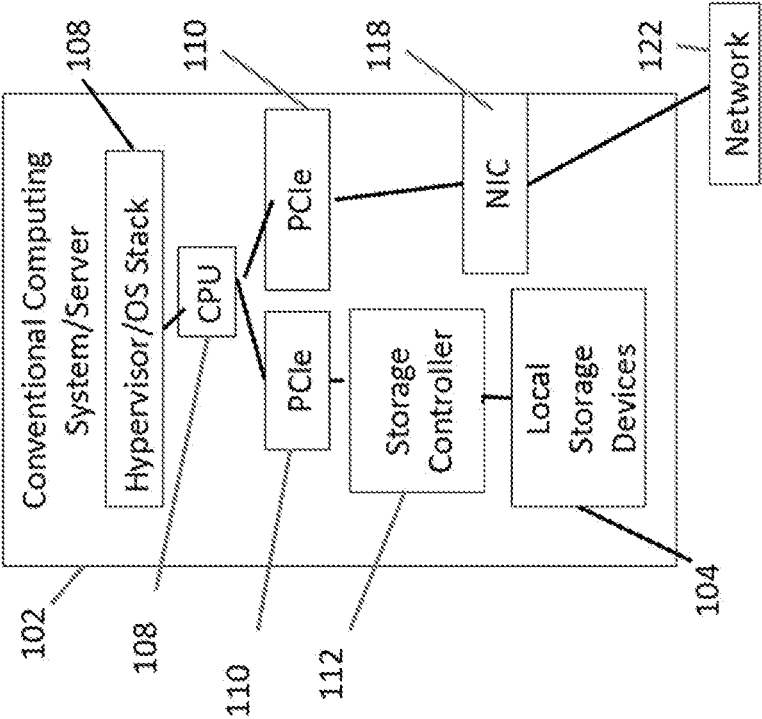


Fig. 1  
PRIOR ART

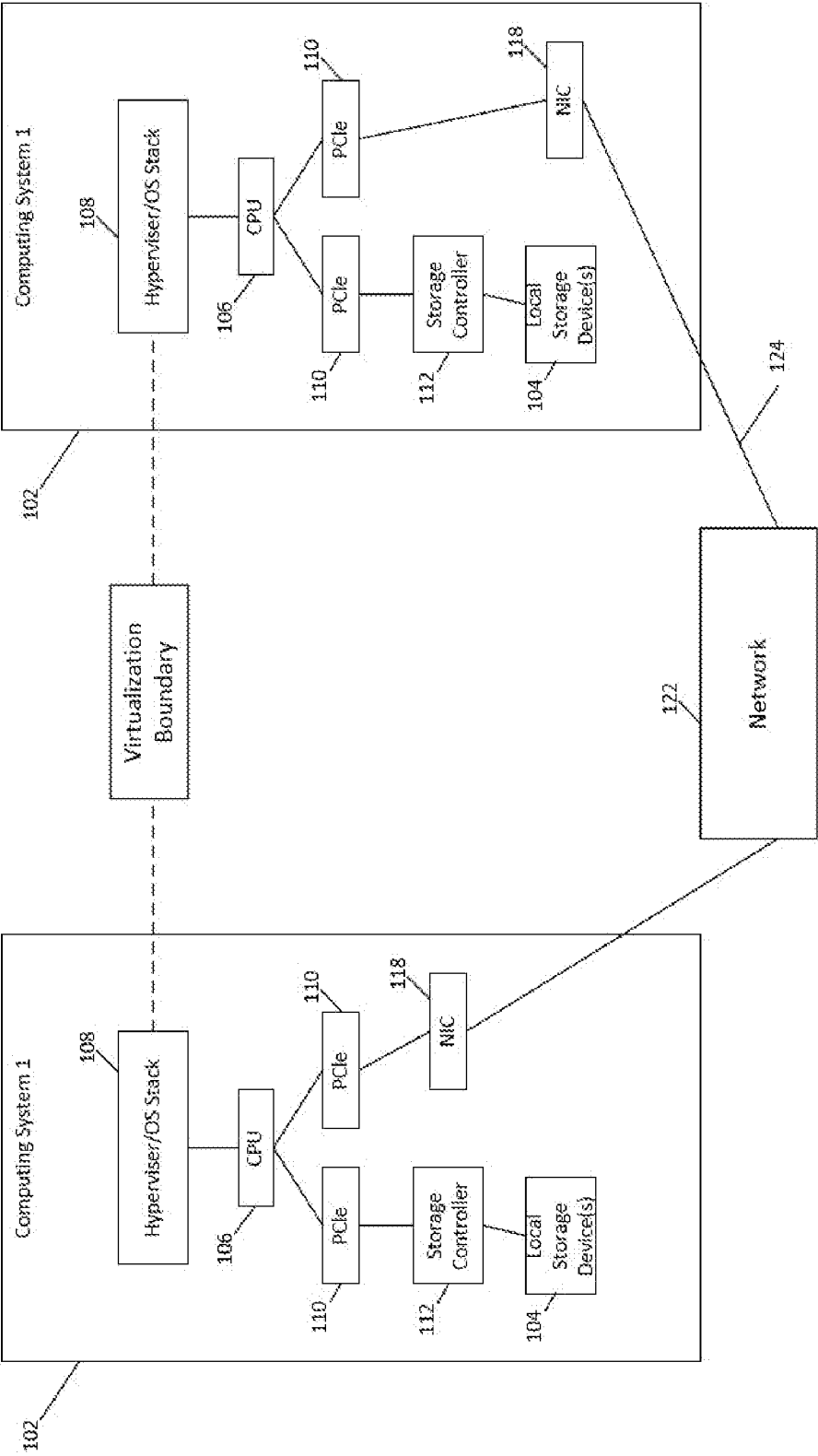


FIG. 2

PRIOR ART

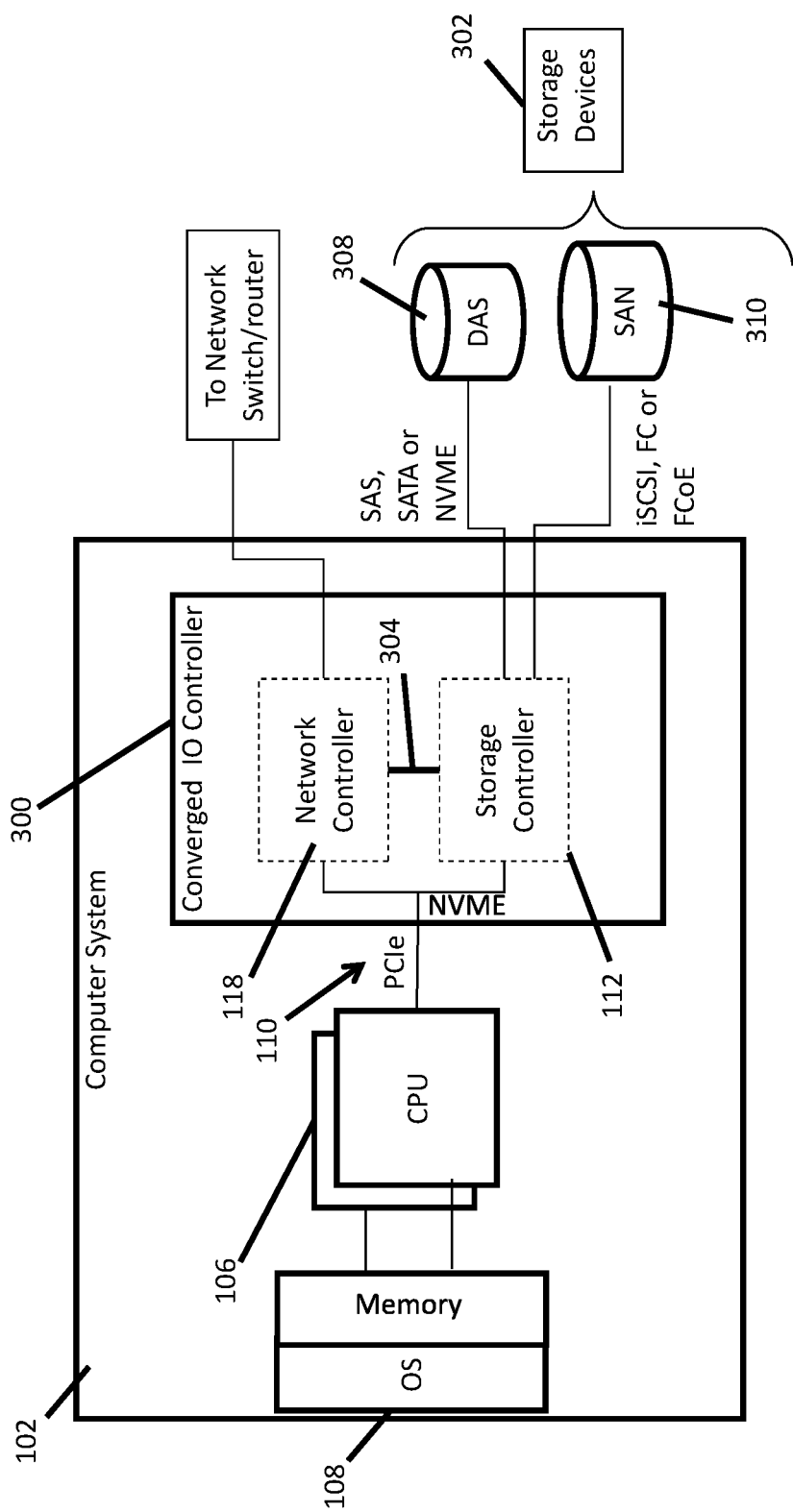


Fig. 3

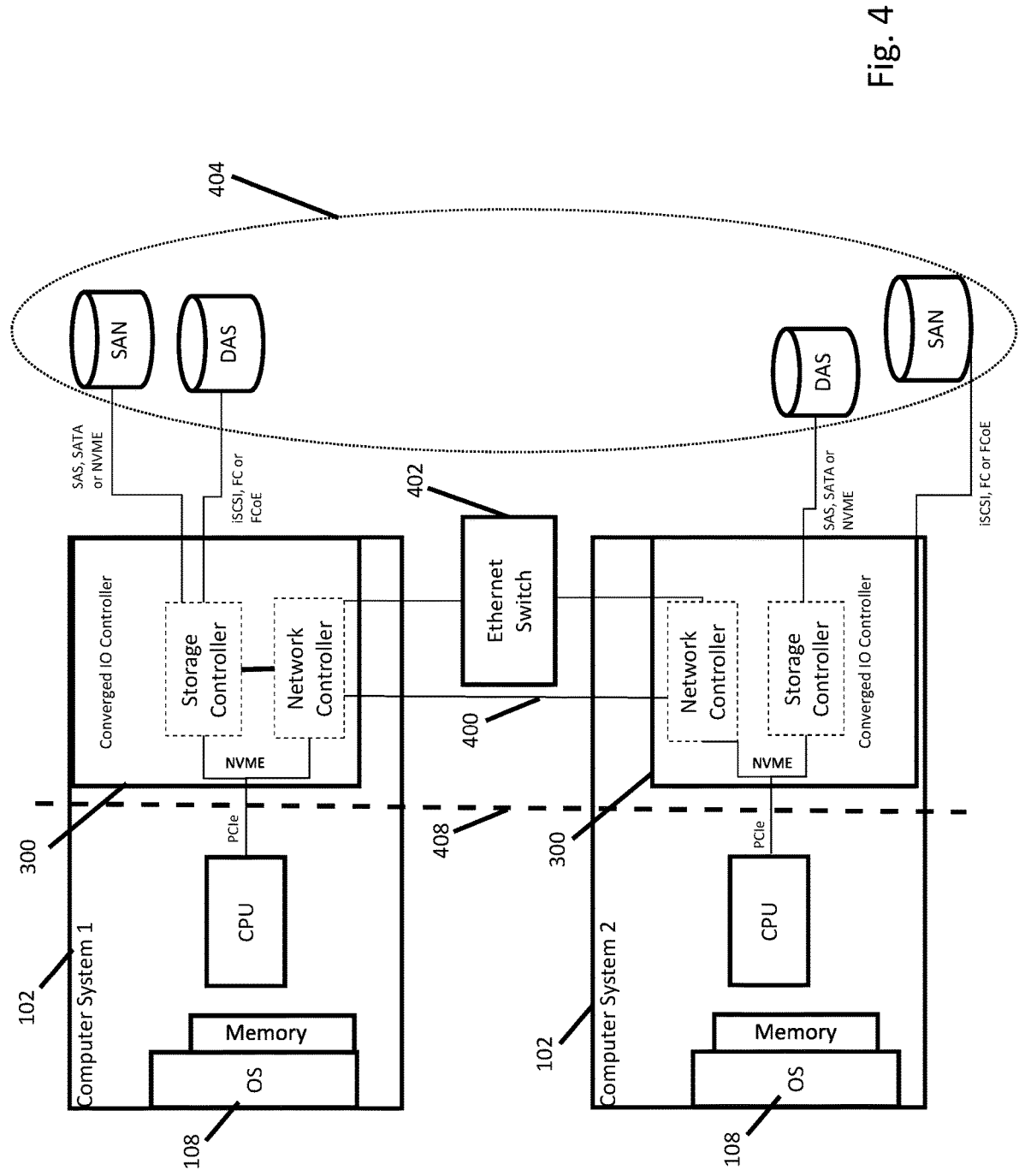


Fig. 4

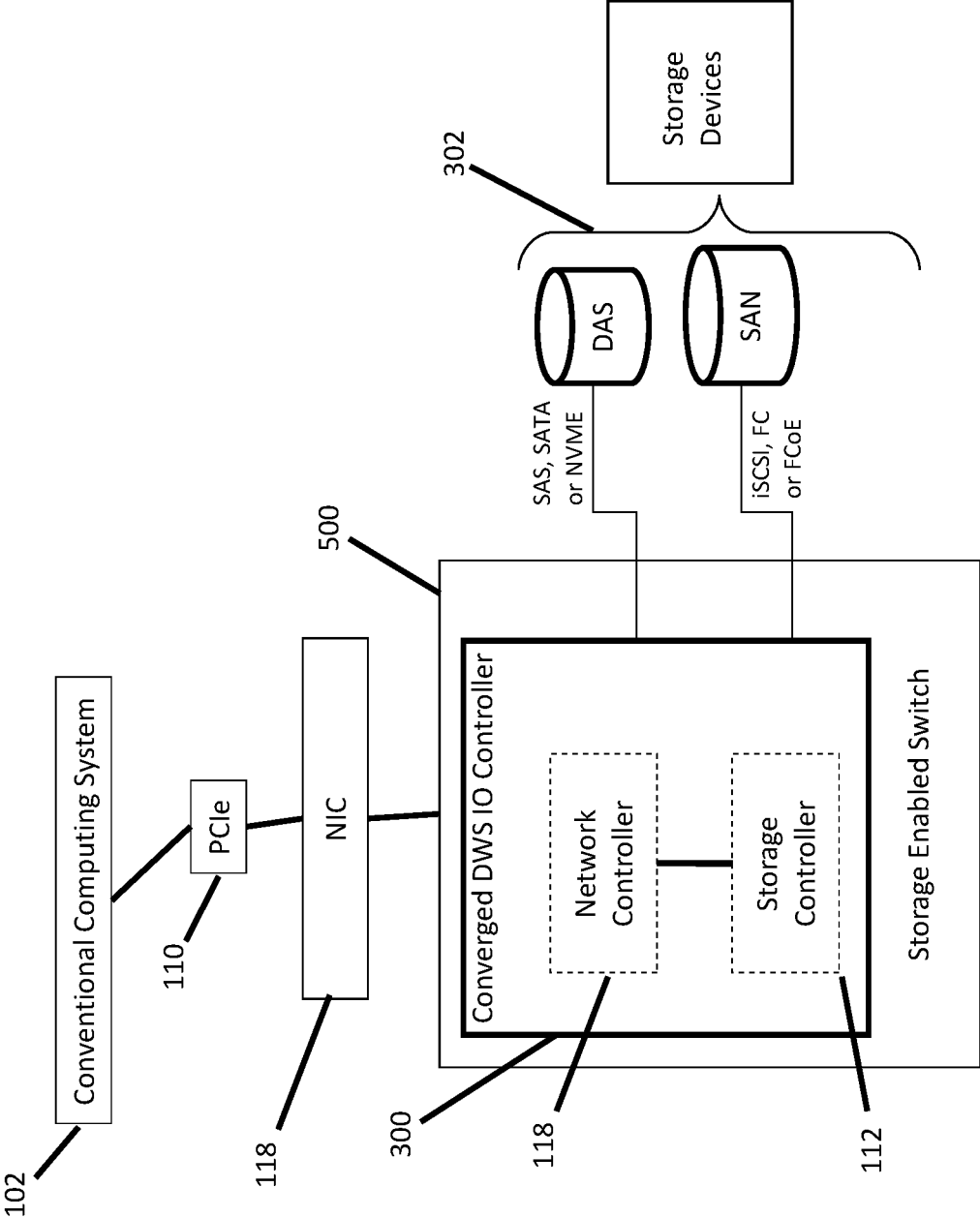
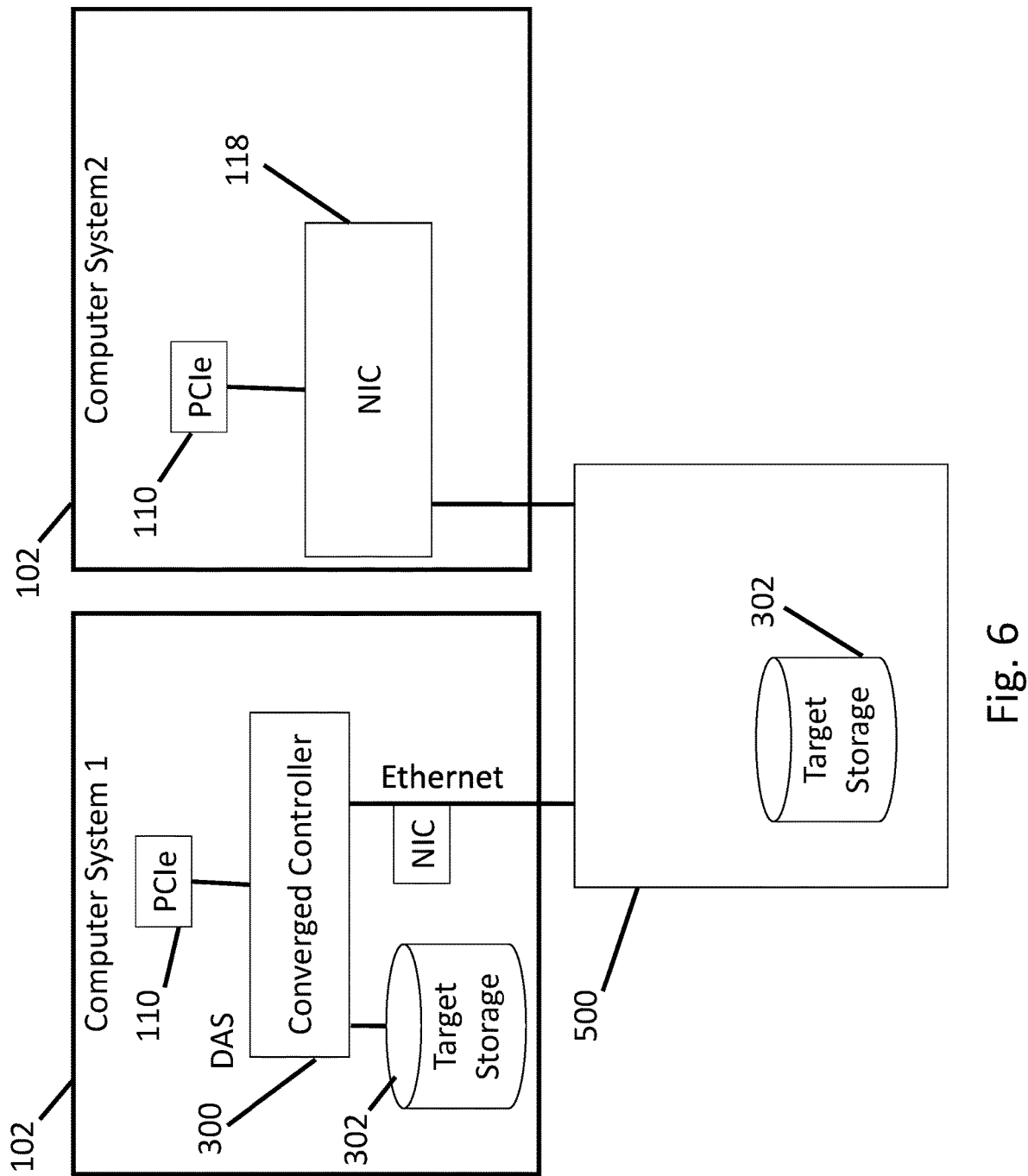
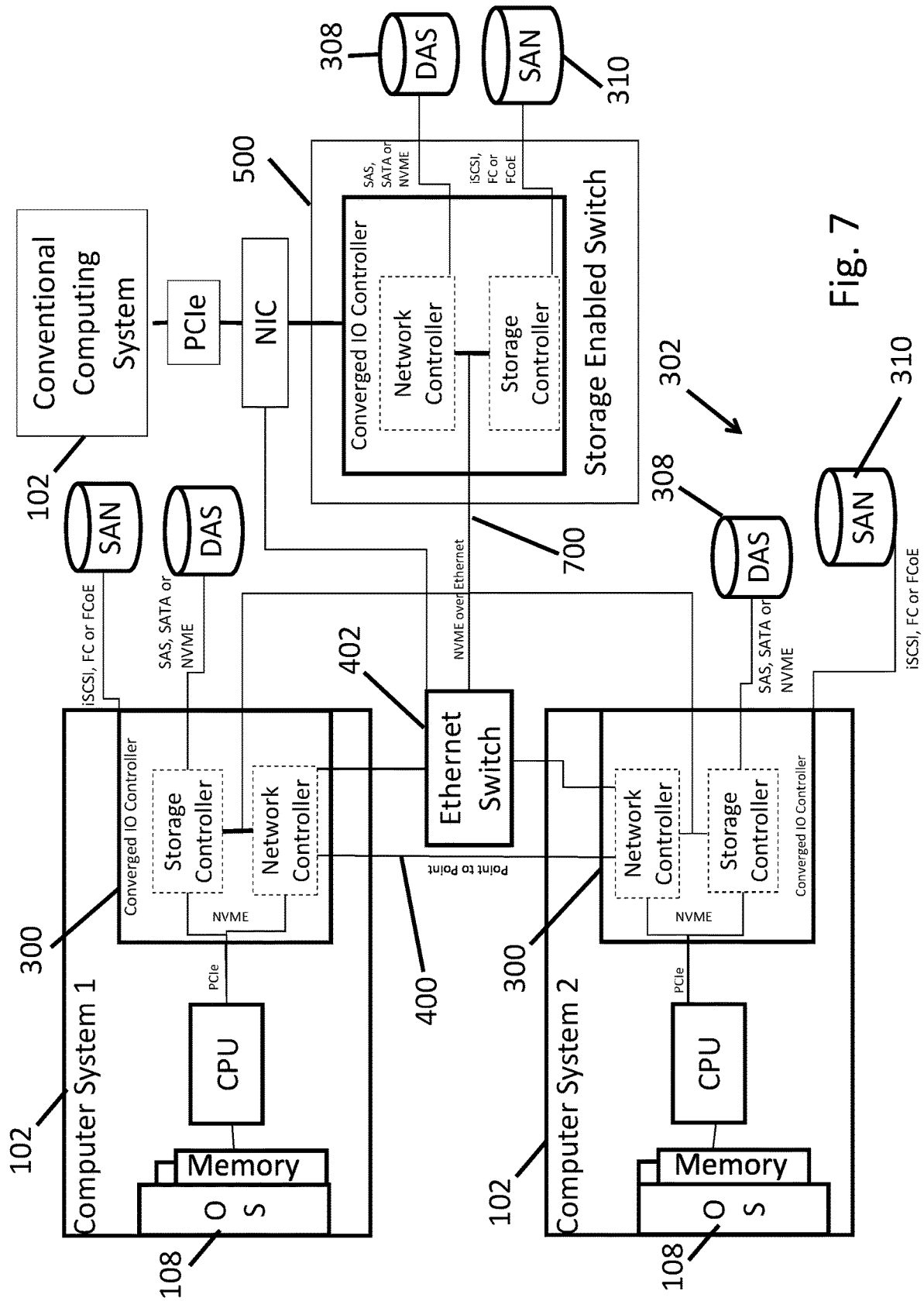


Fig. 5







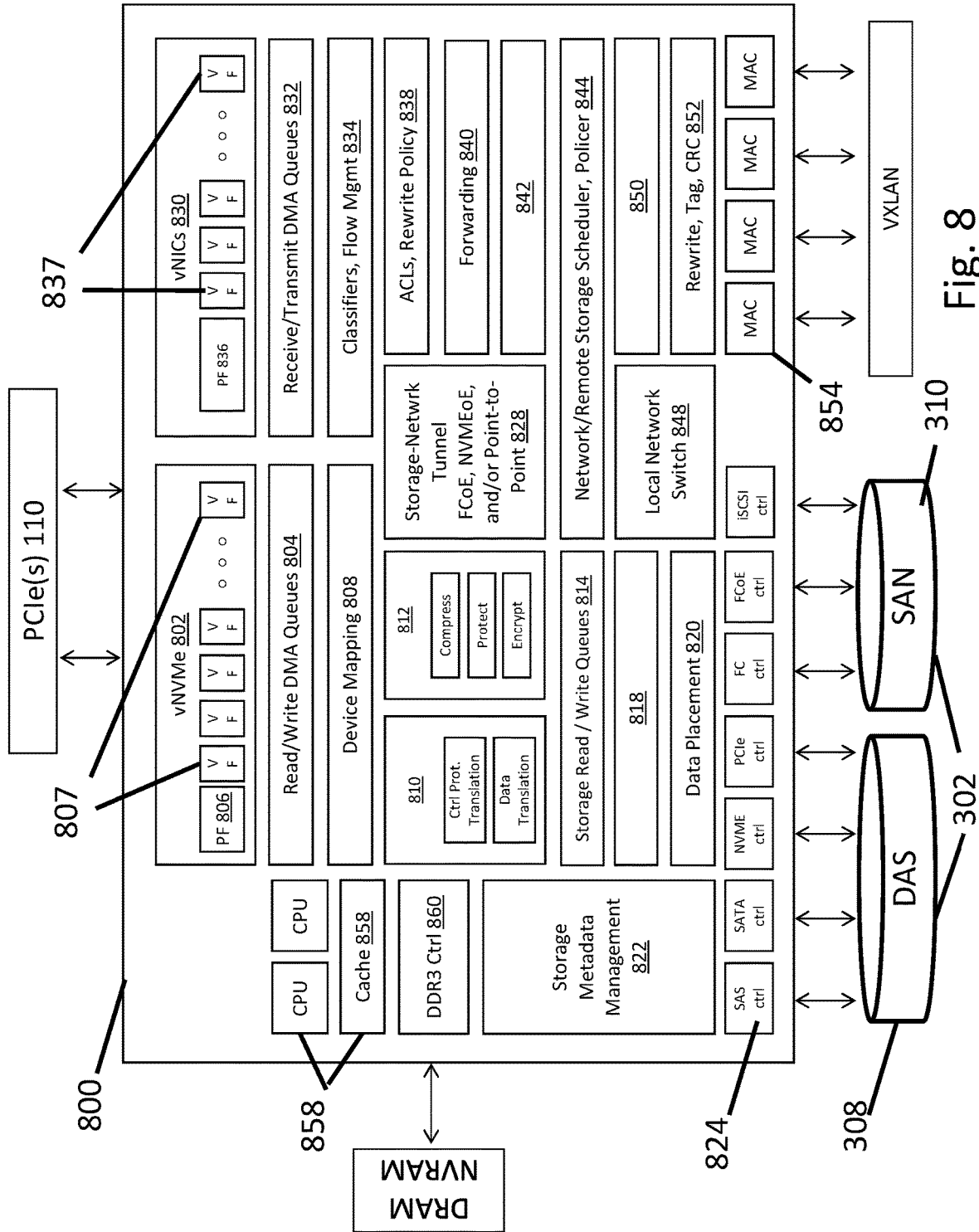


Fig. 8

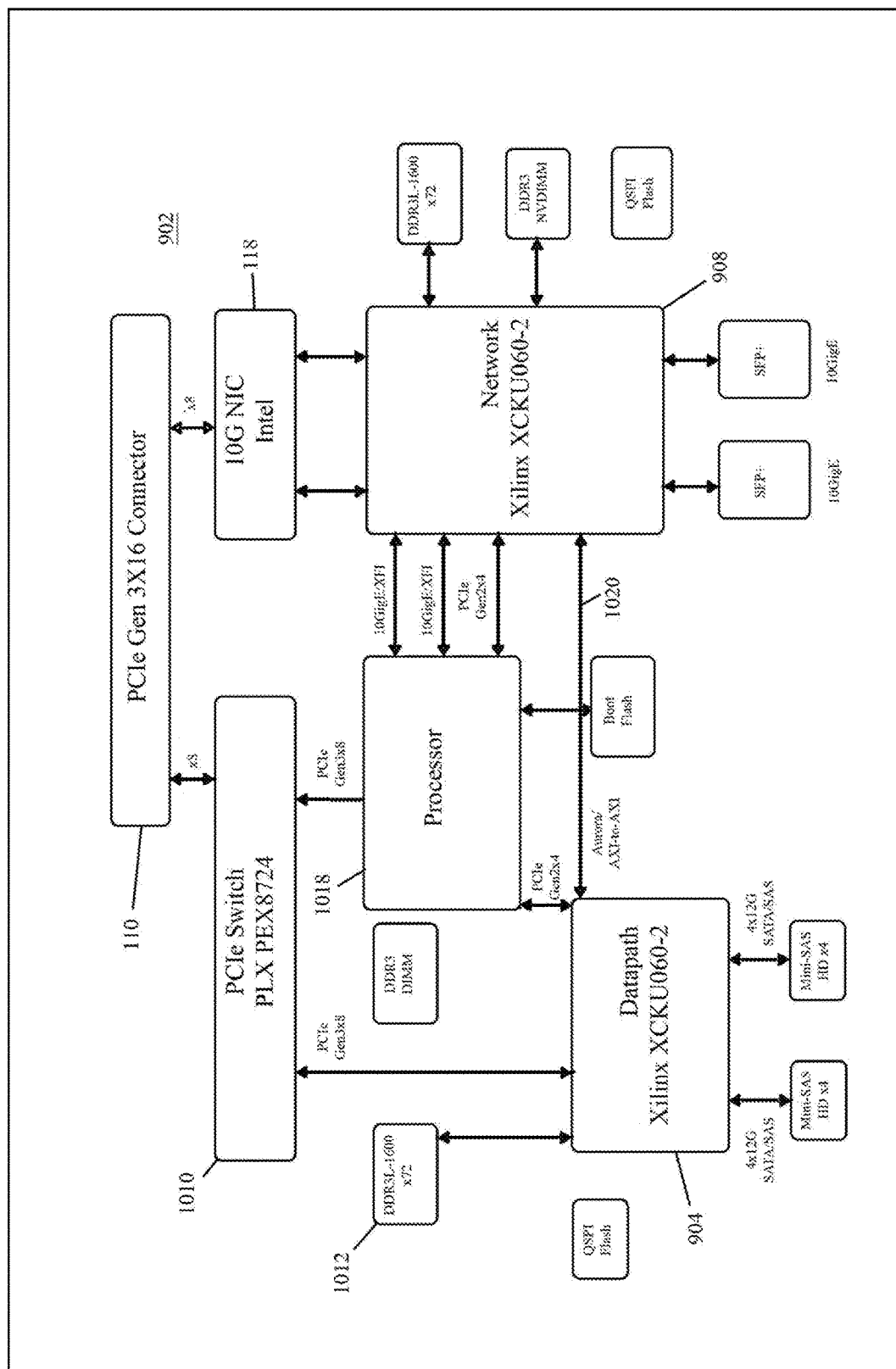


Fig. 9

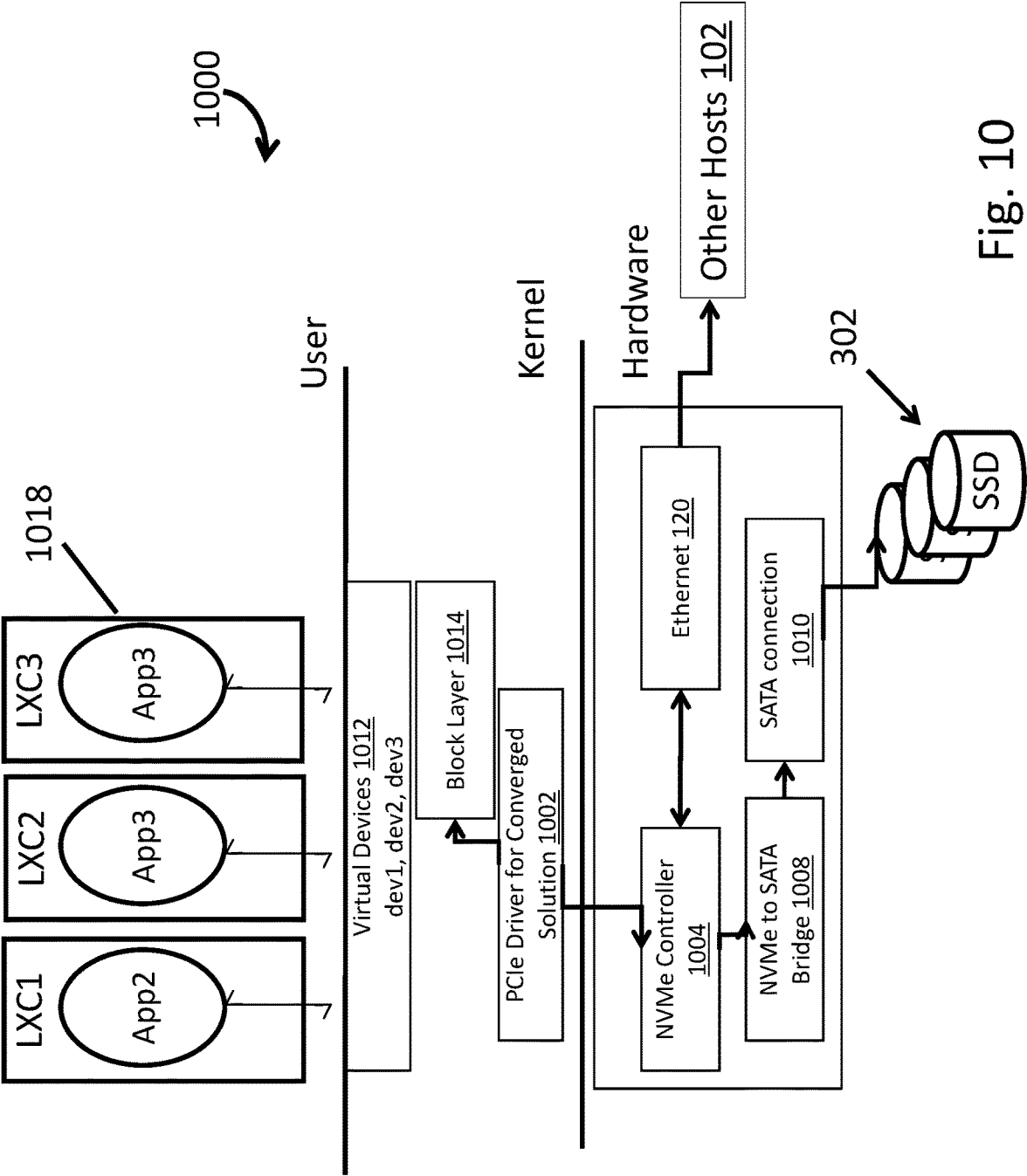


Fig. 10

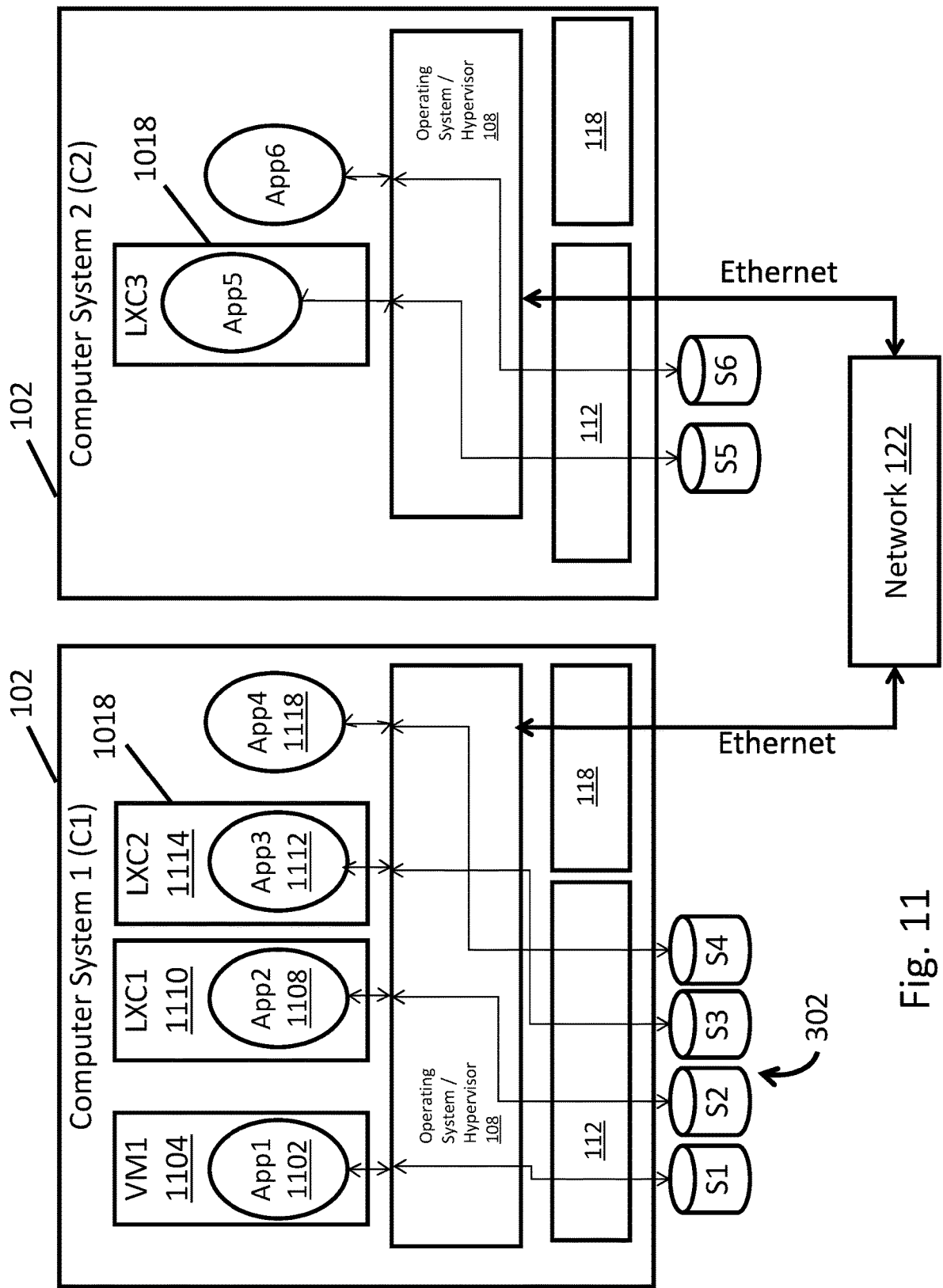


Fig. 11

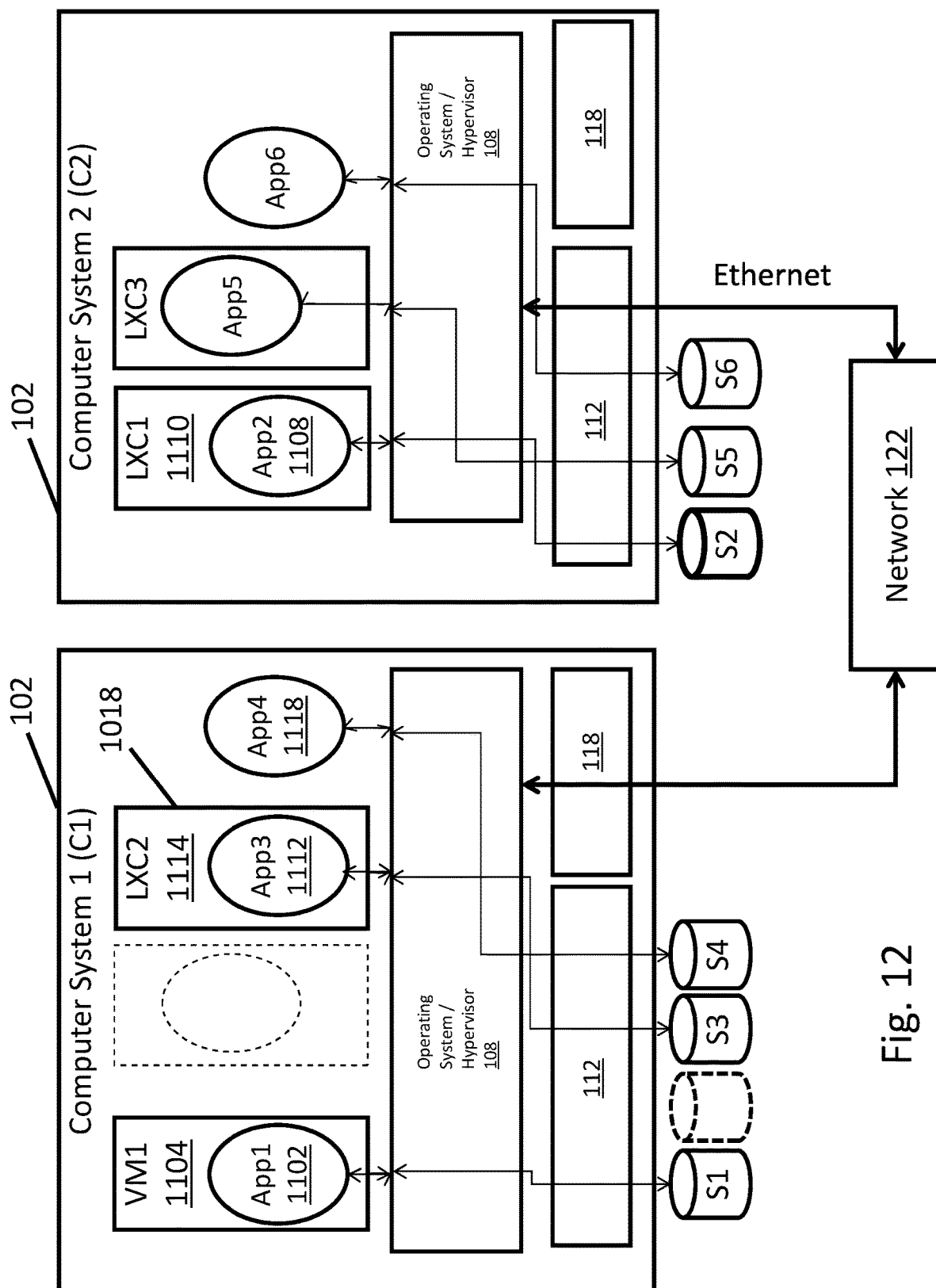


Fig. 12

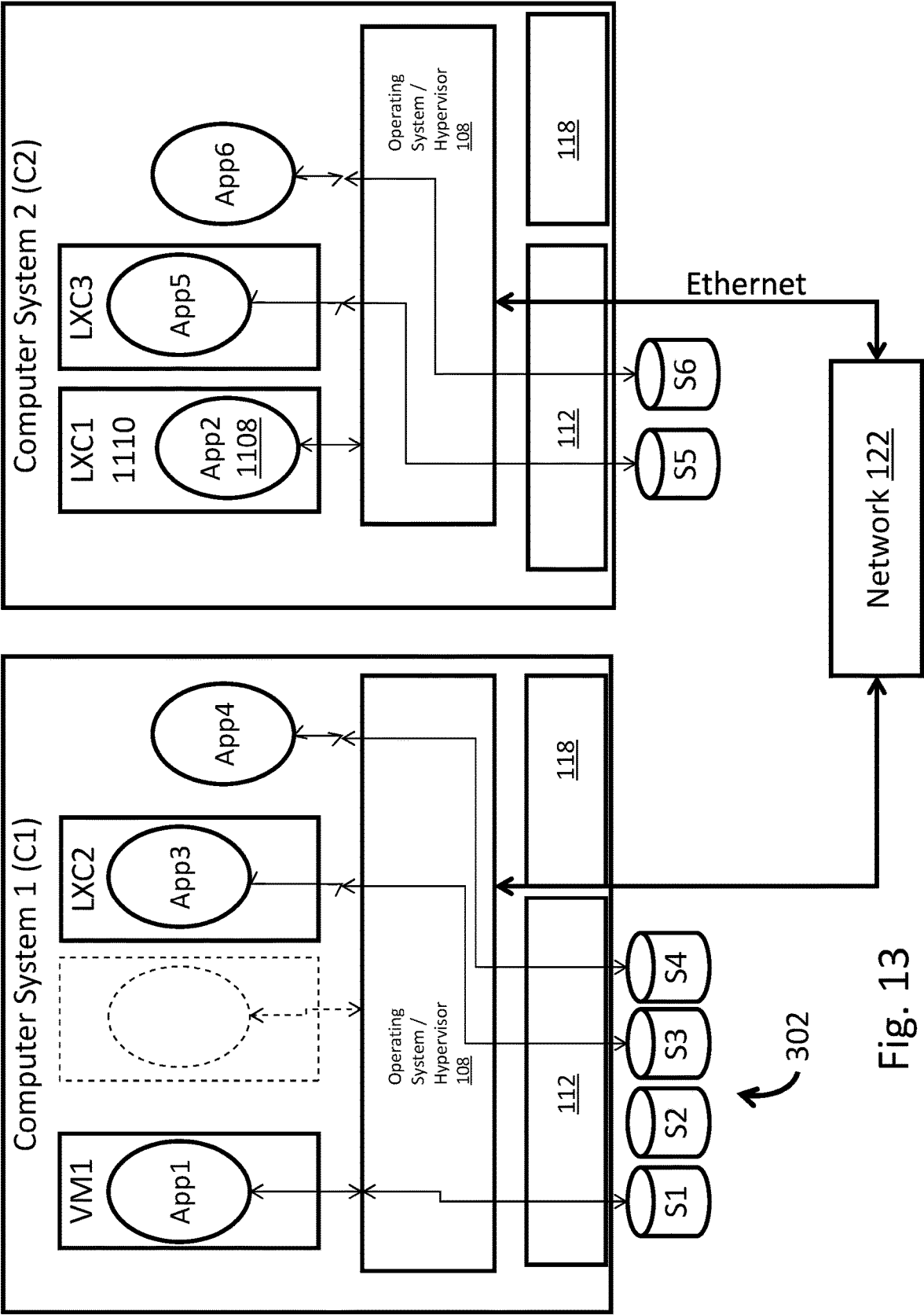


Fig. 13

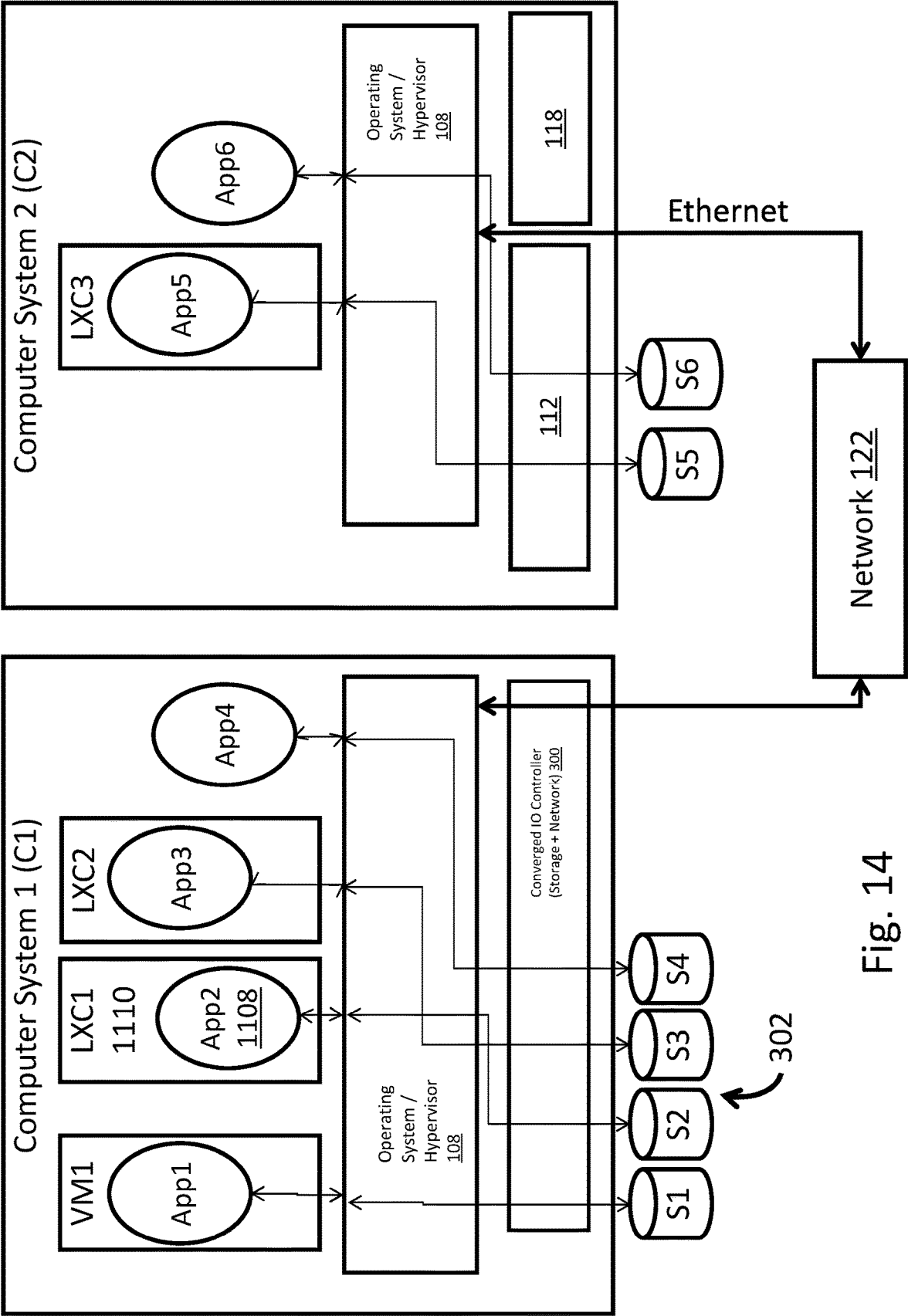


Fig. 14

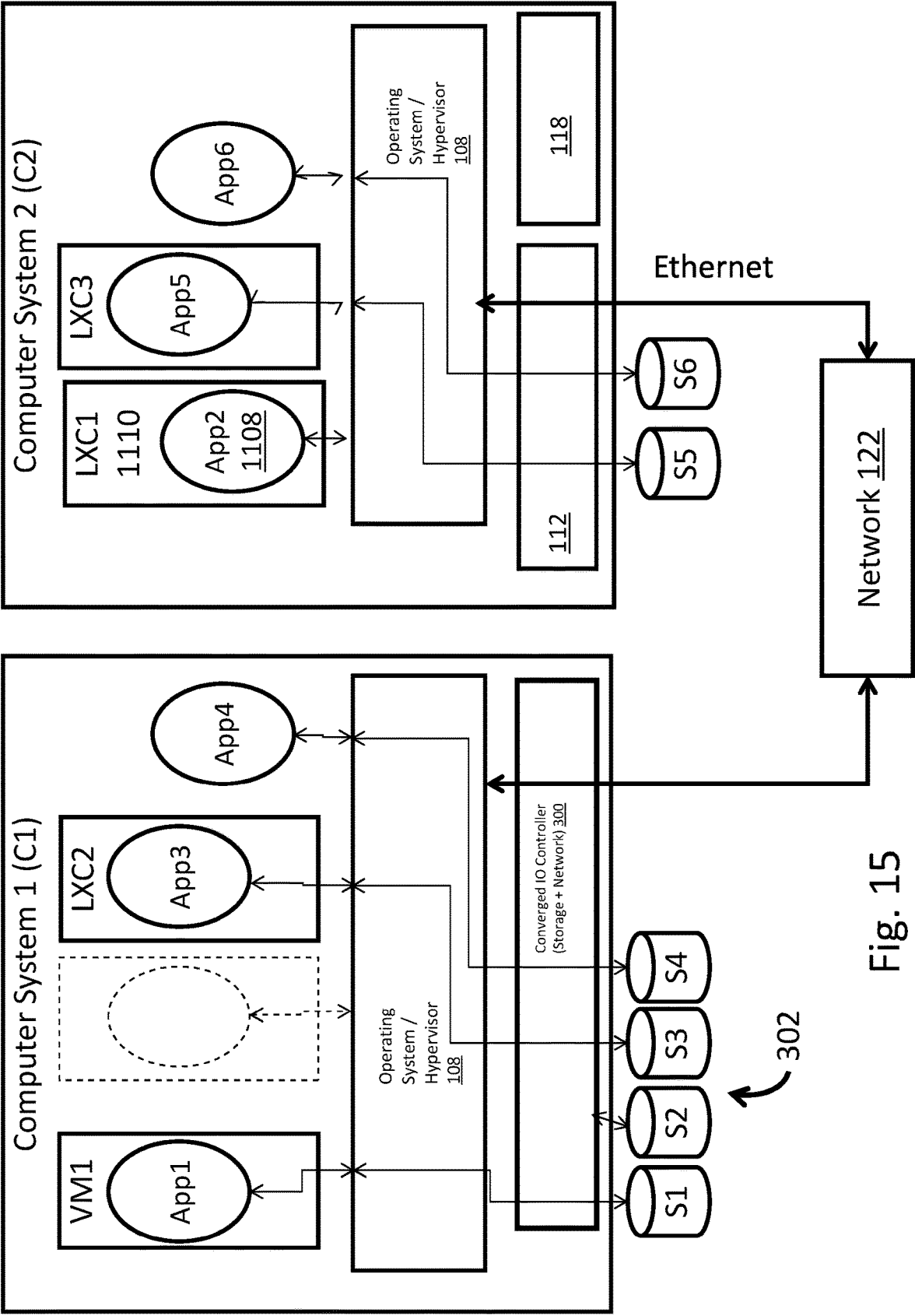


Fig. 15



302

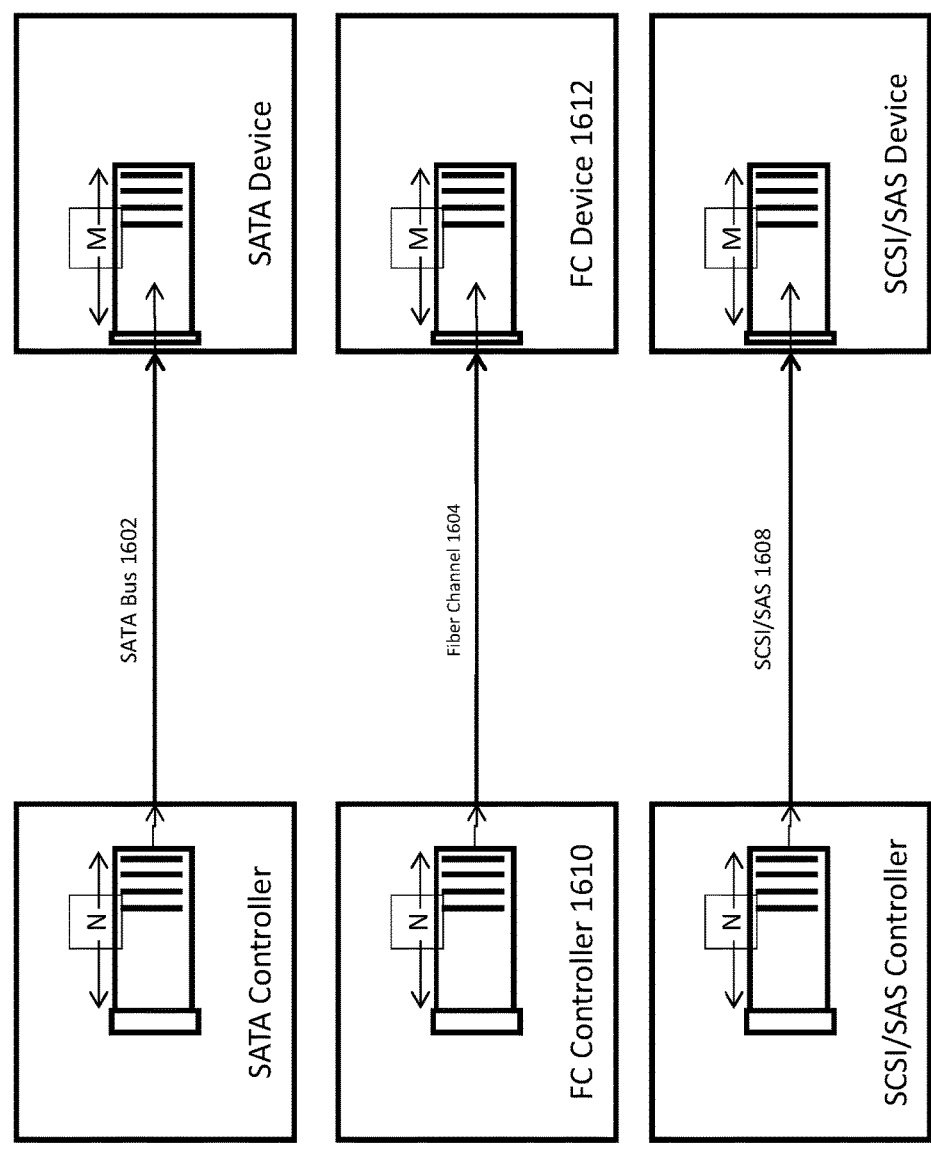


Fig. 16

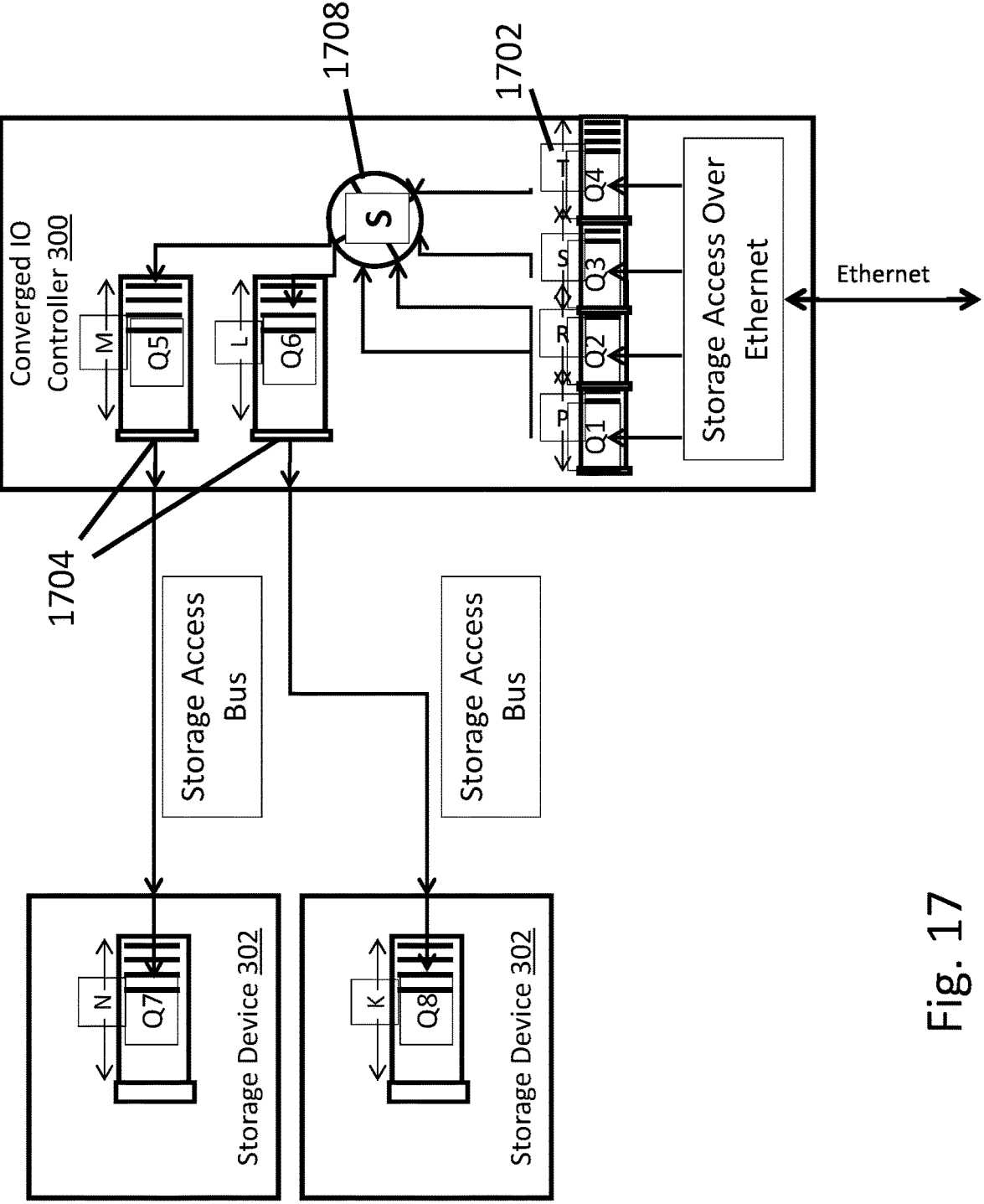


Fig. 17

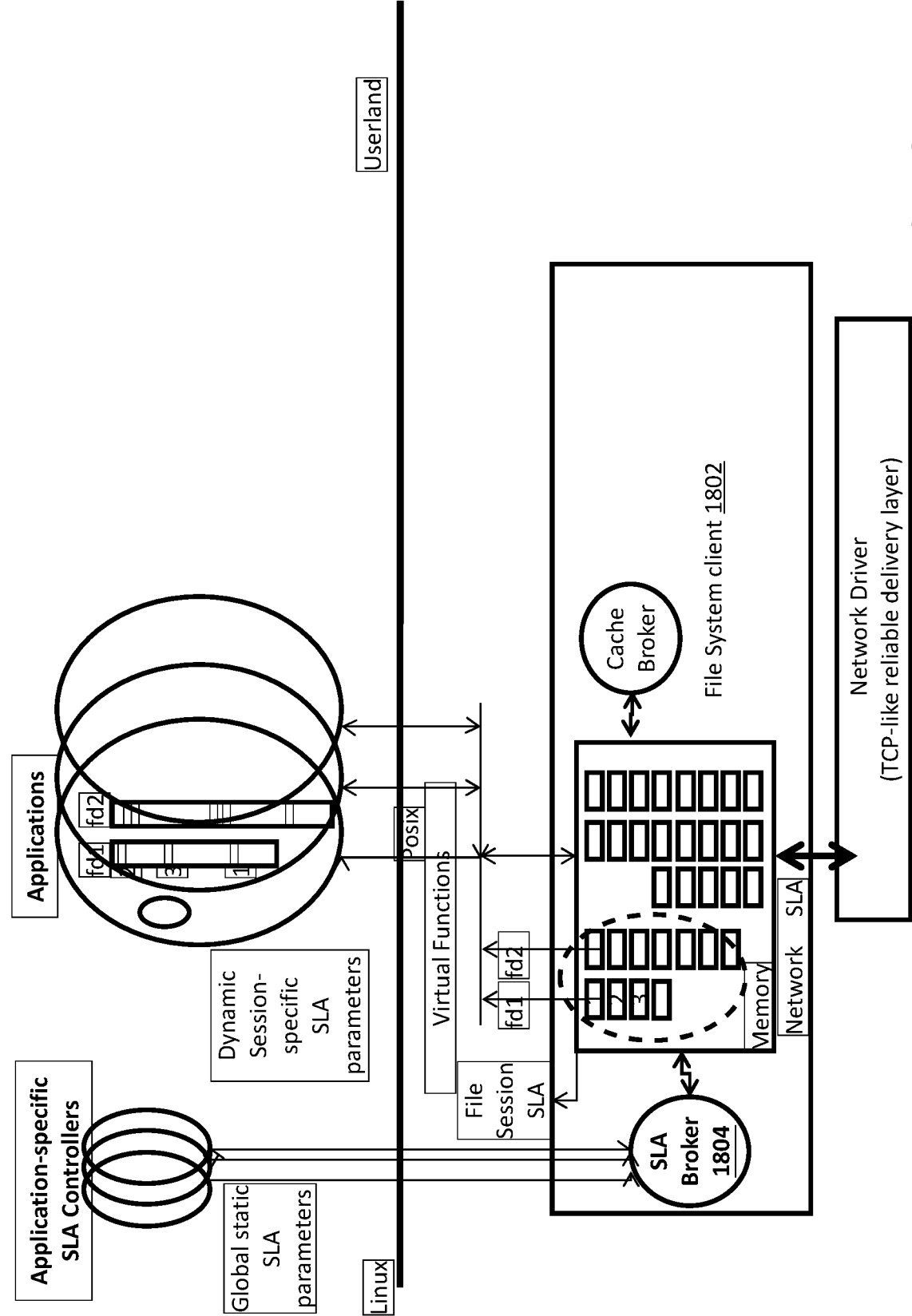


Fig. 18

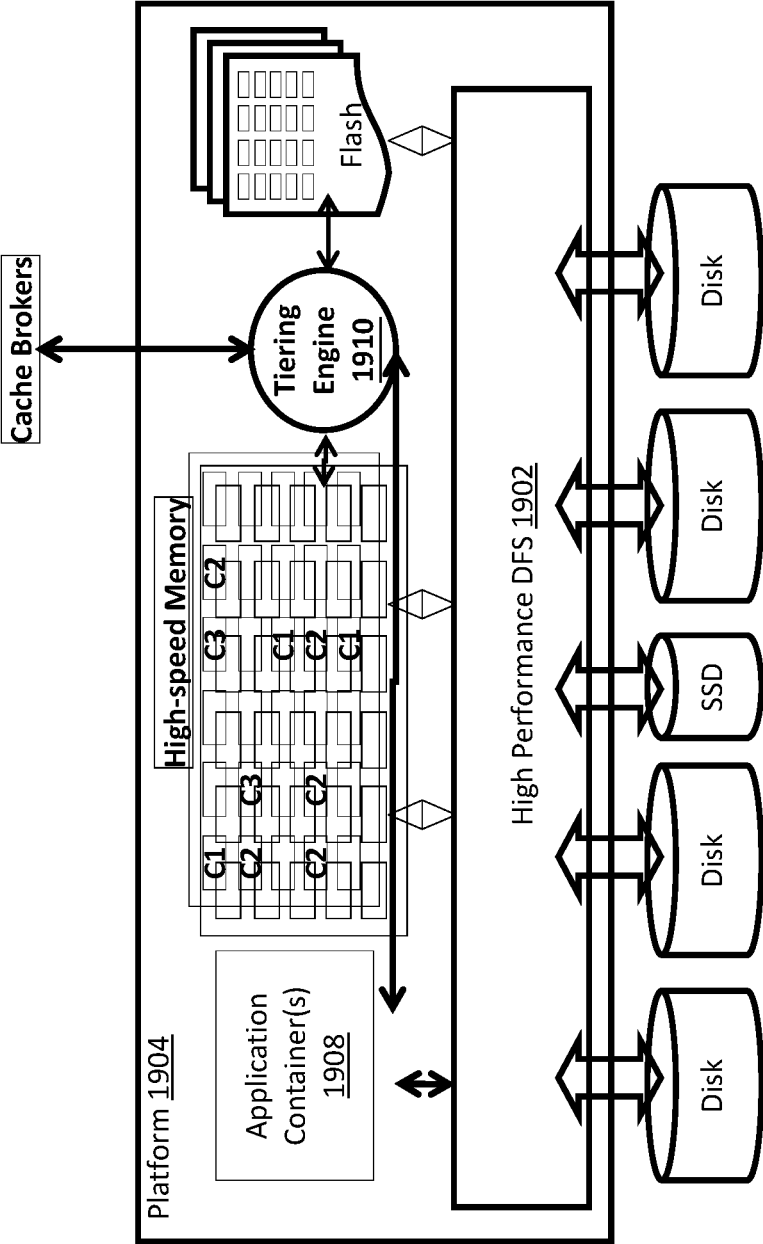


Fig. 19

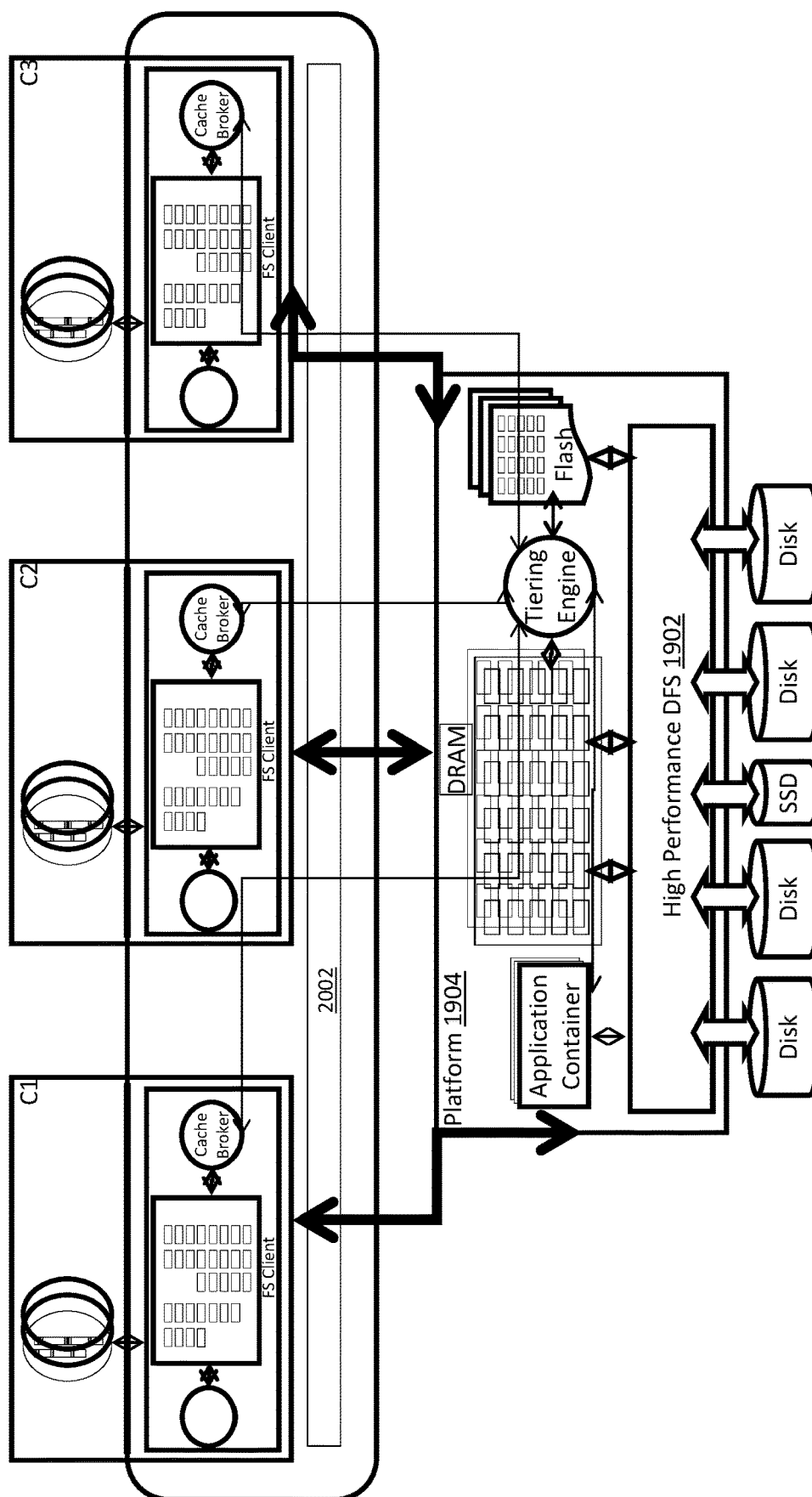


Fig. 20

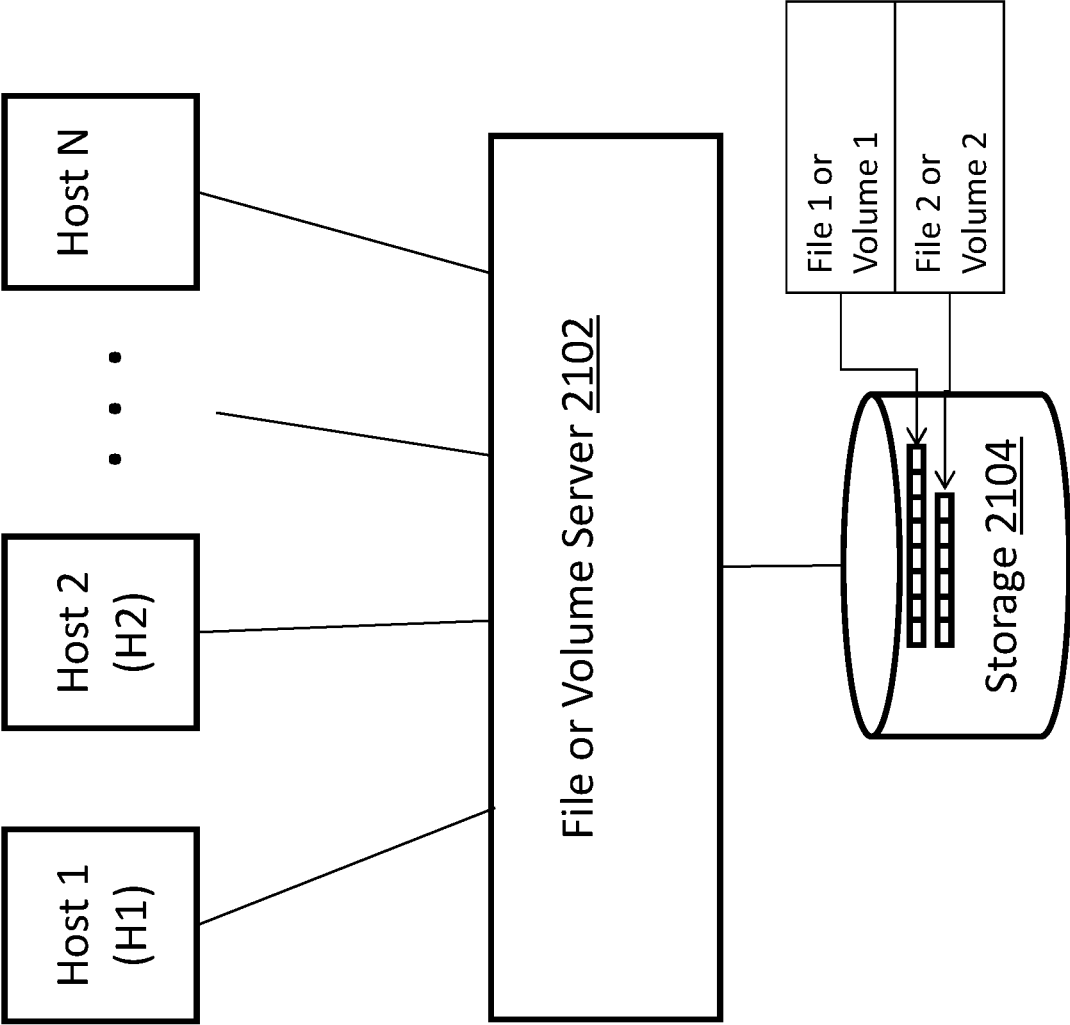


Fig. 21

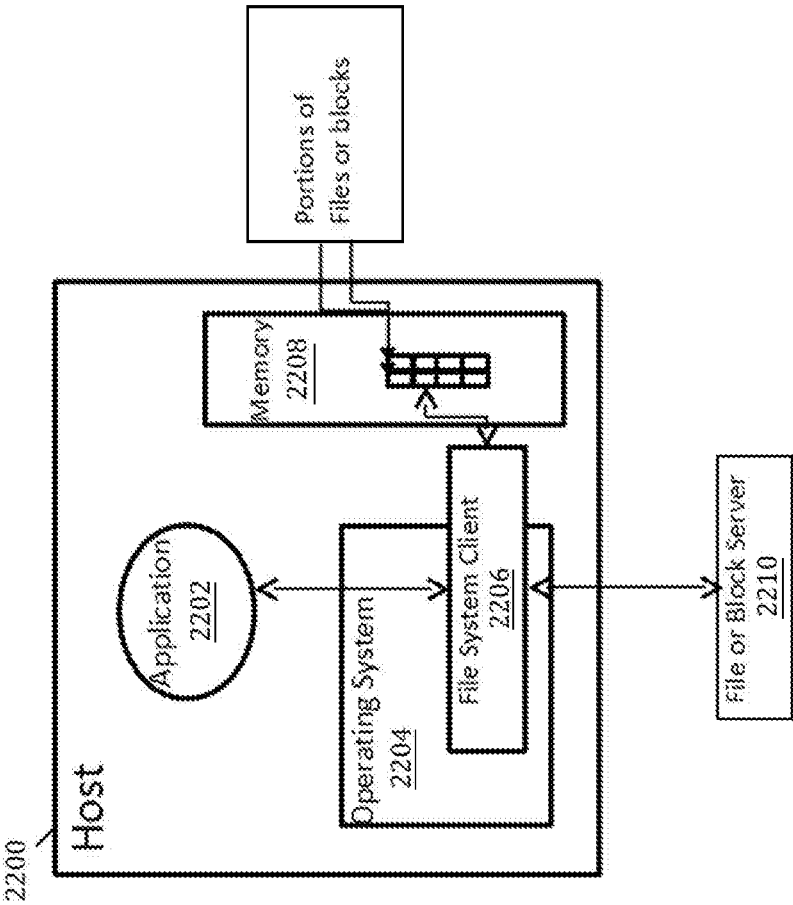


Fig. 22

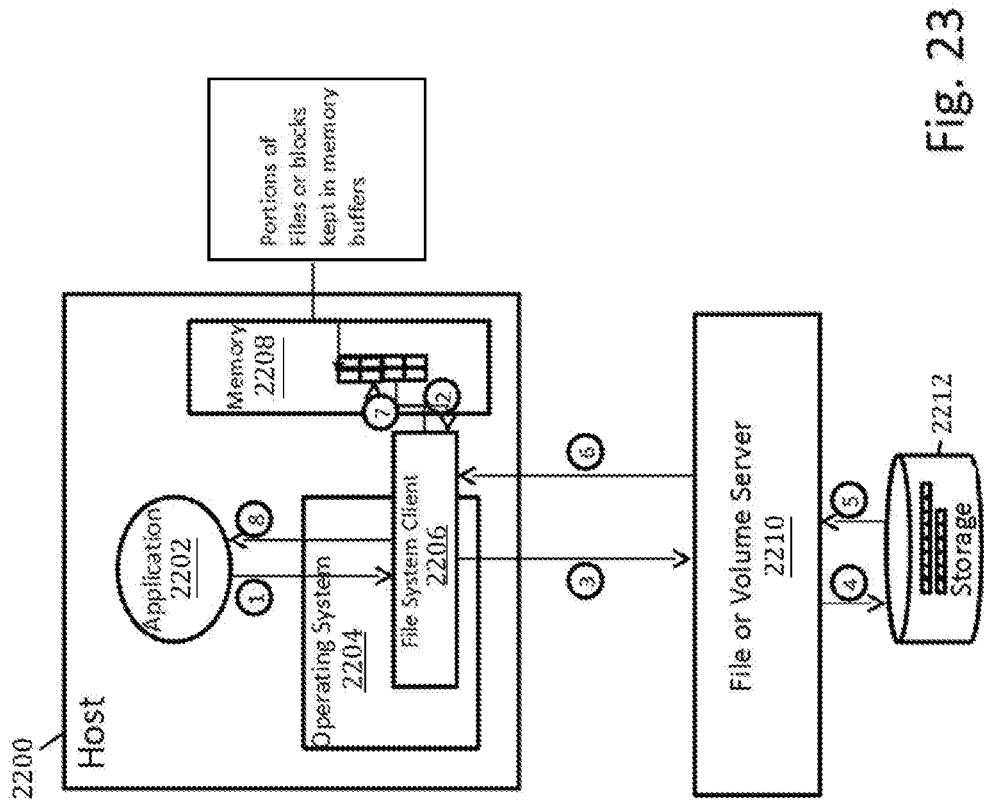
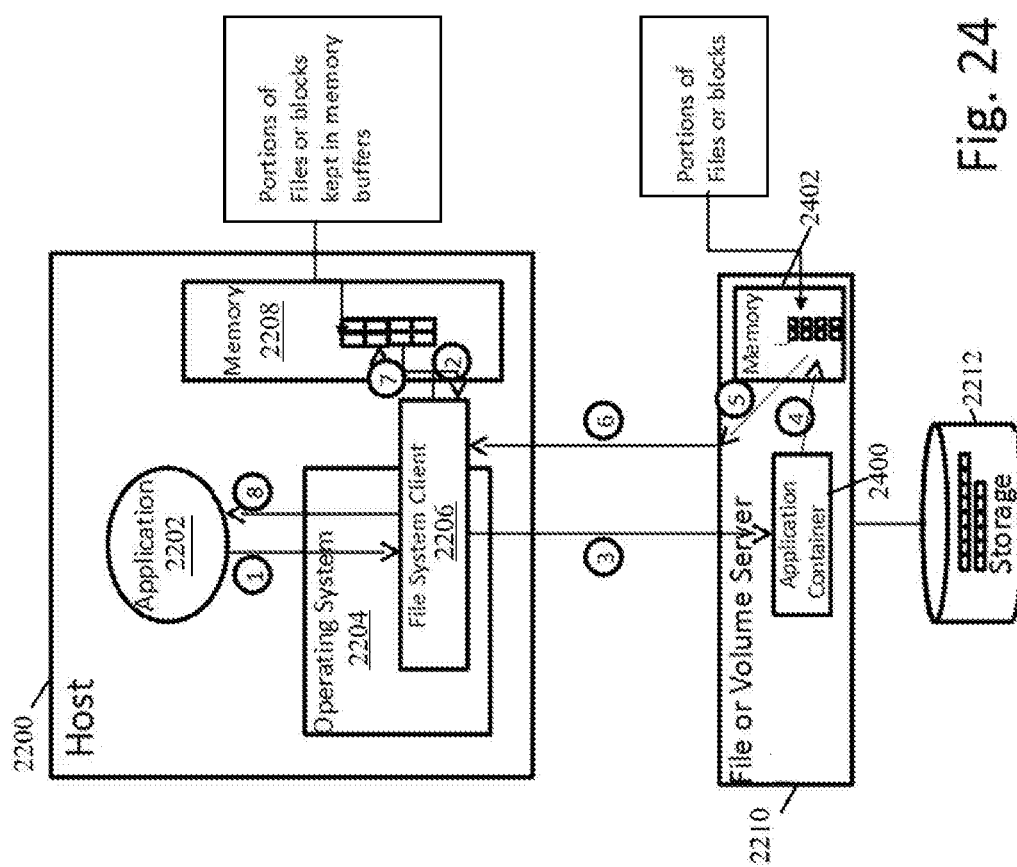


Fig. 23





24  
00  
L

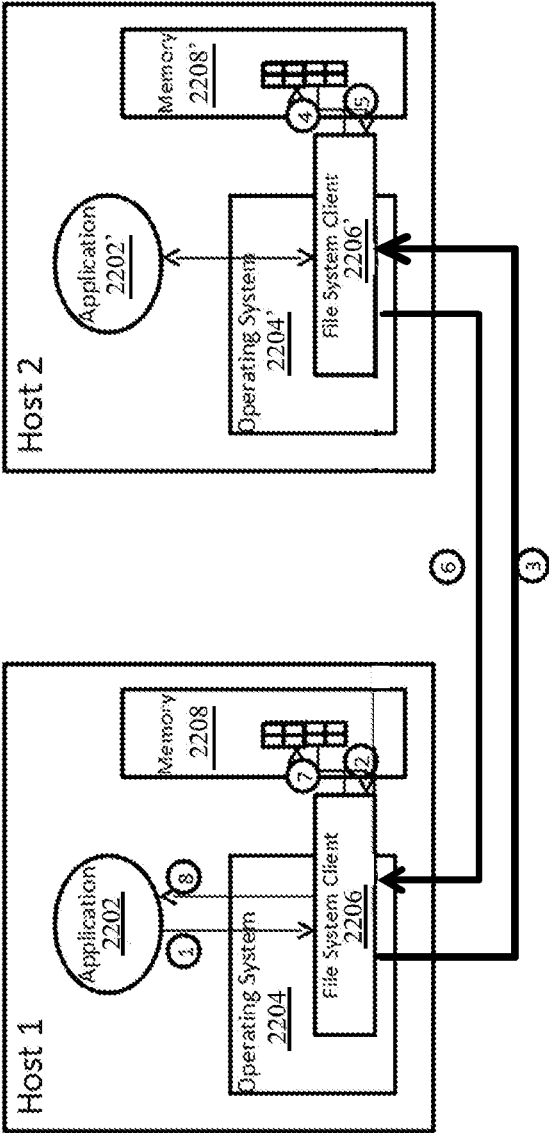


Fig. 25

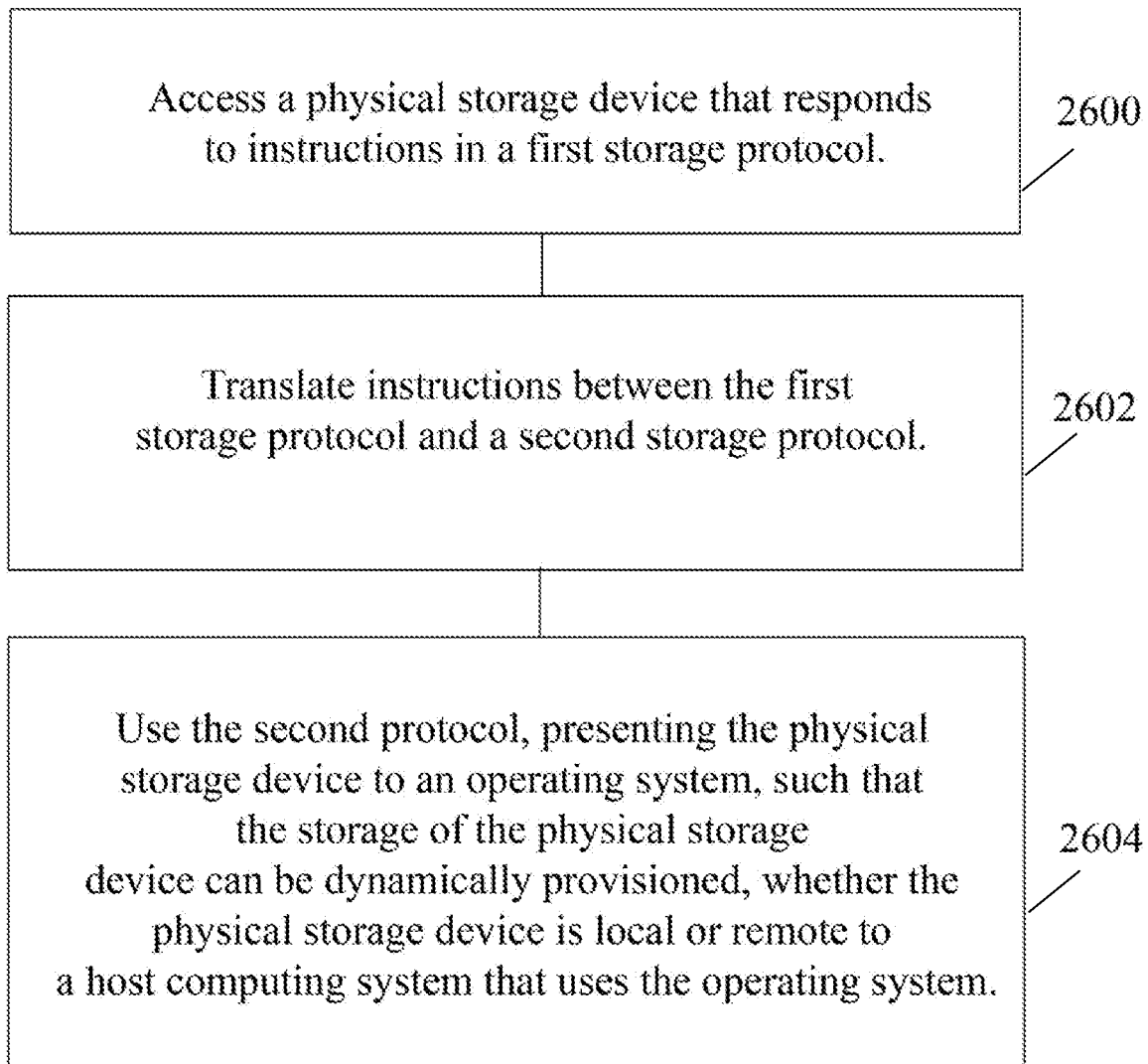


Fig. 26

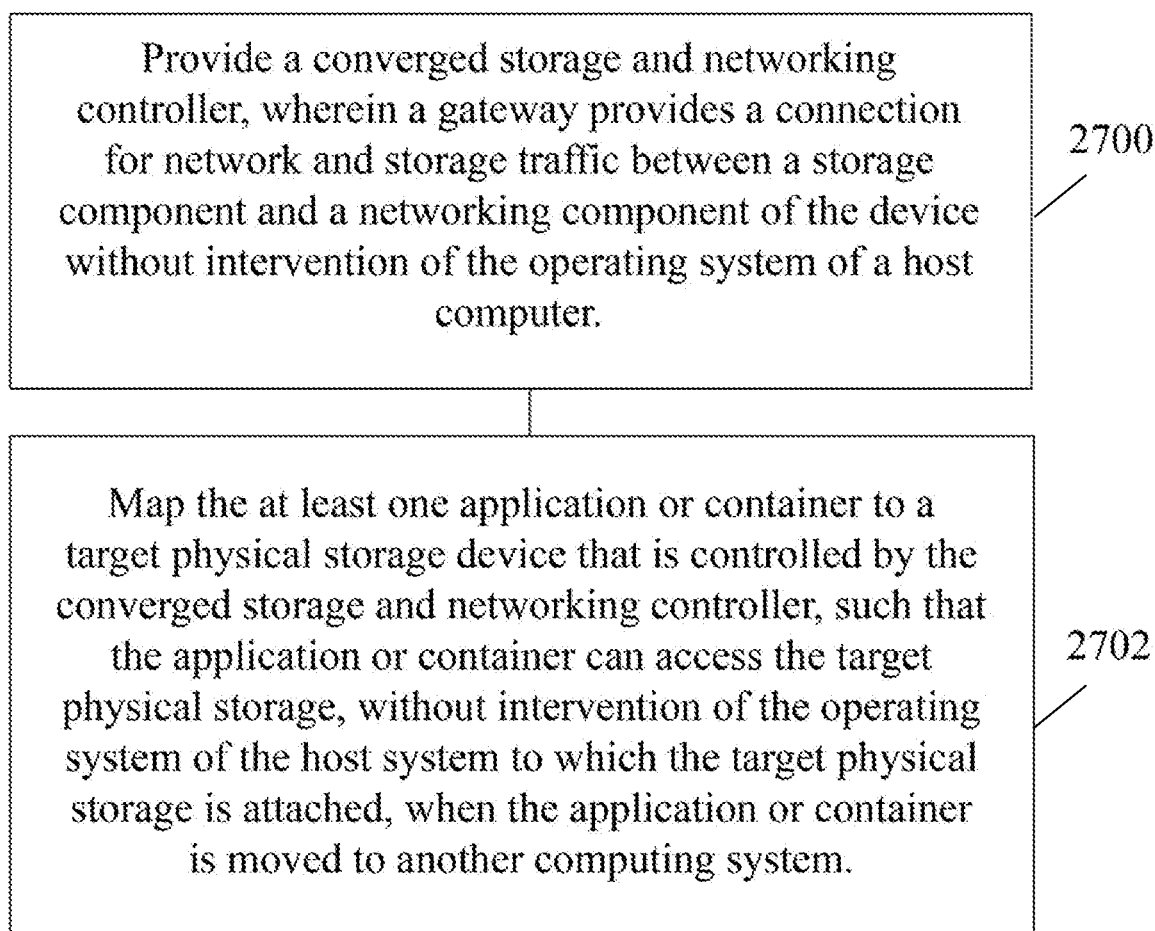


Fig. 27

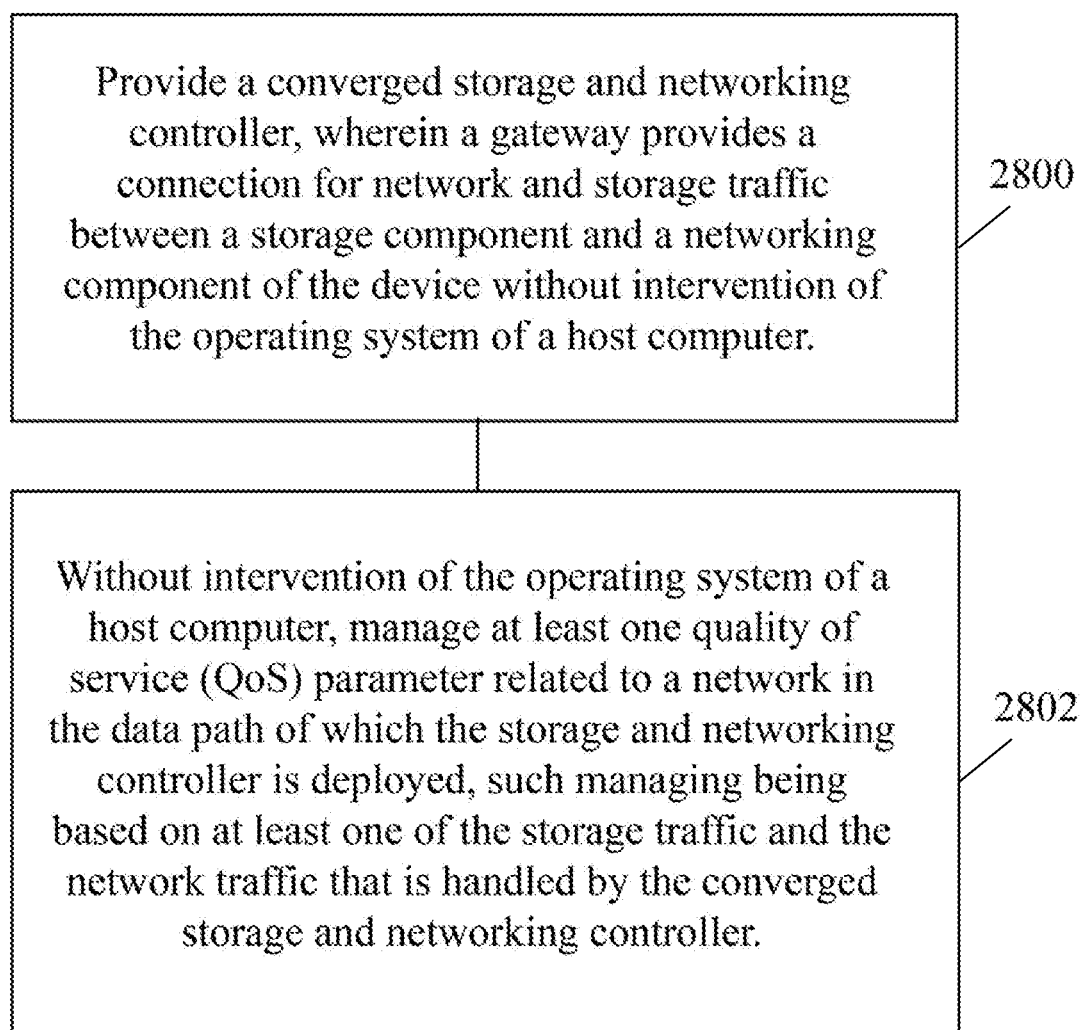


Fig. 28

## GUARANTEEING AVAILABILITY OF TARGET DATA TO REMOTE INITIATORS VIA A HYBRID SOURCE/TARGET CREDIT SCHEME

### CROSS REFERENCE TO RELATED APPLICATIONS

**[0001]** This application is a continuation of U.S. Ser. No. 14/640,717, filed Mar. 6, 2015 (DWIS-0004-U01) and entitled “METHODS AND SYSTEMS FOR CONVERGED NETWORKING AND STORAGE”, which is hereby incorporated by reference in its entirety.

**[0002]** U.S. Ser. No. 14/640,717 claims the benefit of the following provisional applications, each of which is hereby incorporated by reference in its entirety: U.S. patent application 61/950,036, filed Mar. 8, 2014 (DWIS-0002-P01) and entitled “METHOD AND APPARATUS FOR APPLICATION DRIVEN STORAGE ACCESS”; and U.S. patent application 62/017,257, filed Jun. 26, 2014 (DWIS-0003-P01) and entitled “APPARATUS FOR VIRTUALIZED CLUSTER IO”.

### FIELD OF THE INVENTION

**[0003]** This application relates to the fields of networking and data storage, and more particularly to the field of converged networking and data storage devices.

### BACKGROUND OF THE INVENTION

**[0004]** The proliferation of scale-out applications has led to very significant challenges for enterprises that use such applications. Enterprises typically choose between solutions like virtual machines (involving software components like hypervisors and premium hardware components) and so-called “bare metal” solutions (typically involving use of an operating system like Linux™ and commodity hardware). At large scale, virtual machine solutions typically have poor input-output (IO) performance, inadequate memory, inconsistent performance, and high infrastructure cost. Bare metal solutions typically have static resource allocation (making changes in resources difficult and resulting in inefficient use of the hardware), challenges in planning capacity, inconsistent performance, and operational complexity. In both cases, inconsistent performance characterizes the existing solutions. A need exists for solutions that provide high performance in multi-tenant deployments, that can handle dynamic resource allocation, and that can use commodity hardware with a high degree of utilization.

**[0005]** FIG. 1 depicts the general architecture of a computing system **102**, such as a server, functions and modules of which may be involved in certain embodiments disclosed herein. Storage functions (such as access to local storage devices on the server **102**, such as media **104** (e.g., rotating media or flash) and network functions such as forwarding have traditionally been performed separately in either software stacks or hardware devices (e.g., involving a network interface controller **118** or a storage controller **112**, for network functions or storage functions, respectively). Within an operating system stack **108** (which may include an operating system and a hypervisor in some embodiments including all the software stacks associated with storage and networking functions for the computing system), the software storage stack typically includes modules enabling use of various protocols that can be used in storage, such as the

small computer system interface (SCSI) protocol, the serial ATA (SATA) protocol, the non-volatile memory express (NVMe) protocol (a protocol for accessing disk-attached storage (DAS), like solid-state drives (SSDs), through the PCI Express (PCIe) bus **110** of a typical computing system **102**) or the like. The PCIe bus **110** may provide an interconnection between a CPU **106** (with processor(s) and memory) and various IO cards. The storage stack also may include volume managers, etc. Operations within the storage software stack may also include data protection, such as mirroring or RAID, backup, snapshots, deduplication, compression and encryption. Some of the storage functions may be offloaded into a storage controller **112**. The software network stack includes modules, functions and the like for enabling use of various networking protocols, such as Transmission Control Protocol/Internet Protocol (TCP/IP), the domain name system protocol (DNS), the address resolution protocol (ARP), forwarding protocols, and the like. Some of the network functions may be offloaded into a network interface controller **118** (or NIC) or the network fabric switch, such as via an ethernet connection **120**, in turn leading to a network (with various switches, routers and the like). In virtualized environments, a NIC **118** may be virtualized into several virtual NICs as specified by SR-IOV under the PCI Express standard. Although not specified by the PCI Express standard and not as common, storage controllers can also be virtualized in a similar manner. This approach allows virtual entities, such as virtual machines, access to their own private resource.

**[0006]** Referring to FIG. 2, one major problem with hypervisors is with the complexity of IO operations. For example, in order to deal with an operation involving data across two different computers (computer system **1** and computer system **2** in FIG. 2), data must be copied repeatedly, over and over, as it moves among the different software stacks involved in local storage devices **104**, storage controllers **112**, the CPUs **106**, network interface controller **118** and the hypervisor/operating systems **108** of the computers, resulting in large numbers of inefficient data copies for each IO operation whenever an activity is undertaken that involves moving data from one computer to another, changing the configuration of storage, or the like. The route **124** is one of many examples of the complex routes that data may take from one computer to another, moving up and down the software stacks of the two computers. Data that is sought by computing system **2** may be initially located in a local storage device **104**, such as a disk, of computing system **1**, then pulled by a storage controller card **112** (involving an IO operation and copying), send over the PCIe bus **110** (another IO operation) to the CPU **108** where it is handled by a hypervisor or other software component of the OS stack **108** of computing system **1**. Next, the data may be delivered (another IO operation) through the network controller **118** and over the network **122** (another set of IO operations) to computing system **2**. The route continues on computing system **2**, where data may travel through the network controller **118** and to the CPU **106** of computing system **2** (involve additional IO operations), then sent over the PCIe bus **110** to the local storage controller **112** for storage, then back to the hypervisor/OS stack **108** for actual use. These operations may occur across a multiplicity of pairs of computing systems, with each exchange involving this kind of proliferation of IO operations (and many other routes are possible, each involving significant numbers of operations).

Many such complex data replication and transport activities among computing systems are required in scaleout situations, which are increasingly adopted by enterprises. For example, when implementing a scaleout application like MongoDB™, customers must repeatedly run real time queries during rebalancing operations, and perform large scale data loading. Such activities involve very large numbers of IO operations, which result in poor performance in hypervisor solutions. Users of those applications also frequently re-shard (change the shards on which data is deployed), resulting in big problems for bare metal solutions that have static storage resource allocations, as migration of data from one location to another also involves many copying and transport operations, with large numbers of IO operations. As the amount of data used in scaleout applications grows rapidly, and the connectedness among disparate systems increases (such as in cloud deployments involving many machines), these problems grow exponentially. A need exists for storage and networking solutions that reduce the number and complexity of IO operations and otherwise improve the performance and scalability of scaleout applications without requiring expensive, premium hardware.

**[0007]** Referring still to FIG. 2, for many applications and use cases, data (and in turn, storage) needs to be accessed across the network between computing systems **102**. Three high-level steps of this operation include the transfer of data from the storage media of one computing system out of a box, movement across the network **122**, and the transfer of data into a second box (second computing system **102**) to the storage media **104** of that second computing system **102**. First, out of the box transfer, may involve intervention from the storage controller **112**, the storage stack in the OS **108**, the network stack in the OS **108**, and the network interface controller **118**. Many traversals and copying across internal busses (PCIe **110** and memory) as well as CPU **106** processing cycles are spent. This not only degrades performance (creating latency and throughput issues) of the operation, but also adversely affects other applications that run on the CPU. Second, once the data leaves the box, **102** and moves onto the network **122**, it is treated like any other network traffic and needs to be forwarded/routed to its destination. Policies are executed and decisions are made. In environments where a large amount of traffic is moving, congestion can occur in the network **122**, causing degradation in performance as well as problems with availability (e.g., dropped packets, lost connections, and unpredictable latencies). Networks have mechanisms and algorithms to avoid spreading of congestion, such as pause functions, backward congestion notification (BCN), explicit congestion notification (ECN), etc. However, these are reactive methods; that is, they detect formation of congestion points and push back on the source to reduce congestion, potentially resulting in delays and performance impacts. Third, once the data arrives at its “destination” computing system **102**, it needs to be processed, which involves intervention from the network interface controller **118**, the network stack in the OS **108**, the storage stack in the OS **108**, and the storage controller **112**. As with out of the box operations noted above, many traversals and copying across internal busses as well as CPU **106** processing cycles are spent. Further, the final destination of the data may well reside in still a different box. This can be the result of a need for more data protection (e.g., mirroring or across-box RAID) or the need for de-duplication. If so, then the entire sequence of

out-of-the box, across the network, and into the box data transfer needs to be repeated again. As described, limitations of this approach include degradation in raw performance, unpredictable performance, impact on other tenants or operations, availability and reliability, and inefficient use of resources. A need exists for data transfer systems that avoid the complexity and performance impacts of the current approaches.

**[0008]** As an alternative to hypervisors (which provide a separate operating system for each virtual machine that they manage), technologies such as Linux™ containers have been developed (which enable a single operating system to manage multiple application containers). Also, tools such as Dockers have been developed, which provide provisioning for packaging applications with libraries. Among many other innovations described throughout this disclosure, an opportunity exists for leveraging the capabilities of these emerging technologies to provide improved methods and systems for scaleout applications.

## SUMMARY

**[0009]** Provided herein are methods and systems that include a converged storage and network controller in hardware that combines initiator, target storage functions and network functions into a single data and control path, which allows a “cut-through” path between the network and storage, without requiring intervention by a host CPU. For ease of reference, this is referred to variously in this disclosure as a converged hardware solution, a converged device, a converged adaptor, a converged IO controller, a “datawise” controller, or the like throughout this disclosure, and such terms should be understood to encompass, except where context indicates otherwise, a converged storage and network controller in hardware that combines target storage functions and network functions into a single data and control path.

**[0010]** Among other benefits, the converged solution will increase raw performance of a cluster of computing and/or storage resources; enforce service level agreements (SLAs) across the cluster and help guarantee predictable performance; provide a multi-tenant environment where a tenant will not affect its neighbor; provide a denser cluster with higher utilization of the hardware resulting in smaller data center footprint, less power, fewer systems to manage; provide a more scalable cluster; and pool storage resources across the cluster without loss of performance.

**[0011]** The various methods and systems disclosed herein provide high-density consolidation of resources required for scaleout applications and high performance multi-node pooling. These methods and systems provide a number of customer benefits, including dynamic cluster-wide resource provisioning, the ability to guarantee quality-of-service (QoS), Security, Isolation etc. on network and storage functions, and the ability to use shared infrastructure for production and testing/development.

**[0012]** Also provided herein are methods and systems to perform storage functions through the network and to virtualize storage and network devices for high performance and deterministic performance in single or multi-tenant environments.

**[0013]** Also provided herein are methods and systems for virtualization of storage devices, such as those using NVMe

and similar protocols, and the translation of those virtual devices to different physical devices, such as ones using SATA.

**[0014]** The methods and systems disclosed herein also include methods and systems for end-to-end congestion control involving only the hardware on the host (as opposed to the network fabric) that includes remote credit management and a distributed scheduling algorithm at the box level.

**[0015]** Also provided herein are various methods and systems that are enabled by the converged network/storage controller, including methods and systems for virtualization of a storage cluster or of other elements that enable a cluster, such as a storage adaptor, a network adaptor, a container (e.g., a Linux container), a Solaris zone or the like. Among advantages, one aspect of virtualizing a cluster is that containers can become location-independent in the physical cluster. Among other advantages, this allows movement of containers among machines in a vastly simplified process described below.

**[0016]** Provided herein are methods and systems for virtualizing direct-attached storage (DAS), so that the operating system stack **108** still sees a local, persistent device, even if the physical storage is moved and is remotely located; that is, provided herein are methods and systems for virtualization of DAS. In embodiments this may include virtualizing DAS over a fabric, that is, taking a DAS storage system and moving it outside the box and putting it on the network. In embodiments this may include carving DAS into arbitrary name spaces. In embodiments the virtualized DAS is made accessible as if it were actual DAS to the operating system, such as being accessible by the OS **108** over a PCIe bus via NVMe. Thus, provided herein is the ability to virtualize storage (including DAS) so that the OS **108** sees it as DAS, even if the storage is actually accessed over a network protocol such as Ethernet, and the OS **108** is not required to do anything different than would be required with local physical storage.

**[0017]** Provided herein are methods and systems for providing DAS across a fabric, including exposing virtualized DAS to the OS **108** without requiring any modification of the OS **108**.

**[0018]** Also provided herein are methods and systems for virtualization of a storage adaptor (referring to a target storage system).

**[0019]** Provided herein are methods and systems for combining storage initiation and storage targeting in a single hardware system. In embodiments, these may be attached by a PCIe bus **110**. A single root virtualization function (SR-IOV) may be applied to take any standard device and have it act as if it is hundreds of such devices. Embodiments disclosed herein include using SR-IOV to give multiple virtual instances of a physical storage adaptor. SR-IOV is a PCIe standard that virtualizes I/O functions, and while it has been used for network interfaces, the methods and systems disclosed herein extend it to use for storage devices. Thus, provided herein is a virtual target storage system.

**[0020]** Embodiments may include a switch form factor or network interface controller, wherein the methods and systems disclosed herein may include a host agent (either in software or hardware). Embodiments may include breaking up virtualization between a front end and a back end.

**[0021]** Embodiments may include various points of deployment for a converged network and target storage controller. While some embodiments locate the converged

device on a host computing system **102**, in other cases the disk can be moved to another box (e.g., connected by Ethernet to a switch that switches among various boxes below. While a layer may be needed to virtualize, the storage can be separated, so that one can scale storage and computing resources separately. Also, one can then enable blade servers (i.e., stateless servers). Installations that would have formerly involved expensive blade servers and attached to storage area networks (SANs) can instead attach to the switch. In embodiments this comprises a “rackscale” architecture where resources are disaggregated at the rack level.

**[0022]** Methods and systems disclosed herein include methods and systems for virtualizing various types of non-DAS storage as DAS in a converged networking/target storage appliance. In embodiments, one may virtualize whatever storage is desired as DAS, using various front end protocols to the storage systems while exposing storage as DAS to the OS stack **108**.

**[0023]** Methods and systems disclosed herein include virtualization of a converged network/storage adaptor. From a traffic perspective, one may combine systems into one. Combining the storage and network adaptors, and adding in virtualization, gives significant advantages. Say there is a single host **102** with two PCIe buses **110**. To route from the PCIe **110**, you can use a system like RDMA to get to another machine/host **102**. If one were to do this separately, one has to configure the storage and the network RDMA system separately. One has to join each one and configure them at two different places. In the converged scenario, the whole step of setting up QoS, seeing that this is RDMA and that there is another fabric elsewhere is a zero touch process, because with combined storage and networking the two can be configured in a single step. That is, once one knows the storage, one doesn't need to set up the QoS on the network separately.

**[0024]** Method and systems disclosed herein include virtualization and/or indirection of networking and storage functions, embodied in the hardware, optionally in a converged network adaptor/storage adaptor appliance. While virtualization is a level of indirection, protocol is another level of indirection. The methods and systems disclosed herein may convert a protocol suitable for use by most operating systems to deal with local storage, such as NVMe, to another protocol, such as SAS, SATA, or the like. One may expose a consistent interface to the OS **108**, such as NVMe, and in the back end one may convert to whatever storage media is cost-effective. This gives a user a price/performance advantage. If components are cheaper/faster, one can connect any one of them. The back end could be anything, including NVMe.

**[0025]** Provided herein are methods and systems that include a converged data path for network and storage functions in an appliance. Alternative embodiments may provide a converged data path for network and storage functions in a switch.

**[0026]** In embodiments, methods and systems disclosed herein include storage/network tunneling, wherein the tunneling path between storage systems over a network does not involve the operating system of a source or target computer. In conventional systems, one had separate storage and network paths, so accessing storage remotely, required extensive copying to and from memory, I/O buses, etc. Merging the two paths means that storage traffic is going straight onto the network. The OS **108** of each computer sees



only a local disk. Another advantage is simplicity of programming. A user does not need to separately program a SAN, meaning that the methods disclosed herein include a one-step programmable SAN. Rather than requiring discovery and specification of zones, and the like, encryption, attachment, detachment and the like may be centrally, and programmatically done.

**[0027]** Embodiments disclosed herein may include virtualizing the storage to the OS **108** so that the OS **108** sees storage as a local disk. The level of indirection involved in the methods and systems disclosed herein allows the converged system to hide not only the location, but the media type, of storage media. All the OS sees is that there is a local disk, even if the actual storage is located remotely and/or is or a different type, such as a SAN. Thus, virtualization of storage is provided, where the OS **108** and applications do not have to change. One can hide all of the management, policies of tiering, policies of backup, policies of protection and the like that are normally needed to configure complex storage types behind.

**[0028]** Methods and systems are provided for selecting where indirection occurs in the virtualization of storage. Virtualization of certain functions may occur in hardware (e.g., in an adaptor on a host, in a switch, and in varying form factors (e.g., FPGA or ASICs) and in software. Different topologies are available, such as where the methods and systems disclosed herein are deployed on a host machine, on a top of the rack switch, or in a combination thereof. Factors that go into the selection include ease of use. Users who want to run stateless servers may prefer a top of rack. Ones who don't care about that approach might prefer the controller on the host.

**[0029]** Methods and systems disclosed herein include providing NVMe over Ethernet. These approaches can be the basis for the tunneling protocol that is used between devices. NVMe is a suitable DAS protocol that is intended conventionally to go to a local PCIe. Embodiments disclosed herein may tunnel the NVMe protocol traffic over Ethernet. NVMe (non-volatile memory express) is a protocol that in Linux and Windows provides access to PCIe-based Flash Storage. This provides high performance by by-passing the software stacks used in conventional systems.

**[0030]** Embodiments disclosed herein may include providing an NVMe device that is virtualized and dynamically allocated. In embodiments one may piggy back NVMe, but carve up and virtualize and dynamically allocate an NVMe device. In embodiments there is no footprint in the software. The operating system stays the same (just a small driver that sees the converged network/storage card). This results in virtual storage presented like a direct attached disk, but the difference is that now we can pool such devices across the network.

**[0031]** Provided herein are methods and systems for providing the simplicity of direct attached storage (DAS) with the advantages of sharing like in a storage area network (SAN). Each converged appliance in various embodiments disclosed herein may be a host, and any storage drives may be local to a particular host but seen by the other hosts (as in a SAN or other network-accessible storage). The drives in each box enabled by a network/storage controller of the present disclosure behave like a SAN (that is, are available on the network), but the management methods are much simpler. When a storage administrator sets up a SAN, a typical enterprise may have a whole department setting up

zones for a SAN (e.g., a fiber channel switch), such as setting up "who sees what." That knowledge is pre-loaded and a user has to ask the SAN administrator to do the work to set it up. There is no programmability in a typical legacy SAN architecture. The methods and systems disclosed herein provide local units that are on the network, but the local units can still access their storage without having to go through complex management steps like zone definition, etc. These devices can do what a SAN does just by having both network and storage awareness. As such, they represent the first programmatic SAN.

**[0032]** Methods and systems disclosed herein may include persistent, stateful, disaggregated storage enabled by a hardware appliance that provides converged network and storage data management.

**[0033]** Methods and systems disclosed herein may also include convergence of network and storage data management in a single appliance, adapted to support use of containers for virtualization. Such methods and systems are compatible with the container ecosystem that is emerging, but offering certain additional advantages.

**[0034]** Methods and systems are disclosed herein for implementing virtualization of NVMe. Regardless how many sources to how many destinations, as long as the data from the sources is serialized first before going into the hub, then the hub distributes to data to the designated destination sequentially. If so, then data transport resources such as DMA engine can be reduced to only one copy. This may include various use scenarios, in one scenario, for NVMe virtual functions (VFs), if they are all connected to the same PCIe bus, then regardless how many VFs are configured, the data would be coming into this pool of VFs serially, so there is only one DMA engine and only one storage block (for control information) is needed. In another use scenario, for a disk storage system with a pool of discrete disks/controllers, if the data is originated from the physical bus, i.e. PCIe, since the data is serially coming into this pool of disks, then regardless how many disks/controllers are in the pool, the transport resources such as the DMA engine can be reduced to only one instead of one per controller.

**[0035]** In accordance with various exemplary and non-limiting embodiments, a device comprises a converged input/output controller that includes a physical target storage media controller, a physical network interface controller; and a gateway between the storage media controller and the network interface controller, wherein gateway provides a direct connection for storage traffic and network traffic between the storage media controller and the network interface controller.

**[0036]** In accordance with various exemplary and non-limiting embodiments, a method of virtualization of a storage device comprises accessing a physical storage device that responds to instructions in a first storage protocol, translating instructions between the first storage protocol and a second storage protocol and using the second protocol, presenting the physical storage device to an operating system, such that the storage of the physical storage device can be dynamically provisioned, whether the physical storage device is local or remote to a host computing system that uses the operating system.

**[0037]** In accordance with various exemplary and non-limiting embodiments, a method of facilitating migration of at least one of an application and a container comprises providing a converged storage and networking controller,

wherein a gateway provides a connection for network and storage traffic between a storage component and a networking component of the device without intervention of the operating system of a host computer and mapping the at least one application or container to a target physical storage device that is controlled by the converged storage and networking controller, such that the application or container can access the target physical storage, without intervention of the operating system of the host system to which the target physical storage is attached, when the application or container is moved to another computing system.

**[0038]** In accordance with various exemplary and non-limiting embodiments, a method of providing quality of service (QoS) for a network, comprises providing a converged storage and networking controller, wherein a gateway provides a connection for network and storage traffic between a storage component and a networking component of the device without intervention of the operating system, a hypervisor, or other software running on the CPU of a host computer and, also without intervention of the operating system, hypervisor, or other software running on the CPU of a host computer, managing at least one quality of service (QoS) parameter related to a network in the data path of which the storage and networking controller is deployed, such managing being based on at least one of the storage traffic and the network traffic that is handled by the converged storage and networking controller.

**[0039]** QoS may be based on various parameters, such as one or more of a bandwidth parameter, a network latency parameter, an IO performance parameter, a throughput parameter, a storage type parameter and a storage latency parameter. QoS may be maintained automatically when at least one of an application and a container that is serviced by storage through the converged storage and network controller is migrated from a host computer to another computer. Similarly, QoS may be maintained automatically when at least one target storage device that services at least one of an application and a container through the converged storage and network controller is migrated from a first location to another location or multiple locations. For example, storage may be scaled, or different storage media types may be selected, to meet storage needs as requirements are increased. In embodiments, a security feature may be provided, such as encryption of network traffic data, encryption of data in storage, or both. Various storage features may be provided as well, such as compression, protection levels (e.g., RAID levels), use of different storage media types, global de-duplication, and snapshot intervals for achieving at least one of a recovery point objective (RPO) and a recovery time objective (RTO).

#### BRIEF DESCRIPTION OF THE FIGURES

**[0040]** The accompanying figures where like reference numerals refer to identical or functionally similar elements throughout the separate views and which together with the detailed description below are incorporated in and form part of the specification, serve to further illustrate various embodiments and to explain various principles and advantages all in accordance with the systems and methods disclosed herein.

**[0041]** FIG. 1 illustrates a general architecture in accordance with an exemplary and non-limiting embodiment.

**[0042]** FIG. 2 illustrates a computer system in accordance with an exemplary and non-limiting embodiment.

**[0043]** FIG. 3 illustrates a converged solution in accordance with an exemplary and non-limiting embodiment.

**[0044]** FIG. 4 illustrates two computing systems enabled by a converged solution in accordance with an exemplary and non-limiting embodiment.

**[0045]** FIG. 5 illustrates a converged controller in accordance with an exemplary and non-limiting embodiment.

**[0046]** FIG. 6 illustrates a deployment of a converged controller in accordance with an exemplary and non-limiting embodiment.

**[0047]** FIG. 7 illustrates a plurality of systems in accordance with an exemplary and non-limiting embodiment.

**[0048]** FIG. 8 illustrates a block diagram of a field-programmable gate array (FPGA) in accordance with an exemplary and non-limiting embodiment.

**[0049]** FIG. 9 illustrates an architecture of a controller card in accordance with an exemplary and non-limiting embodiment.

**[0050]** FIG. 10 illustrates a software stack in accordance with an exemplary and non-limiting embodiment.

**[0051]** FIGS. 11-15 illustrate the movement of an application container across multiple systems in accordance with an exemplary and non-limiting embodiment.

**[0052]** FIG. 16 illustrates packet transmission in accordance with an exemplary and non-limiting embodiment.

**[0053]** FIG. 17 illustrates a storage access scheme in accordance with an exemplary and non-limiting embodiment.

**[0054]** FIG. 18 illustrates the operation of a file system in accordance with an exemplary and non-limiting embodiment.

**[0055]** FIG. 19 illustrates the operation of a distributed file server in accordance with an exemplary and non-limiting embodiment.

**[0056]** FIG. 20 illustrates a high performance distributed file server (DFS) in accordance with an exemplary and non-limiting embodiment.

**[0057]** FIG. 21 illustrates a system in accordance with an exemplary and non-limiting embodiment.

**[0058]** FIG. 22 illustrates a host in accordance with an exemplary and non-limiting embodiment.

**[0059]** FIG. 23 illustrates an application accessing a block of data in accordance with an exemplary and non-limiting embodiment.

**[0060]** FIG. 24 illustrates an application accessing a block of data in accordance with an exemplary and non-limiting embodiment.

**[0061]** FIG. 25 illustrates a system in accordance with an exemplary and non-limiting embodiment.

**[0062]** FIG. 26 illustrates a method according to an exemplary and non-limiting embodiment.

**[0063]** FIG. 27 illustrates a method according to an exemplary and non-limiting embodiment.

**[0064]** FIG. 28 illustrates a method according to an exemplary and non-limiting embodiment.

**[0065]** Skilled artisans will appreciate that elements in the figures are illustrated for simplicity and clarity and have not necessarily been drawn to scale. For example, the dimensions of some of the elements in the figures may be exaggerated relative to other elements to help to improve understanding of embodiments of the systems and methods disclosed herein.

## DETAILED DESCRIPTION OF THE INVENTION

[0066] The present disclosure will now be described in detail by describing various illustrative, non-limiting embodiments thereof with reference to the accompanying drawings and exhibits. The disclosure may, however, be embodied in many different forms and should not be construed as being limited to the illustrative embodiments set forth herein. Rather, the embodiments are provided so that this disclosure will be thorough and will fully convey the concept of the disclosure to those skilled in the art. The claims should be consulted to ascertain the true scope of the disclosure.

[0067] Before describing in detail embodiments that are in accordance with the systems and methods disclosed herein, it should be observed that the embodiments reside primarily in combinations of method steps and/or system components related to converged networking and storage. Accordingly, the system components and method steps have been represented where appropriate by conventional symbols in the drawings, showing only those specific details that are pertinent to understanding the embodiments of the systems and methods disclosed herein so as not to obscure the disclosure with details that will be readily apparent to those of ordinary skill in the art.

[0068] Referring to FIG. 3, the converged solution 300 may include three important aspects and may be implemented in a hardware device that includes a combination of hardware and software modules and functions. First, a cut-through data path 304 may be provided between a network controller 118 and a storage controller 112, so that access of the storage to and from the network can be direct, without requiring any intervention of the OS stack 108, the PCIe bus 110, or the CPU 106. Second, cut through storage stack access, such as to storage devices 302, may be provided, such as access of the storage to and from entities on the local host, which allows bypassing of complex legacy software stacks for storage access, such as SCSI/SAS/SATA stacks. Third, end-to-end congestion management and flow control of the network may be provided, such as by a mechanism to reserve and schedule the transfer of data across the network, which guarantees the availability of the target's data to remote initiators and minimizes the congestion of the traffic as it flows through intermediate network fabric switches. The first and second aspects remove software stacks (hence the CPU 106 and memory) from the path of the data, eliminating redundant or unnecessary movement and processing. End-to-end congestion management and flow control delivers a deterministic and reliable transport of the data.

[0069] As noted above, one benefit of the converged solution 300 is that the operating system stack 108 connects to the converged solution 300 over a conventional PCIe 110 or a similar bus, so that the OS stack 108 sees the converged solution 300, and any storage that it controls through the cut-through to storage devices 302, as one or more local, persistent devices, even if the physical storage is remotely located. Among other things, this comprises the capability for virtualization of DAS 308, which may include virtualizing DAS 308 over a fabric, that is, taking a DAS 308 storage system and moving it outside the computing system 102 and putting it on the network. The storage controller 112 of the converged solution 300 may connect to and control DAS 308 on the network 122 via various known protocols,

such as SAS, SATA, or NVMe. In embodiments virtualization may include carving DAS 308 into arbitrary name spaces. In embodiments the virtualized DAS 308 is made accessible as if it were actual, local, physical DAS to the operating system, such as being accessible by the OS 108 over a PCIe bus 110 to the storage controller 112 of the converged solution 300 via a standard protocol such as NVMe. Again, the OS 108 sees the entire solution 300 as a local, physical device, such as DAS. Thus, provided herein is the ability to virtualize storage (including DAS and other storage types, such as SAN 310) so that the OS 108 sees any storage type as DAS, even if the storage is actually accessed over a network 122, and the OS 108 is not required to do anything different than would be required with local physical storage. In the case where the storage devices 302 are SAN 310 storage, the storage controller 112 of the converged solution may control the SAN 310 through an appropriate protocol used for storage area networks, such as the Internet Small Computing System Interface (iSCSI), Fibre Channel (FC), or Fibre Channel over Ethernet (FCoE). Thus, the converged solution 300 provides a translation for the OS stack 108 from any of the other protocols used in storage, such as Ethernet, SAS, SATA, NVMe, iSCSI, FC or FCoE, among others, to a simple protocol like NVMe that makes the disparate storage types and protocols appear as local storage accessible over PCIe 110. This translation in turns enables virtualization of a storage adaptor (referring to any kind of target storage system). Thus, methods and systems disclosed herein include methods and systems for virtualizing various types of non-DAS storage as DAS in a converged networking/target storage appliance 300. In embodiments, one may virtualize whatever storage is desired as DAS, using various protocols to the storage systems while exposing storage as DAS to the OS stack 108. Thus, provided herein are methods and systems for virtualization of storage devices, such as those using NVMe and similar protocols, and the translation of those virtual devices to different physical devices, such as ones using SATA.

[0070] Storage/network tunneling 304, where the tunneling path between storage systems over the network 122 does not involve the operating system of a source or target computer enables a number of benefits. In conventional systems, one has separate storage and network paths, so accessing storage remotely required extensive copying to and from memory, I/O buses, etc. Merging the two paths means that storage traffic is going straight onto the network. An advantage is simplicity of programming. A user does not need to separately program a SAN 310, meaning that the methods disclosed herein enable a one-step programmable SAN 310. Rather than requiring discovery and specification of zones, and the like, configuration, encryption, attachment, detachment and the like may be centrally, and programmatically done. As an example, a typical SAN is composed of "initiators," "targets," and a switch fabric, which connects the initiators and targets. Typically which initiators see which targets are defined/controlled by the fabric switches, called "zones." Therefore, if an initiator moves or a target moves, zones need to be updated. The second control portion of a SAN typically lies with the "targets." They can control which initiator port can see what logical unit numbers (LUNs) (storage units exposed by the target). This is typically referred to as LUN masking and LUN mapping. Again, if an initiator moves locations, one has to re-program the "Target". Consider now that in such an environment if an

application moves from one host to another (such as due to a failover, load re-balancing, or the like) the zoning and LUN masking/mapping needs to be updated. Alternatively, one could pre-program the SAN, so that every initiator sees every target. However, doing so results in an un-scalable and un-secure SAN. In the alternate solution described throughout this disclosure, such a movement of an application, a container, or a storage device does NOT require any SAN re-programming, resulting in a zero touch solution. The mapping maintained and executed by the converged solution 300 allows an application or a container, the target storage media, or both, to be moved (including to multiple locations) and scaled independently, without intervention by the OS, a hypervisor, or other software running on the host CPU.

[0071] The fact that the OS 108 sees storage as a local disk allows simplified virtualization of storage. The level of indirection involved in the methods and systems disclosed herein allows the converged system 300 to hide not only the location, but the media type, of storage media. All the OS 108 sees is that there is a local disk, even if the actual storage is located remotely and/or is of a different type, such as a SAN 310. Thus, virtualization of storage is provided through the converged solution 300, where the OS 108 and applications do not have to change. One can hide all of the management, policies of tiering, policies of backup, policies of protection and the like that are normally needed to configure complex storage types behind.

[0072] The converged solution 300 enables the simplicity of direct attached storage (DAS) with the advantages of a storage area network (SAN). Each converged appliance 300 in various embodiments disclosed herein may act as a host, and any storage devices 302 may be local to a particular host but seen by the other hosts (as is the case in a SAN 310 or other network-accessible storage). The drives in each box enabled by a network/storage controller of the present disclosure behave like a SAN 310 (e.g., are available on the network), but the management methods are much simpler. When a storage administrator normally sets up a SAN 310, a typical enterprise may have a whole department setting up zones for a SAN 310 (e.g., a fiber channel switch), such as setting up “who sees what.” That knowledge must be pre-loaded, and a user has to ask the SAN 310 administrator to do the work to set it up. There is no programmability in a typical legacy SAN 310 architecture. The methods and systems disclosed herein provide local units that are on the network, but the local units can still access their storage without having to go through complex management steps like zone definition, etc. These devices can do what a SAN does just by having both network and storage awareness. As such, they represent the first programmatic SAN.

[0073] The solution 300 can be described as a “Converged IO Controller” that controls both the storage media 302 and the network 122. This converged controller 300 is not just a simple integration of the storage controller 112 and the network controller (NIC) 118. The actual functions of the storage and network are merged such that storage functions are performed as the data traverses to and from the network interface. The functions may be provided in a hardware solution, such as an FPGA (one or more) or ASIC (one or more) as detailed below.

[0074] Referring to FIG. 4, two or more computing systems 102 enabled by converged solutions 300 may serve as hosts for respective storage targets, where by merging storage and network and controlling both interfaces, direct

access to the storage 302 can be achieved remotely over the network 122 without traversing internal busses or CPU/software work, such as by a point-to-point path 400 or by an Ethernet switch 402 to another computer system 102 that is enabled by a converged solution 300. The highest performance (high IOPs and low latency) can be achieved. Further, storage resources 302 can now be pooled across the cluster. In FIG. 4, this is conceptually illustrated by the dotted oval 404.

[0075] In embodiments, the converged solution 300 may be included on a host computing system 102, with the various components of a conventional computing system as depicted in FIG. 1, together with the converged IO controller 300 as described in connection with FIG. 3. Referring to FIG. 5, in alternative embodiments, the converged controller 300 may be disposed in a switch, such as a top of the rack switch, thus enabling a storage enabled switch 500. The switch may reside on the network 122 and be accessed by a network controller 118, such as of a conventional computing system 102.

[0076] Referring to FIG. 6, systems may be deployed in which a converged controller 300 is disposed both on one or more host computing systems 102 and on a storage enabled switch 500, which may be connected to systems 102 that are enabled by converged solutions 300 and to non-enabled systems 102. As noted above, target storage 302 for the converged controller(s) 300 on the host computing system 102 and on the storage enabled switch 500 can be visible to each other across the network, such as being treated as a unified resource, such as to virtualization solutions. In sum, intelligence, including handling converged network and storage traffic on the same device, can be located in a host system, in a switch, or both in various alternative embodiments of the present disclosure.

[0077] Embodiments disclosed herein may thus include a switch form factor or a network interface controller, or both which may include a host agent (either in software or hardware). These varying deployments allow breaking up virtualization capabilities, such as on a host and/or on a switch and/or between a front end and a back end. While a layer may be needed to virtualize certain functions, the storage can be separated, so that one can scale storage and computing resources separately. Also, one can then enable blade servers (i.e., stateless servers). Installations that would have formerly involved expensive blade servers and attached storage area networks (SANs) can instead attach to the storage enabled switch 500. In embodiments this comprises a “rackscale” architecture, where resources are disaggregated at the rack level.

[0078] Methods and systems are provided for selecting where indirection occurs in the virtualization of storage. Virtualization of certain functions may occur in hardware (e.g., in a converged adaptor 300 on a host 102, in a storage enabled switch 500, in varying hardware form factors (e.g., FPGAs or ASICs) and in software. Different topologies are available, such as where the methods and systems disclosed herein are deployed on a host machine 102, on a top of the rack switch 500, or in a combination thereof. Factors that go into the selection of where virtualization should occur include ease of use. Users who want to run stateless servers may prefer a top of rack storage enabled switch 500. Ones who don't care about that approach might prefer the converged controller 300 on the host 102.

[0079] FIG. 7 shows a more detailed view of a set of systems that are enabled with converged controllers 300, including two computer systems 102 (computer system 1 and computer system 2), as well as a storage enabled switch 500. Storage devices 302, such as DAS 308 and SAN 310 may be controlled by the converged controller 300 or the storage enabled switch 500. DAS 308 may be controlled in either case using SAS, SATA or NVMe protocols. SAN 310 may be controlled in either case using iSCSI, FC or FCoE. Connections among hosts 102 that have storage controllers 300 may be over a point-to-point path 400, over an Ethernet switch 402, or through a storage enabled switch 500, which also may provide a connection to a conventional computing system. As noted above, the multiple systems with intelligent converged controllers 300 can each serve as hosts and as storage target locations that the other hosts see, thereby providing the option to be treated as a single cluster of storage for purposes of an operating system 108 of a computing system 102.

[0080] Method and systems disclosed herein include virtualization and/or indirection of networking and storage functions, embodied in the hardware converged controller 300, optionally in a converged network adaptor/storage adaptor appliance 300. While virtualization is a level of indirection, protocol is another level of indirection. The methods and systems disclosed herein may convert a protocol suitable for use by most operating systems to deal with local storage, such as NVMe, to another protocol, such as SAS, SATA, or the like. One may expose a consistent interface to the OS 108, such as NVMe, and on the other side of the converged controller 300 one may convert to whatever storage media 302 is cost-effective. This gives a user a price/performance advantage. If components are cheaper/faster, one can connect any one of them. The side of the converged controller 300 could face any kind of storage, including NVMe. Furthermore the storage media type may be any of the following including, but not limited, to HDD, SSD (based on SLC, MLC, or TLC Flash), RAM etc or a combination thereof.

[0081] In embodiments, a converged controller may be adapted to virtualize NVMe virtual functions, and to provide access to remote storage devices 302, such as ones connected to a storage-enabled switch 500, via NVMe over an Ethernet switch 402. Thus, the converged solution 300 enables the use of NVMe over Ethernet 700, or NVMeoE. Thus, methods and systems disclosed herein include providing NVMe over Ethernet. These approaches can be the basis for the tunneling protocol that is used between devices, such as the host computing system 102 enabled by a converged controller 300 and/or a storage enabled switch 500. NVMe is a suitable DAS protocol that is intended conventionally to go to a local PCIe 110. Embodiments disclosed herein may tunnel the NVMe protocol traffic over Ethernet. NVMe (non-volatile memory express) is a protocol that in Linux and Windows provides access to PCIe-based Flash. This provides high performance via by-passing the software stacks used in conventional systems, while avoiding the need to translate from NVMe (as used by the OS stack 108) and the traffic tunneled over Ethernet to other devices.

[0082] FIG. 8 is a block diagram of an FPGA 800, which may reside on an IO controller card and enable an embodiment of a converged solution 300. Note that while a single FPGA 800 is depicted, the various functional blocks could

be organized into multiple FPGAs, into one or more customer Application Specific Integrated Circuits (ASICs), or the like. For example, various networking blocks and various storage blocks could be handled in separate (but interconnected) FPGAs or ASICs. References throughout this disclosure to an FPGA 800 should be understood, except where context indicates otherwise, to encompass these other forms of hardware that can enable the functional capabilities reflected in FIG. 8 and similar functions. Also, certain functional groups, such as for networking functions and/or storage functions, could be embodied in merchant silicon.

[0083] The embodiment of the FPGA 800 of FIG. 8 has four main interfaces. First, there is PCIe interface, such as to the PCIe bus 110 of a host computer 102. Thus, the card is a PCIe end point. Second, there is a DRAM/NVRAM interface. For example, a DDR interface may be provided to external DRAM or NVRAM, used by the embedded CPUs, meta-data and data structures, and packet/data buffering. Third, there is a storage interface to media, such as DAS 308 and SAN 310. Storage interfaces can include ones for SAS, SATA, NVMe, iSCSI, FC and/or FCoE, and could in embodiments be any interface to rotating media, flash, or other persistent form of storage, either local or over a cut-through to a network-enabled storage like SAN 310. Fourth, a network interface is provided, such as Ethernet to a network fabric. The storage interfaces and the network interfaces can be used, in part, to enable NVMe over Ethernet.

[0084] The internal functions of the FPGA 800 may include a number of enabling features for the converged solution 300 and other aspects of the present disclosure noted throughout. A set of virtual endpoints (vNVMe) 802 may be provided for the host. Analogous to the SR-IOV protocol that is used for the network interface, this presents virtual storage targets to the host. In this embodiment of the FPGA 800, NVMe has benefits of low software overhead, which in turn provides high performance. A virtual NVMe device 802 can be dynamically allocated/de-allocated/moved and resized. As with SR-IOV, there is one physical function (PF) 806 that interfaces with a PCIe driver 110 (see below), and multiple virtual functions 807 (VF) in which each appears as an NVMe device.

[0085] Also provided in the FPGA 800 functions are one or more read and write direct memory access (DMA) queues 804, referred to in some cases herein as a DMA engine 804. These may include interrupt queues, doorbells, and other standard functions to perform DMA to and from the host computing system 102.

[0086] A device mapping facility 808 on the FPGA 800 may determine the location of the virtual NVMe devices 802. The location options would be local (ie—attached to one of the storage media interfaces 824 shown), or remote on another host 102 of a storage controller 300. Access to a remote vNVMe device requires going through a tunnel 828 to the network 122.

[0087] A NVMe virtualization facility 810 may translate NVMe protocol instructions and operations to the corresponding protocol and operations of the backend storage media 302, such as SAS or SATA (in the case of use of NVMe on the backend storage medium 302, no translation may be needed) where DAS 308 is used, or such as iSCSI, FC or FCoE in the case where SAN 310 storage is used in the backend. References to the backend here refer to the other side of the converged controller 300 from the host 102.

[0088] A data transformation function **812** may format the data as it is stored onto the storage media **302**. These operations could include re-writes, transformation, compression, protection (such as RAID), encryption and other functions that involve changing the format of the data in any way as necessary to allow it to be handled by the applicable type of target storage medium **308**. In some embodiments, storage medium **308** may be remote.

[0089] In embodiments, storage read and write queues **814** may include data structures or buffering for staging data during a transfer. In embodiments, temporary memory, such as DRAM of NVRAM (which may be located off the FPGA **800**) may be used for temporary storage of data.

[0090] A local storage scheduler and shaper **818** may prioritize and control access to the storage media **302**. Any applicable SLA policies for local storage may be enforced in the scheduler and shaper **818**, which may include strict priorities, weighted round robin scheduling, IOP shapers, and policers, which may apply on a per queue, per initiator, per target, or per c-group basis, and the like.

[0091] A data placement facility **820** may implement an algorithm that determines how the data is laid out on the storage media **302**. That may involve various placement schemes known to those of skill in the art, such as striping across the media, localizing to a single device **302**, using a subset of the devices **302**, or localizing to particular blocks on a device **302**.

[0092] A storage metadata management facility **822** may include data structures for data placement, block and object i-nodes, compression, deduplication, and protection. Metadata may be stored either in off-FPGA **800** NVRAM/DRAM or in the storage media **302**.

[0093] A plurality of control blocks **824** may provide the interface to the storage media. These may include SAS, SATA, NVMe, PCIe, iSCSI, FC and/or FCoE, among other possible control blocks, in each case as needed for the appropriate type of target storage media **302**.

[0094] A storage network tunnel **828** of the FPGA **800** may provide the tunneling/cut-through capabilities described throughout this disclosure in connection with the converged solution **300**. Among other things, the tunnel **828** provides the gateway between storage traffic and network traffic. It includes encapsulation/de-encapsulation or the storage traffic, rewrite and formatting of the data, and end-to-end coordination of the transfer of data. The coordination may be between FPGAs **800** across nodes within a host computing system **102** or in more than one computing system **102**, such as for the point-to-point path **400** described in connection with FIG. 4. Various functions, such as sequence numbers, packet loss, time-outs, and retransmissions may be performed. Tunneling may occur over Ethernet, including by FCoE or NVMeoE.

[0095] A virtual network interface card facility **830** may include a plurality of SR-IOV endpoints to the host **102**, presented as virtual network interface cards. One physical function (PF) **836** may interface with a PCIe driver **110** (see software description below), and multiple virtual functions (VF) **837**, in which each appear as a network interface card (NIC) **118**.

[0096] A set of receive/transmit DMA queues **832** may include interrupt queues, doorbells, and other standard functions to perform DMA to and from the host **102**.

[0097] A classifier and flow management facility **834** may perform standard network traffic classification, typically to IEEE standard **802.1Q** class of service (COS) mappings or other priority levels.

[0098] An access control and rewrite facility **838** may handle access control lists (ACLs) and rewrite policies, including access control lists typically operating on Ethernet tuples (MAC SA/DA, IP SA/DA, TCP ports, etc.) to reclassify or rewrite packets.

[0099] A forwarding function **840** may determine destination of the packet, such as through layer 2 (L2) or layer 3 (L3) mechanisms.

[0100] A set of network receive and transmit queues **842** may handle data structures or buffering to the network interface. Off-FPGA **800** DRAM may be used for packet data.

[0101] A network/remote storage scheduler and policer **844** may provide priorities and control access to the network interface. SLA policies for remote storage and network traffic may be enforced here, which may include strict priorities, weighted round robin, IOP and bandwidth shapers, and policers on a per queue, per initiator, per target, per c-group, or per network flow basis, and the like.

[0102] A local network switch **848** may forward packets between queues in the FPGA, so that traffic does not need to exit the FPGA **800** to the network fabric **122** if the destination is local to the FPGA **800** or the host **102**.

[0103] An end-to-end congestion control/credit facility **850** may prevent network congestion. This is accomplished with two algorithms. First there may be an end-to-end reservation/credit mechanism with a remote FPGA **800**. This may be analogous to a SCSI transfer ready function, where the remote FPGA **800** permits the storage transfer if it can immediately accept the data. Similarly, the local FPGA **800** allocates credits to remote FPGAs **800** as they request a transfer. SLA policies for remote storage may also be enforced here. Second there may be a distributed scheduling algorithm, such as an iterative round-robin algorithm, such as the iSLIP algorithm for input-queues proposed in the publication "The iSLIP Scheduling Algorithm for Input-Queues Switches", by Nick McKeown, IEEE/ACM TRANSACTIONS ON NETWORKING, VOL. 7, NO. 2, APRIL 1999. The algorithm may be performed cluster wide using the intermediate network fabric as the crossbar.

[0104] A rewrite, tag, and CRC facility **852** may encapsulate/de-encapsulate the packet with the appropriate tags and CRC protection.

[0105] A set of interfaces **854**, such as MAC interfaces, may provide an interface to Ethernet.

[0106] A set of embedded CPU and cache complexes **858** may implement a process control plan, exception handling, and other communication to and from the local host and network remote FPGAs **800**.

[0107] A memory controller **860**, such as a DDR controller, may act as a controller for the external DRAM/NVRAM.

[0108] As a result of the integration of functions provided by the converged solution **300**, as embodied in one example by the FPGA **800**, provided herein are methods and systems for combining storage initiation and storage targeting in a single hardware system. In embodiments, these may be attached by a PCIe bus **110**. A single root virtualization function (SR-IOV) or the like may be applied to take any standard device (e.g., any storage media **302** device) and have it act as if it is hundreds of such devices. Embodiments

disclosed herein include using a protocol like SR-IOV to give multiple virtual instances of a physical storage adaptor. SR-IOV is a PCIe standard that virtualizes I/O functions, and while it has been used for network interfaces, the methods and systems disclosed herein extend it to use for storage devices. Thus, provided herein is a virtualized target storage system. In embodiments the virtual target storage system may handle disparate media as if the media are a disk or disks, such as DAS 310.

[0109] Enabled by embodiments like the FPGA 800, embodiments of the methods and systems disclosed herein may also include providing an NVMe device that is virtualized and dynamically allocated. In embodiments one may piggyback the normal NVMe protocol, but carve up, virtualize and dynamically allocate the NVMe device. In embodiments there is no footprint in the software. The operating system 108 stays the same or nearly the same (possibly having a small driver that sees the converged network/storage card 300). This results in virtual storage that looks like a direct attached disk, but the difference is that now we can pool such storage devices 302 across the network 122.

[0110] Methods and systems are disclosed herein for implementing virtualization of NVMe. Regardless how many sources are related to how many destinations, as long as the data from the sources is serialized first before going into the hub, then the hub distributes to data to the designated destination sequentially. If so, then data transport resources such as the DMA queues 804, 832 can be reduced to only one copy. This may include various use scenarios. In one scenario, for NVMe virtual functions (VFs), if they are all connected to the same PCIe bus 110, then regardless how many VFs 807 are configured, the data would be coming into this pool of VFs 807 serially, so there is only one DMA engine 804, and only one storage block (for control information) is needed.

[0111] In another use scenario, for a disk storage system with a pool of discrete disks/controllers, if the data is originated from the physical bus, i.e. PCIe 110, since the data is serially coming into this pool of disks, then regardless how many disks/controllers are in the pool, the transport resources such as the DMA engine 804 can be reduced to only one instead of one per controller.

[0112] Methods and systems disclosed herein may also include virtualization of a converged network/storage adaptor 300. From a traffic perspective, one may combine systems into one. Combining the storage and network adaptors, and adding in virtualization, gives significant advantages. Say there is a single host 102 with two PCIe buses 110. To route from the PCIe 110, you can use a system like remote direct memory access (RDMA) to get to another machine/host 102. If one were to do this separately, one has to configure the storage and the network RDMA systems separately. One has to join each one and configure them at two different places. In the converged solution 300, the whole step of setting up QoS, seeing that this is RDMA and that there is another fabric elsewhere is a zero touch process, because with combined storage and networking the two can be configured in a single step. That is, once one knows the storage, one doesn't need to set up the QoS on the network separately. Thus, single-step configuration of network and storage for RDMA solutions is enabled by the converged solution 300.

[0113] Referring again to FIG. 4, remote access is enabled by the FPGA 800 or similar hardware as described in

connection with FIG. 8. The virtualization boundary is indicated in FIG. 4 by the dotted line 408. To the left of this line, virtual storage devices (e.g., NVMe 802) and virtual network interfaces 830 are presented to the operating system 108. The operating system cannot tell these are virtual devices. To the right of the virtualization boundary 408 are physical storage devices 302 (e.g., using SATA or other protocols noted above) and physical network interfaces. Storage virtualization functions are implemented by the vNVMe 802 and the NVMe virtualization facility 810 of FIG. 8. Network virtualization functions are implemented by the vNIC facility 830. Location of the physical storage media is also hidden from the operating system 108. Effectively, the physical disks 302 across servers can be pooled and accessed remotely. The operating system 108 issues a read or write transaction to the storage media 302 (it is a virtual device, but the operation system 108 sees it as a physical device). If the physical storage media 302 happens to be remote, the read/write transaction is mapped to the proper physical location, encapsulated, and tunneled through Ethernet. This process may be implemented by the device mapping facility 808, the NVMe virtualization facility 810, the data transformation facility 812 and the storage-network tunnel 828 of FIG. 8. The target server (second computing system) un-tunnels the storage read/write and directly accesses its local storage media 302. If the transaction is a write, the data is written to the media 302. If the transaction is a read, the data is prepared, mapped to the origin server, encapsulated, and tunneled through Ethernet. The transaction completion arrives at the origin operating system 102. In a conventional system, these steps would require software intervention in order to process the storage request, data formatting, and network access. As shown, all of these complex software steps are avoided.

[0114] Referring to FIG. 9, a simplified block diagram is provided of an architecture of a controller card 902, as one embodiment of a converged solution 300 as described throughout this disclosure. The controller card 902 may be, for example, a standard, full-height, half-length PCIe card, such as a Gen3  $\times$ 16 card. However, a non-standard card size is acceptable, preferably sized so that it can fit into various types of targeted chassis. The PCIe form factor limits the stack up and layers used on the PCB.

[0115] The controller card 902 may be used as an add-on card on a commodity chassis, such as a 2RU, 4 node chassis. Each node of the chassis (called a sled) is typically 1RU and 6.76" wide. The motherboard typically may provide a PCIe Gen3  $\times$ 16 connector near the back. A riser card may be used to allow the Controller card 902 to be installed on top of the motherboard; thus, the clearance between the card and the motherboard may be limited to roughly on slot width.

[0116] In embodiments, the maximum power supplied by the PCIe connector is 75 W. The controller card 902 may consume about 60 W or less.

[0117] The chassis may provide good airflow, but the card should expect a 10C rise in ambient temperature, because in this example the air will be warmed by dual Xeon processors and 16 DIMMs. The maximum ambient temperature for most servers is 35C, so the air temperature at the controller card 902 will likely be 45C or higher in some situations. Custom heat sinks and baffles may be considered as part of the thermal solution.

[0118] There are two FPGAs in the embodiment of the controller card **902** depicted in FIG. 9, a datapath FPGA, or datapath chip **904**, and a networking FPGA, or networking chip **908**.

[0119] The datapath chip **904** provides connectivity to the host computer **102** over the PCIe connector **110**. From the host processor's point of view, the controller card **902** looks like multiple NVMe devices. The datapath chip **904** bridges NVMe to standard SATA/SAS protocol and in this embodiment controls up to six external disk drives over SATA/SAS links. Note that SATA supports up to 6.0 Gbps, while SAS supports up to 12.0 Gbps.

[0120] The networking chip **908** switches the two 10G Ethernet ports of the NIC device **118** and the eCPU **1018** to two external 10G Ethernet ports. It also contains a large number of data structures for used in virtualization.

[0121] The motherboard of the host **102** typically provides a PCIe Gen3  $\times 16$  interface that can be divided into two separate PCIe Gen3  $\times 8$  busses in the Intel chipset. One of the PCIe Gen3  $\times 8$  bus **110** is connected to the Intel NIC device **118**. The second PCIe Gen3  $\times 8$  bus **110** is connected to a PLX PCIe switch chip **1010**. The downstream ports of the switch chip **1010** are configured as two PCIe Gen3  $\times 8$  busses **110**. One of the busses **110** is connected to the eCPU while the second is connected to the datapath chip **904**.

[0122] The datapath chip **904** uses external memory for data storage. A single  $\times 72$  DDR3 channel **1012** should provide sufficient bandwidth for most situations. The networking chip **908** also uses external memory for data storage, and a single  $\times 72$  DDR3 channel is likely to be sufficient for most situations. In addition, the data structures require the use of non-volatile memory, such as one that provides high performance and sufficient density, such as Non-volatile DIMM (NVDIMM, which typically has a built-in power switching circuit and super-capacitors as energy storage elements for data retention).

[0123] The eCPU **1018** communicates with the networking **908** using two sets of interfaces. It has a PCIe Gen2 $\times 4$  interface for NVMe-like communication. The eCPU **1018** also has two 10G Ethernet interfaces that connect to the networking chip **908**, such as through its L2 switch.

[0124] An AXI bus **1020** (a bus specification of the ARM chipset) will be used throughout the internal design of the two chips **904**, **908**. To allow seamless communication between the datapath chip **904** and the networking chip **908**, the AXI bus **1020** is used for chip-to-chip connection. The Xilinx Aurora™ protocol, a serial interface, may be used as the physical layer.

[0125] The key requirements for FPGA configuration are that (1) The datapath chip **904** must be ready before PCIe configuration started (QSPI Flash memory (serial flash memory with quad SPI bus interface) may be fast enough) and (2) the chips are preferably field upgradeable. The Flash memory for configuration is preferably large enough to store at least 3 copies of the configuration bitstream. The bitstream refers to the configuration memory pattern used by Xilinx™ FPGAs. The bitstream is typically stored in non-volatile memory and is used to configure the FPGA during initial power-on. The eCPU **1018** may be provided with a facility to read and write the configuration Flash memories. New bitstreams may reside with the processor of the host **102**. Security and authentication may be handled by the eCPU **1018** before attempting to upgrade the Flash memories.

[0126] In a networking subsystem, the Controller card **902** may handle all network traffic between the host processor and the outside world. The Networking chip **908** may intercept all network traffics from the NIC **118** and externally.

[0127] The Intel NIC **118** in this embodiment connects two 10GigE, XFI interfaces **1022** to the Networking chip **908**. The embedded processor will do the same. The Networking chip **908** will perform an L2 switching function and route Ethernet traffic out to the two external 10GigE ports. Similarly, incoming 10GigE traffic will be directly to either the NIC **118**, the eCPU **1018**, or internal logic of the Networking chip **908**.

[0128] The controller card **902** may use SFP+ optical connectors for the two external 10G Ethernet ports. In other embodiments, the card may support 10GBASE-T using an external PHY and RJ45 connectors; but a separate card may be needed, or a custom paddle card arrangement may be needed to allow switching between SFP+ and RJ45.

[0129] All the management of the external port and optics, including the operation of the LEDs, may be controlled by the Networking chip **908**. Thus, signals such as PRST, I2C/MDIO, etc may be connected to the Networking chip **908** instead of the NIC **118**.

[0130] In a storage subsystem, the Datapath chip **904** may drive the mini-SAS HD connectors directly. In embodiments such as depicted in FIG. 10, the signals may be designed to operate at 12 Gbps to support the latest SAS standard.

[0131] To provide efficient use of board space, two  $\times 4$  mini-SAS HD connectors may be used. All eight sets of signals may be connected to the Datapath chip **904**, even though only six sets of signals might be used at any one time.

[0132] On the chassis, high-speed copper cables may be used to connect the mini-SAS HD connectors to the motherboard. The placement of the mini-SAS HD connectors may take into account the various chassis' physical space and routing of the cables.

[0133] The power to the controller card **902** may be supplied by the PCIe  $\times 16$  connector. No external power connection needs to be used. Per PCIe specification, the PCIe  $\times 16$  connector may supply only up to 25 W of power after power up. The controller card **902** may be designed such that it draws less than 25 W until after PCIe configuration. Thus, a number of interfaces and components may need to be held in reset after initial power up. The connector may supply up to 75 W of power after configuration, which may be arranged such that the 75 W is split between the 3.3V and 12V rails.

[0134] FIG. 10 shows a software stack **1000**, which includes a driver **1002** to interface to the converged solution **300**, such as one enabled by the FPGA **800**. The NVMe controller **1004** is the set of functions of the hardware (e.g., FPGA **800**) that serves the function of an NVMe controller and allocates virtual devices **1012** to the host. In FIG. 10, dev1, dev2, dev3 are examples of virtual devices **1012** which are dynamically allocated to containers **1018** LXC1, LXC2, and LXC3, respectively. The NVMe to SATA bridge **1008** is the part of the hardware sub-system (e.g., FPGA **800**) that converts and maps virtual devices **1012** (dev1, dev2, dev3) to storage devices **302** (e.g., SSDs in the figure). The connection **1010** is the part of the hardware system that provides a SATA connection (among other possible connection options noted above). The Ethernet link **120**, which can expose virtual devices **1012** (i.e dev1, dev2, dev3) to other



host(s) **102** connected via the Ethernet link **120** using a storage tunneling protocol. The PCI-E (NVMe driver) **1002** may program and drive the hardware subsystem for the storage side. This driver **1002** may run on the host as part of the operating system (e.g., Linux OS in this example). The block layer **1014** may be a conventional SCSI sub-system of the Linux operating system, which may interface with the converged solution PCIe driver **1002** to expose virtual storage devices **1012**. The containers **1018** (LXC1, LXC2, LXC3) may request and dynamically be allocated virtual storage devices **1012** (dev1, dev2 and dev3, respectively).

[0135] FIGS. 11 through 15 show an example of the movement of an application container **1018** (e.g., a Linux container) across multiple systems **102**, first in the absence of a converged solution **300** and then in the presence of such a converged solution **300**. FIG. 11 shows an example of two conventional computer systems **102** with conventional storage controllers **112** and network controllers **118** hosting virtualized software in an OS/Hypervisor stack **108**. Computer System 1 (C1) has a configuration similar to the one shown in FIG. 1 with CPU, memory and conventional storage controller **112** and network controller **118**. The system runs an operating system **108**, such as Linux™, Microsoft Windows™, etc., and/or hypervisor software, such as Xen, VMware, etc. to provide support for multiple applications natively or over virtualized environments, such as via virtual machines or containers. In this computer system **102**, application App1 **1102** is running inside a virtual machine VM1 **1104**. Applications App2 **1108** and App3 **1112** are running within virtualized containers LXC1 **1110** and LXC2 **1114** respectively. In addition to these, application App4 **1118** is running natively over the Operating System **108**. Although typically, a practical scenario might have only virtual machines or containers or native applications (not all three), here it is depicted in a combined fashion deliberately to cover all cases of virtualized environments. Computer System 2 (C2) **102** has similar configuration supporting App5 and App6 in a container and natively, respectively. Each of these applications access their storage devices **302** independent of each other, namely App1 uses **51**, App2 uses **S2**, etc. These storage devices **302** (designated **S1-S6**) are not limited to being independent physical entities. They could be logically carved out of one or more physical storage elements as deemed necessary. As one can see, (represented by the arrow from each storage device **302** to an application), the data flow between the storage **302** and the application **1102**, **1108**, **1112**, **1118** passes through the storage controller **112** and the operating system/hypervisor stack **108** before it reaches the application, entailing the challenges described in connection with FIG. 1.

[0136] Referring to FIG. 12, when an application or a container is moved from C1 to C2, its corresponding storage device needs to be moved too. The movement could be needed due to the fact that C1 might be running out of resources (such as CPU, memory, etc.) to support the existing applications (App1-App4) over a period of time, such as because of behavioral changes within these applications.

[0137] Typically, it is easier to accomplish the movement within a reasonable amount of time as long as the application states and the storage are reasonable in terms of size. Typically storage-intense applications may use large amounts (e.g., multiple terabytes) of storage, in which case,

it may not be practical to move the storage **302** within an acceptable amount of time. In that case, storage may continue to stay where it was and software-level shunting/tunneling would be undertaken to access the storage remotely, as shown in FIG. 13.

[0138] As shown in FIG. 13, App2 **1108**, after its movement to computer system C2, continues to access its original storage **S2** located on computer system C1 by traversing through Operating Systems or Hypervisors **108** of both the systems C1 and C2. This is because the mapping of storage access through the network controllers **118** to that storage controller **112** and its attached storage devices **302** is done by the Operating System or Hypervisor software stack **108** running within the main CPU.

[0139] As shown in FIG. 13 after its movement to C2, App2 **1108** continues to access its original storage **S2** located in C1 by traversing through Operating Systems or Hypervisors **108** of both the systems C1 and C2. This is because, the mapping of storage access through the network controllers **118** from C2 to C1 and over to that storage controller **112** of C1 is done by the Operating System or Hypervisor software **108** running within the main CPU of each computer system.

[0140] Consider a similar scenario when a converged controller **300** is applied as shown in the FIG. 14. As one can see, the scenario is almost identical to FIG. 11, except the Converged IO Controller **300** replaces the separate storage controller **112** and network controller **118**. In this case, when App2 **1108** along with its container LXC1 is moved to C2 (as shown in FIG. 15), the storage **S2** is not moved, and the access is optimized by avoiding the traversal through any software (Operating System, Hypervisor **108** or any other) running in main CPU present in computing system C1.

[0141] Thus, provided herein is a novel way of bypassing the main CPU where a storage device is located, which in turn (a) allows one to reduce latency and bandwidth significantly in accessing a storage across multiple computer systems and (b) vastly simplifies and improves situations in which an application needs to be moved away from a machine on which its storage is located.

[0142] Ethernet networks behave on a best effort basis and hence lossy in nature as well as bursty. Any packet could be lost forever or buffered and delivered in bursty and delayed manner along with other packets. Whereas, typical storage centric applications are sensitive to losses and bursts, it is important that when storage traffic is sent over Ethernet networks.

[0143] Conventional storage accesses over their buses/networks involve reliable and predictable methods. For example, Fibre Channel networks employ credit based flow control to limit number of accesses made by end systems. And the number of credits given to an end system is based on whether the storage device has enough command buffers to receive and fulfill storage requests in predictable amount of time fulfilling required latency and bandwidth needs. The figure below shows some credit schemes adopted by different types of buses such as SATA, Fibre Channel (FC), SCSI, SAS, etc.

[0144] Referring to FIG. 16, Ethernet networks behave on a best effort basis and hence tend to be lossy in nature, as well as bursty. Any packet could be lost forever or buffered and delivered in a delayed manner, in a congestion-inducing burst, along with many other packets. Typical storage-centric applications are sensitive to losses and bursts, so it is

important when storage traffic is sent over buses and Ethernet networks, that those involve reliable and predictable methods for maintaining integrity. For example, Fibre Channel networks conventionally employ credit-based flow control to limit the number of accesses made by end systems at any one time. The number of credits given to an end system can be based on whether the storage device **302** has enough command buffers to receive and fulfill storage requests in a predictable amount of time that satisfies required latency and bandwidth requirements. FIG. 16 shows some of the credit schemes adopted by different types of buses such as a SATA bus **1602**, Fibre Channel (FC) **1604**, and SCSI/SAS connection **1608**, among other types of such schemes.

**[0145]** As one can see, for example, an FC controller **1610** may have its own buffering up to a limit of 'N' storage commands before sending them to an FC-based storage device **1612**, but the FC device **1612** might have a different buffer limit, say 'M' in this example, which could be greater than, equal to, or less than 'N'. A typical credit-based scheme uses target level (in this example, one of the storage devices **302**, such as the FC Device **1602**, is the target) aggregate credits, information about which is propagated to various sources (in this example, the controller, such as the FC Controller **1610**, is the source) which are trying to access the target **302**. For example, if two sources are accessing a target that has a queue depth of 'N,' then sum of the credits given to the sources would not exceed 'N,' so that at any given time the target will not receive more than 'N' commands. The distribution of credits among the sources may be arbitrary, or it may be based on various types of policies (e.g., priorities based on cost/pricing, SLAs, or the like). When the queue is serviced, by fulfilling the command requests, credits may be replenished at the sources as appropriate. By adhering to this kind of credit-based storage access, losses that would result from queues at the target being overwhelmed can be avoided.

**[0146]** Typical storage accesses over Ethernet, such as FCOE, iSCSI, and the like, may extend the target-oriented, credit-based command fulfillment to transfers over Ethernet links. In such cases, they may be target device-oriented, rather than being source-oriented. Provided herein are new credit based schemes that can instead be based on which or what kind of source should get how many credits. For example, the converged solution **300** described above, which directly interfaces the network to the storage, may employ a multiplexer to map a source-oriented, credit-based scheduling scheme to a target device oriented credit based scheme, as shown in FIG. 17.

**[0147]** As shown in FIG. 17, four sources are located over Ethernet and there are two target storage devices **302**. Typical target-oriented, credit-based schemes would expose two queues (one per target), or two connections per source to each of the targets. Instead, as shown in FIG. 17, the queues (Q1,Q2,Q3,Q4) **1702** are on a per-source basis, and they mapped/multiplexed to two target-oriented queues (Q5, Q6) **1704** across the multiplexor (S) **1708**. By employing this type of source-oriented, credit-based scheme, one may guarantee access bandwidth and predictable access latency, independent of the number of target storage devices **302**. As an example, one type of multiplexing is to make sure queue size 'P' of Q1 does not exceed 'L+M' of Q5 and Q6, so that Q1 is not overwhelmed by its source.

**[0148]** In embodiments, methods and systems to provide access to blocks of data from a storage device **302** is

described. In particular, a novel approach to allowing an application to access its data, fulfilling a specific set of access requirements is described.

**[0149]** **[000149]** As used herein, the term "application-driven data storage" (ADS) encompasses storage that provides transparency to any application in terms of how the application's data is stored, accessed, transferred, cached and delivered to the application. ADS may allow applications to control these individual phases to address the specific needs of the particular application. As an example, an application might be comprised of multiple instances of itself, as well as multiple processes spread across multiple Linux nodes across the network. These processes might access multiple files in shared or exclusive manners among them. Based on how the application wants to handle these files, these processes may want to access different portions of the files more frequently, may need quick accesses or use once and throw away. Based on these criteria, it might want to prefetch and/or retain specific portions of a file in different tiers of cache and/or storage for immediate access on per session or per file basis as it wishes. These application specific requirements cannot be fulfilled in a generic manner such as disk striping of entire file system, prefetching of read-ahead sequential blocks, reserving physical memory in the server or LRU or FIFO based caching of file contents.

**[0150]** Application-driven data storage I/O is not simply applicable to the storage entities alone. It impacts the entire storage stack in several ways. First, it impacts the storage I/O stack in the computing node where the application is running comprising the Linux paging system, buffering, underlying File system client, TCP/IP stack, classification, QoS treatment and packet queuing provided by the networking hardware and software. Second, it impacts the networking infrastructure that interconnects the application node and its storage, comprising Ethernet segments, optimal path selections, buffering in switches, classification and QoS treatment of latency-sensitive storage traffic as well as implosion issues related to storage I/O. Also, it impacts the storage infrastructure which stores and maintains the data in terms of files comprising the underlying file layout, redundancy, access time, tiering between various types of storage as well as remote repositories.

**[0151]** Methods and systems disclosed herein include ones relating to the elements affecting a typical application within an application node and how a converged solution **300** may change the status quo to address certain critical requirements of applications.

**[0152]** Conventional Linux stacks may consist of simple building blocks, such generic memory allocation, process scheduling, file access, memory mapping, page caching, etc. Although these are essential for an application to run on Linux, this is not optimal for certain categories of applications that are input/output (IO) intensive, such as NoSQL. NoSQL applications are very IO intensive, and it is harder to predict their data access in a generic manner. If applications have to be deployed in a utility-computing environment, it is not ideal for Linux to provide generic minimal implementations of these building blocks. It is preferred for these building blocks to be highly flexible and have application-relevant features that can be controllable from the application(s).

**[0153]** Although every application has its own specific requirements, in an exemplary embodiment, the NoSQL class of applications has the following requirements which,

when addressed by the Linux stack, could greatly improve the performance of NoSQL applications and other IO intensive applications. The requirements are first, the use of file level priority. The Linux file system should provide access level priority between different files at a minimum. For example, an application process (consisting of multiple threads) accessing two different files with one file given higher priority over the other (such as one database/table/index over the other). This would enable the precious storage I/O resources be preferentially utilized based on the data being accessed. One would argue that this could be indirectly addressed by running one thread/process be run at a higher or lower priority, but those process level priorities are not communicated over to file system or storage components. Process or thread level priorities are meant only for utilizing CPU resources. Moreover, it is possible that same thread might be accessing these two files and hence will be utilizing the storage resources at two different levels based on what data (file) being accessed. Second, there may be a requirement for access level preferences. A Linux file system may provide various preferences (primarily SLA) during a session of a file (opened file), such as priority between file sessions, the amount of buffering of blocks, the retention/life time preferences for various blocks, alerts for resource thresholds and contentions, and performance statistics. As an example, when a NoSQL application such as MongoDB or Cassandra would have two or more threads for writes and reads, where if writes may have to be given preference over reads, a file session for write may have to be given preference over a file session for read for the same file. This capability enables two sessions of the same file to have two different priorities.

**[0154]** Many of the NoSQL applications store different types of data into the same file; for example, MongoDB stores user collections as well as (b-tree) index collections in the same set of database files. MongoDB may want to keep the index pages (btree and collections) in memory in preference over user collection pages. When these files are opened, MongoDB may want to influence the Linux, File and storage systems to treat the pages according to MongoDB policies as opposed to treating these pages in a generic FIFO or LRU basis agnostic of the application's requirements.

**[0155]** Resource alerts and performance statistics enable an NoSQL database to understand the behavior of the underlying File and storage system and could service its database queries accordingly or trigger actions to be carried out such as sharding of the database or reducing/increasing of File I/O preference for other jobs running in the same host (such as backup, sharding, number read/write queries serviced, etc.). For example, performance stats about min, max and average number of IOPs and latencies as well as top ten candidate pages thrashed in and out of host memory in a given period of time would enable an application to fine tune itself dynamically adjusting the parameters noted above.

**[0156]** A requirement may also exist for caching and tiering preferences. A Linux file system may need to have a dynamically configurable caching policy while applications are accessing their files. Currently, Linux file systems typically pre-fetch contiguous blocks of a file, hoping that applications are reading the file in a sequential manner like a stream. Although it is true for many legacy applications like web servers and video streamers, emerging NoSQL applications do not follow sequential reads strictly. These

applications read blocks randomly. As an example, MongoDB stores the document keys in index tables in b-tree, laid out flat on a portion of a file, which, when a key is searched, accesses the blocks randomly until it locates the key. Moreover, these files are not dedicated to such b-tree based index tables alone. These files are shared among various types of tables (collections) such as user documents and system index files. Because of this, a Linux file system cannot predict what portions of the file need to be cached, read ahead, swapped out for efficient memory usage, etc.

**[0157]** In embodiments of the methods and systems described herein, there is a common thread across various applications in their requirements for storage. In particular, latency and IOPs for specific types of data at specific times and places of need are very impactful on performance of these applications.

**[0158]** For example, to address the host level requirements listed above, disclosed herein are methods and systems for a well fine-tuned file-system client that enables applications to completely influence and control the storing, retrieving, retaining and tiering of data according to preference within the host and elsewhere.

**[0159]** As shown in FIG. 18, a File System (FS) client **1802** keeps separate buffer pools for separate sessions of a file (fd1 and fd2). It also pre-allocates and maintains aggregate memory pools for each application or set of processes. The SLA-Broker **1804** may be exercised by the application either internally within the process/thread where the file I/O is carried out or externally from another set of processes, to influence the FS Client **1802** to provide appropriate storage I/O SLAs dynamically. Controlling the SLA from an external process enables a legacy application with no knowledge of these newer storage control features immediately without modifying the application itself.

**[0160]** Methods and systems disclosed herein may provide extensive tiering services for data retrieval across network and hosts. As one can see in FIG. 19 below, a High Performance Distributed File Server (DFS) **1902** enables application to run in the Platform **1904** in a containerized form to determine and execute what portions of files should reside in which media (DRAM, NVRAM, SSD or HDDs) in cached form storage form dynamically. These application containers **1908** can determine other storage policies such as whether a file has to be striped, mirrored, RAIDed and disaster recovered (DR'ed) as well.

**[0161]** The methods and systems disclosed herein also provide extensive caching service, wherein an application container in the High Performance DFS **1902** can proactively retrieve specific pages of a file from local storage and/or remote locations and push these pages to specific places for fast retrieval later when needed. For instance, the methods and systems may local memory and SSD usages of the hosts running the application and proactively push pages of an application's interest into any of these hosts' local memory/SSD. The methods and systems may use the local tiers of memory, SSD and HDD provisioned for this purpose in the DFS platform **1904** for very low latency retrieval by the application at a later time of its need.

**[0162]** The use of extending the cache across hosts of the applications is immense. For example, in MongoDB when the working set temporarily grows beyond its local host's memory, thrashing happens, and it significantly reduces the query handling performance. This is because when a needed file data page is discarded in order to bring in a new page to

satisfy a query and subsequently, if the original page has to be brought back, the system has to reread the page afresh from the disk subsystem, thereby incurring huge latency in completing a query. Application-driven storage access helps these kinds of scenarios by keeping a cache of the discarded page elsewhere in the network (in another application host's memory/SSD or in local tiers of storage of the High Performance DFS system 1902) temporarily until MongoDB requires the page again and thereby significantly reducing the latency in completing the query.

[0163] Referring to FIG. 20, High Performance DFS 1902 takes advantage of DRAM and SSD resources across the application hosts in a single, unified RAM and SSD-based tier/cache 2002, in order to cache and serve the application data as necessary and as influenced and controlled by the application.

[0164] A system comprising of a set of hosts (H1 through HN), a file or block server 2102 and a storage subsystem 2104 is disclosed herein as shown in the FIG. 21. A host H1-HN is typically a computer running an application that needs access to data permanently or temporarily stored in storage. The file or volume server 2102 may be a data organizer and a data server, typically running a hardware comprising a central processing unit (CPU), memory and special hardware to connect to external devices such as networking and storage devices. The file or volume server 2102 organizes user data in terms of multiple fixed or variable number of bytes called blocks. It stores these blocks of data in an internal or external storage. A random, but logically related, sequence of blocks is organized into a file or a volume. One or more Hosts H1-HN can access these files or volumes through an application programming interface (API) or any other protocol. A file or volume server can serve one or more files and volumes to one or more hosts. It is to be noted that a host and a file or volume server can be in two different physical entities connected directly or through a network or they could be logically located together in a single physical computer.

[0165] Storage 2104 may be a collection of entities capable of retaining a piece of data temporarily or permanently. This is typically comprised of static or dynamic random access memory (RAM), solid state storage (SSD), hard disk drive (HDD) or a combination of all of these. Storage could be an independent physical entity connected to a File or volume server 2102 through a link or a network. It could also be integrated with file or volume server 2102 in a single physical entity. Hence, hosts H1-HN, file or volume server 2102 and storage 2104 could be physically collocated in a single hardware entity.

[0166] A host is typically comprised of multiple logical entities as shown in FIG. 22. An application 2202 typically runs in a host and would access its data elements through an API provided by its local operating system 2204 or any other entity in place of it. The operating system 2204 typically has a standard API interface to interface to a file system through their file system client 2206. A file system client 2206 is a software entity running within the host to interface with a file or volume server 2210 either located remotely or locally. It would provide the data elements needed by application 2202, which are present in a single or multiple files or volumes, by retrieving them from file or volume server 2210 and keeping them in the host's memory 2208 until the application completes its processing of the elements of data. In a typical application scenario, a specific piece of data

would be read and/or modified multiple number of times as required. It is also typical that an entire file or volume, consisting of multiple data elements, is potentially much larger than the size of local memory 2208 in certain types of applications. This makes operating system 2204 and file system client 2206 more complicated in its implementation in order to decide what blocks of data to be retained in or evicted from memory 2208 based on the prediction that the application 2202 may or may not access them in future. So far, the existing implementations execute some generic and application-independent methods, such as first-in-first-out (FIFO) or least-recently-used (LRU), to retain or evict the blocks of data in memory in order to bring in new blocks of data from file or volume server 2210. Moreover, when a memory occupied by a block of data is to be reclaimed for storing another block of data, the original data is simply erased without the consideration for its future use. Normally, the disk subsystem is very slow and incurs high latency when a block of data is read from it and transferred by file or volume server 2210 to file system client 2206 to memory 2208. So, when the original block of data is erased, the application might have to wait longer if it tries to access the original data in near future. The main problem with this kind of implementation is that none of the modules in the path of data access, namely operating system 2204, file system client 2206, memory 2208, block server 2210 and storage have any knowledge of what, when and how often a block of data is going to be accessed by application 2202.

[0167] An example scenario depicting an application 2202 accessing a block of data from storage 2212 is shown in FIG. 23. The numbered circles are to show the steps involved in the process of accessing a block of data. These steps are explained below. First, application 2202 uses API of file or Operating System 2204 to access a block of data operating system 2204 invokes an equivalent API for file system client 2206 to access the same. Second, file system client 2206 tries to find if the data exists in its local memory buffers dedicated for this purpose. If found, steps (3) through (7) below are skipped. Third, sends a command to retrieve the data from block server 2210. Fourth, block server 2210 sends a read command to storage 2212 to read the block of data from the storage. Fifth, storage 2212 returns the block of data to block server 2210 after reading it from the storage. Sixth, block server 2210 returns the block of data to file system client 2206. Seventh, file system client 2206 saves the data in a memory buffer in memory 2208 for any future access. Eighth, file system client 2206 returns the requested data to the application 2202.

[0168] In the methods and systems disclosed herein, in order to address performance requirements related to data access by most newer class of applications in the area of NoSQL and BigData, it is proposed that the components in the data block access comprising operating system 2204, file system client 2206, memory 2208, block server 2210 and storage 2212 be controlled by any application 2202. Namely, we propose the following. First, enable operating system 2204 to provide additional API to allow applications to control file system client 2206. Second, enhance file system client 2206 to support the following: (a) allow application 2202 to create a dedicated pool of memory in memory 2208 for a particular file or volume, in the sense, a file or volume will have a dedicated pool of memory buffers to hold data specific to it which are not shared or removed for the purposes of other files or volumes; (b) allow application

**2202** to create a dedicated pool of memory in memory **2208** for a particular session with a file or volume such that two independent sessions with a file or volume will have independent memory buffers to hold their data. As an example, a critically important file session may have large number of memory buffers in memory **2208**, so that the session can take advantage of more data being present for quicker and frequent access, whereas a second session with the same file may be assigned with very few buffers and hence it might have to incur more delay and reuse of its buffers to access various parts of the file; (c) allow application **2202** to create an extended pool of buffers beyond memory **2208** across other hosts or block server **2210** for quicker access. This enables blocks of data be kept in memory **2208** of other hosts as well as any memory **2402** present in the file or block server **2210**; (d) allow application **2202** to make any block of data to be more persistent in memory **2208** relative to other blocks of data for a file, volume or a session. This allows an application to pick and choose a block of data to be always available for immediate access and not let operating system **2204** or file system client **2206** to evict it based on their own eviction policies; and (e) allow application **2202** to make any block of data to be less persistent in memory **2208** relative to other blocks of data for a file, volume or a session. This allows an application to let know operating system **2204** and file system client **2206** to evict and reuse the buffer of the data block as and when they choose to. This helps in retaining other normal blocks of data for longer period of time. Third, enable block server **2210** to host application specific modules in terms of application container **2400** as shown in the FIG. 24 with the following capabilities: (a) enable application container **2400** to fetch blocks of data of interest to application **2202** ahead of time and store them in local memory **2402** for later quick access and avoid the latency penalty associated with storage **2212** and (b) enable storing of evicted pages from memory **2208** of hosts in local memory **2402** for any later access by application **2202**.

**[0169]** The application driven feature of (2)(c) above needs further explanation. There are two scenarios. The first one involves block of data being retrieved from the memory of block server **2210**. The other scenario involves retrieving the same from another host. Assuming the exact same block data has been read from storage **2212** by two hosts (H1) and (H2), the methods and systems disclosed herein provide a system such as depicted in FIG. 25. When a block of data is noticed to be present in another host (H2), it is directly retrieved from it by file system client **2206** instead asking block server **2210** to retrieve it from storage **2212**, which will be slower and incurs high latency.

**[0170]** In embodiments, if file system client **2206** decides to evict a block of data from (D1) because of storing a more important block of data in its place, file system client **2206** could send the evicted block of data to file system client **2206'** to be stored in memory **2208'** on its behalf.

**[0171]** It should be noted that the abovementioned techniques can be applied to achieving fast failover in case of failure(s) of Hosts. Furthermore the caching techniques described above; especially pertaining to RAM can be used to achieve failover with a warm cache. FIG. 25 shows an example of a fast failover system with a warm cache. The end result is that during a failure of a node, the end application on a new node will not undergo a time period

before the cache (in RAM) is warmed and thereby incur a period of lower application performance.

**[0172]** Provided herein is a system and method with a processor and a file server with an application specific module to control the storage access according to the application's needs.

**[0173]** Also provided herein is a system and method with a processor and a data (constituting blocks of fixed size bytes, similar or different objects with variable number of bytes) storage enabling an application specific module to control the storage access according to the application's needs.

**[0174]** Also provided herein is a system and method which retrieves a stale file or storage data block, previously maintained for the purposes of an application's use, from a host's memory and/or its temporary or permanent storage element and stores it in another host's memory or and/or its temporary or permanent storage element, for the purposes of use by the application at a later time.

**[0175]** Also provided herein is a system and method which retrieves any file or storage data block, previously maintained for the purposes of an application's use, from a host's memory and/or its temporary or permanent storage element and stores it in another host's memory or and/or its temporary or permanent storage element, for the purposes of use by the application at a later time.

**[0176]** Also provided herein is a system and method which utilizes memory and/or its temporary or permanent storage element of a host to store any file or storage data block which would be subsequently accessed by an application running in another host for the purposes of reducing latency of data access.

**[0177]** File or storage data blocks, previously maintained for the purposes of an application's use, from a host's memory and/or its temporary or permanent storage element, may be stored in another host's memory or and/or its temporary or permanent storage element, for the purposes of use by the application at a later time.

**[0178]** The mechanism of transferring a file or storage data block, previously maintained for the purposes of an application's use, from a host's memory and/or its temporary or permanent storage element to another host over a network.

**[0179]** In accordance with various exemplary and non-limiting embodiments, there is disclosed a device comprising a converged input/output controller that includes a physical target storage media controller, a physical network interface controller and a gateway between the storage media controller and the network interface controller, wherein gateway provides a direct connection for storage traffic and network traffic between the storage media controller and the network interface controller.

**[0180]** In accordance with some embodiments, the device may further comprise a virtual storage interface that presents storage media controlled by the storage media controller as locally attached storage, regardless of the location of the storage media. In accordance with yet other embodiments, the device may further comprise a virtual storage interface that presents storage media controlled by the storage media controller as locally attached storage, regardless of the type of the storage media. In accordance with yet other embodiments, the device may further comprise a virtual storage interface that facilitates dynamic provisioning of the storage media, wherein the physical storage may be either local or remote.

[0181] In accordance with yet other embodiments, the device may further comprise a virtual network interface that facilitates dynamic provisioning of the storage media, wherein the physical storage may be either local or remote. In accordance with yet other embodiments, the device may be adapted to be installed as a controller card on a host computing system, in particular, wherein the gateway operates without intervention by the operating system of the host computing system.

[0182] In accordance with yet other embodiments, the device may include at least one field programmable gate array providing at least one of the storage functions and the network functions of the device. In accordance with yet other embodiments, the device may be configured as a network-deployed switch. In accordance with yet other embodiments, the device may further comprise a functional component of the device for translating storage media instructions between a first protocol and at least one other protocol.

[0183] With reference to FIG. 26, there is illustrated an exemplary and non-limiting method of virtualization of a storage device. First, at step 2600 there is accessed a physical storage device that responds to instructions in a first storage protocol. Next, at step 2602, there are translated instructions between the first storage protocol and a second storage protocol. Lastly, at step 2604, using the second protocol, the physical storage device is presented to an operating system, such that the storage of the physical storage device can be dynamically provisioned, whether the physical storage device is local or remote to a host computing system that uses the operating system.

[0184] In accordance with various embodiments, the first protocol is at least one of a SATA protocol, an NVMe protocol, a SAS protocol, an iSCSI protocol, a fiber channel protocol and a fiber channel over Ethernet protocol. In other embodiments, the second protocol is an NVMe protocol.

[0185] In some embodiments, the method may further comprise providing an interface between an operating system and a device that performs the translation of instructions between the first and second storage protocols and/or providing an NVMe over Ethernet connection between the device that performs the translation of instructions and a remote, network-deployed storage device.

[0186] With reference to FIG. 27, there is illustrated an exemplary and non-limiting method of facilitating migration of at least one of an application and a container. First, at step 2700, there is provided a converged storage and networking controller, wherein a gateway provides a connection for network and storage traffic between a storage component and a networking component of the device without intervention of the operating system of a host computer. Next, at step 2702, the at least one application or container is mapped to a target physical storage device that is controlled by the converged storage and networking controller, such that the application or container can access the target physical storage, without intervention of the operating system of the host system to which the target physical storage is attached, when the application or container is moved to another computing system.

[0187] In accordance with various embodiments, the migration is of a Linux container or a scaleout application.

[0188] In accordance with yet other embodiments, the target physical storage is a network-deployed storage device that uses at least one of an iSCSI protocol, a fiber channel

protocol and a fiber channel over Ethernet protocol. In yet other embodiments, the target physical storage is a disk attached storage device that uses at least one of a SAS protocol, a SATA protocol and an NVMe protocol.

[0189] With reference to FIG. 28, there is illustrated an exemplary and non-limiting method of providing quality of service (QoS) for a network. First, at step 2800, there is provided a converged storage and networking controller, wherein a gateway provides a connection for network and storage traffic between a storage component and a networking component of the device without intervention of the operating system of a host computer. Next, at step 2802, without intervention of the operating system of a host computer, there is managed at least one quality of service (QoS) parameter related to a network in the data path of which the storage and networking controller is deployed, such managing being based on at least one of the storage traffic and the network traffic that is handled by the converged storage and networking controller.

[0190] While only a few embodiments of the present disclosure have been shown and described, it will be obvious to those skilled in the art that many changes and modifications may be made thereunto without departing from the spirit and scope of the present disclosure as described in the following claims. All patent applications and patents, both foreign and domestic, and all other publications referenced herein are incorporated herein in their entireties to the full extent permitted by law.

[0191] The methods and systems described herein may be deployed in part or in whole through a machine that executes computer software, program codes, and/or instructions on a processor. The present disclosure may be implemented as a method on the machine, as a system or apparatus as part of or in relation to the machine, or as a computer program product embodied in a computer readable medium executing on one or more of the machines. In embodiments, the processor may be part of a server, cloud server, client, network infrastructure, mobile computing platform, stationary computing platform, or other computing platform. A processor may be any kind of computational or processing device capable of executing program instructions, codes, binary instructions and the like. The processor may be or may include a signal processor, digital processor, embedded processor, microprocessor or any variant such as a co-processor (math co-processor, graphic co-processor, communication co-processor and the like) and the like that may directly or indirectly facilitate execution of program code or program instructions stored thereon. In addition, the processor may enable execution of multiple programs, threads, and codes. The threads may be executed simultaneously to enhance the performance of the processor and to facilitate simultaneous operations of the application. By way of implementation, methods, program codes, program instructions and the like described herein may be implemented in one or more thread. The thread may spawn other threads that may have assigned priorities associated with them; the processor may execute these threads based on priority or any other order based on instructions provided in the program code. The processor, or any machine utilizing one, may include non-transitory memory that stores methods, codes, instructions and programs as described herein and elsewhere. The processor may access a non-transitory storage medium through an interface that may store methods, codes, and instructions as described herein and elsewhere. The

storage medium associated with the processor for storing methods, programs, codes, program instructions or other type of instructions capable of being executed by the computing or processing device may include but may not be limited to one or more of a CD-ROM, DVD, memory, hard disk, flash drive, RAM, ROM, cache and the like.

**[0192]** A processor may include one or more cores that may enhance speed and performance of a multiprocessor. In embodiments, the process may be a dual core processor, quad core processors, other chip-level multiprocessor and the like that combine two or more independent cores (called a die).

**[0193]** The methods and systems described herein may be deployed in part or in whole through a machine that executes computer software on a server, client, firewall, gateway, hub, router, or other such computer and/or networking hardware. The software program may be associated with a server that may include a file server, print server, domain server, internet server, intranet server, cloud server, and other variants such as secondary server, host server, distributed server and the like. The server may include one or more of memories, processors, computer readable media, storage media, ports (physical and virtual), communication devices, and interfaces capable of accessing other servers, clients, machines, and devices through a wired or a wireless medium, and the like. The methods, programs, or codes as described herein and elsewhere may be executed by the server. In addition, other devices required for execution of methods as described in this application may be considered as a part of the infrastructure associated with the server.

**[0194]** **[000194]** The server may provide an interface to other devices including, without limitation, clients, other servers, printers, database servers, print servers, file servers, communication servers, distributed servers, social networks, and the like. Additionally, this coupling and/or connection may facilitate remote execution of program across the network. The networking of some or all of these devices may facilitate parallel processing of a program or method at one or more location without deviating from the scope of the disclosure. In addition, any of the devices attached to the server through an interface may include at least one storage medium capable of storing methods, programs, code and/or instructions. A central repository may provide program instructions to be executed on different devices. In this implementation, the remote repository may act as a storage medium for program code, instructions, and programs.

**[0195]** **[000195]** The software program may be associated with a client that may include a file client, print client, domain client, internet client, intranet client and other variants such as secondary client, host client, distributed client and the like. The client may include one or more of memories, processors, computer readable media, storage media, ports (physical and virtual), communication devices, and interfaces capable of accessing other clients, servers, machines, and devices through a wired or a wireless medium, and the like. The methods, programs, or codes as described herein and elsewhere may be executed by the client. In addition, other devices required for execution of methods as described in this application may be considered as a part of the infrastructure associated with the client.

**[0196]** The client may provide an interface to other devices including, without limitation, servers, other clients, printers, database servers, print servers, file servers, communication servers, distributed servers and the like. Addi-

tionally, this coupling and/or connection may facilitate remote execution of program across the network. The networking of some or all of these devices may facilitate parallel processing of a program or method at one or more location without deviating from the scope of the disclosure. In addition, any of the devices attached to the client through an interface may include at least one storage medium capable of storing methods, programs, applications, code and/or instructions. A central repository may provide program instructions to be executed on different devices. In this implementation, the remote repository may act as a storage medium for program code, instructions, and programs.

**[0197]** The methods and systems described herein may be deployed in part or in whole through network infrastructures. The network infrastructure may include elements such as computing devices, servers, routers, hubs, firewalls, clients, personal computers, communication devices, routing devices and other active and passive devices, modules and/or components as known in the art. The computing and/or non-computing device(s) associated with the network infrastructure may include, apart from other components, a storage medium such as flash memory, buffer, stack, RAM, ROM and the like. The processes, methods, program codes, instructions described herein and elsewhere may be executed by one or more of the network infrastructural elements. The methods and systems described herein may be adapted for use with any kind of private, community, or hybrid cloud computing network or cloud computing environment, including those which involve features of software as a service (SaaS), platform as a service (PaaS), and/or infrastructure as a service (IaaS).

**[0198]** The methods, program codes, and instructions described herein and elsewhere may be implemented on a cellular network has sender-controlled contact media content item multiple cells. The cellular network may either be frequency division multiple access (FDMA) network or code division multiple access (CDMA) network. The cellular network may include mobile devices, cell sites, base stations, repeaters, antennas, towers, and the like. The cell network may be a GSM, GPRS, 3G, EVDO, mesh, or other networks types.

**[0199]** The methods, program codes, and instructions described herein and elsewhere may be implemented on or through mobile devices. The mobile devices may include navigation devices, cell phones, mobile phones, mobile personal digital assistants, laptops, palmtops, netbooks, pagers, electronic books readers, music players and the like. These devices may include, apart from other components, a storage medium such as a flash memory, buffer, RAM, ROM and one or more computing devices. The computing devices associated with mobile devices may be enabled to execute program codes, methods, and instructions stored thereon. Alternatively, the mobile devices may be configured to execute instructions in collaboration with other devices. The mobile devices may communicate with base stations interfaced with servers and configured to execute program codes. The mobile devices may communicate on a peer-to-peer network, mesh network, or other communications network. The program code may be stored on the storage medium associated with the server and executed by a computing device embedded within the server. The base station may include a computing device and a storage medium. The

storage device may store program codes and instructions executed by the computing devices associated with the base station.

**[0200]** The computer software, program codes, and/or instructions may be stored and/or accessed on machine readable media that may include: computer components, devices, and recording media that retain digital data used for computing for some interval of time; semiconductor storage known as random access memory (RAM); mass storage typically for more permanent storage, such as optical discs, forms of magnetic storage like hard disks, tapes, drums, cards and other types; processor registers, cache memory, volatile memory, non-volatile memory; optical storage such as CD, DVD; removable media such as flash memory (e.g. USB sticks or keys), floppy disks, magnetic tape, paper tape, punch cards, standalone RAM disks, Zip drives, removable mass storage, off-line, and the like; other computer memory such as dynamic memory, static memory, read/write storage, mutable storage, read only, random access, sequential access, location addressable, file addressable, content addressable, network attached storage, storage area network, bar codes, magnetic ink, and the like.

**[0201]** The methods and systems described herein may transform physical and/or intangible items from one state to another. The methods and systems described herein may also transform data representing physical and/or intangible items from one state to another.

**[0202]** The elements described and depicted herein, including in flow charts and block diagrams throughout the figures, imply logical boundaries between the elements. However, according to software or hardware engineering practices, the depicted elements and the functions thereof may be implemented on machines through computer executable media has sender-controlled contact media content item a processor capable of executing program instructions stored thereon as a monolithic software structure, as standalone software modules, or as modules that employ external routines, code, services, and so forth, or any combination of these, and all such implementations may be within the scope of the present disclosure. Examples of such machines may include, but may not be limited to, personal digital assistants, laptops, personal computers, mobile phones, other handheld computing devices, medical equipment, wired or wireless communication devices, transducers, chips, calculators, satellites, tablet PCs, electronic books, gadgets, electronic devices, devices has sender-controlled contact media content item artificial intelligence, computing devices, networking equipment, servers, routers and the like. Furthermore, the elements depicted in the flow chart and block diagrams or any other logical component may be implemented on a machine capable of executing program instructions. Thus, while the foregoing drawings and descriptions set forth functional aspects of the disclosed systems, no particular arrangement of software for implementing these functional aspects should be inferred from these descriptions unless explicitly stated or otherwise clear from the context. Similarly, it will be appreciated that the various steps identified and described above may be varied, and that the order of steps may be adapted to particular applications of the techniques disclosed herein. All such variations and modifications are intended to fall within the scope of this disclosure. As such, the depiction and/or description of an order for various steps should not be understood to require a particular order of execution for those steps, unless

required by a particular application, or explicitly stated or otherwise clear from the context.

**[0203]** The methods and/or processes described above, and steps associated therewith, may be realized in hardware, software or any combination of hardware and software suitable for a particular application. The hardware may include a general-purpose computer and/or dedicated computing device or specific computing device or particular aspect or component of a specific computing device. The processes may be realized in one or more microprocessors, microcontrollers, embedded microcontrollers, programmable digital signal processors or other programmable device, along with internal and/or external memory. The processes may also, or instead, be embodied in an application specific integrated circuit, a programmable gate array, programmable array logic, or any other device or combination of devices that may be configured to process electronic signals. It will further be appreciated that one or more of the processes may be realized as a computer executable code capable of being executed on a machine-readable medium.

**[0204]** The computer executable code may be created using a structured programming language such as C, an object oriented programming language such as C++, or any other high-level or low-level programming language (including assembly languages, hardware description languages, and database programming languages and technologies) that may be stored, compiled or interpreted to run on one of the above devices, as well as heterogeneous combinations of processors, processor architectures, or combinations of different hardware and software, or any other machine capable of executing program instructions.

**[0205]** Thus, in one aspect, methods described above and combinations thereof may be embodied in computer executable code that, when executing on one or more computing devices, performs the steps thereof. In another aspect, the methods may be embodied in systems that perform the steps thereof, and may be distributed across devices in a number of ways, or all of the functionality may be integrated into a dedicated, standalone device or other hardware. In another aspect, the means for performing the steps associated with the processes described above may include any of the hardware and/or software described above. All such permutations and combinations are intended to fall within the scope of the present disclosure.

**[0206]** While the disclosure has been disclosed in connection with the preferred embodiments shown and described in detail, various modifications and improvements thereon will become readily apparent to those skilled in the art. Accordingly, the spirit and scope of the present disclosure is not to be limited by the foregoing examples, but is to be understood in the broadest sense allowable by law.

**[0207]** The use of the terms “a” and “an” and “the” and similar referents in the context of describing the disclosure (especially in the context of the following claims) is to be construed to cover both the singular and the plural, unless otherwise indicated herein or clearly contradicted by context. The terms “comprising,” “has sender-controlled contact media content item,” “including,” and “containing” are to be construed as open-ended terms (i.e., meaning “including, but not limited to,”) unless otherwise noted. Recitation of ranges of values herein are merely intended to serve as a shorthand method of referring individually to each separate value falling within the range, unless otherwise indicated herein, and each separate value is incorporated into the



specification as if it were individually recited herein. All methods described herein can be performed in any suitable order unless otherwise indicated herein or otherwise clearly contradicted by context. The use of any and all examples, or exemplary language (e.g., “such as”) provided herein, is intended merely to better illuminate the disclosure and does not pose a limitation on the scope of the disclosure unless otherwise claimed. No language in the specification should be construed as indicating any non-claimed element as essential to the practice of the disclosure.

**[0208]** While the foregoing written description enables one of ordinary skill to make and use what is considered presently to be the best mode thereof, those of ordinary skill will understand and appreciate the existence of variations, combinations, and equivalents of the specific embodiment, method, and examples herein. The disclosure should therefore not be limited by the above described embodiment, method, and examples, but by all embodiments and methods within the scope and spirit of the disclosure.

**[0209]** All documents referenced herein are hereby incorporated by reference.

What is claimed is:

1. A converged controller for interfacing a set of sources and a set of targets with credit-based flow control, the controller comprising:

- a plurality of source-oriented queues, each source-oriented queue connected to a different source of the set of sources;
- a plurality of target-oriented queues, each target-oriented queue connected to a different target of the set of targets and configured with a number of target access credits; and
- a multiplexer for selectively coupling a source-oriented queue of the plurality of source-oriented queues to at least one target-oriented queue of the plurality of target-oriented queues, wherein the coupling enables a number of data accesses between a source connected to the source-oriented queue and a subset of the set of targets connected to the at least one target-oriented queue according to the credit-based flow control;

wherein the credit-based flow control limits the number of data accesses according to a number of credits allocated to the source connected to the source-oriented queue; and

wherein the number of credits is computed from the number of target access credits of the at least one target-oriented queue.

2. The controller of claim 1, wherein the number of credits allocated to the source connected to the source-oriented queue is less than or equal to a depth of the source-oriented queue.

3. The controller of claim 2, wherein the depth of each of the plurality of source-oriented queues is less than or equal to a total depth of all the plurality of target-oriented queues.

4. The controller of claim 3, wherein at least one of the set of targets is a direct connected data storage.

5. The controller of claim 3, wherein at least one of the set of sources is an ethernet device.

6. The controller of claim 1, wherein the number of credits allocated to the source connected to the source-oriented queue is based at least in part on a size of command buffers of the subset of targets.

7. The controller of claim 3, wherein each of the plurality of target-oriented queues are sized according to a size of a command buffer of a connected target.

8. The controller of claim 1, wherein credits are allocated to the source in response a data transfer request from the source.

9. The controller of claim 1, further comprising a physical storage media controller, a physical network interface controller and a direct connection therebetween for performing data accesses between the source connected to the source-oriented queue and the subset of targets connected to the at least one target-oriented queue.

10. A method for source-oriented credit-based scheduling of data flow :

providing a set of target access credits to a plurality of target-oriented queues for accessing target resources; mapping with a multiplexer a source-oriented queue of a plurality of source-oriented queues to a portion of the plurality of target-oriented queues;

providing a set of source access credits for the source-oriented queue of the plurality of source-oriented queues responsive to a request from at least one of a plurality of source resources connected to the plurality of source-oriented queues to access the target resources; and

limiting a maximum number of source access credits for the source-oriented queue of the plurality of source-oriented queues based on a total count of target access credits provided to the portion of the plurality of target-oriented queues.

11. The method of claim 10, wherein providing the set of target access credits further comprises limiting the set of target access credits to a size that is less than or equal to a total depth of the plurality of target-oriented queues.

12. The method of claim 10, wherein at least one of the target resources is a direct connected data storage.

13. The method of claim 10, wherein at least one of the plurality of source resources is an ethernet device.

14. The method of claim 10, wherein limiting the maximum number of source access credits further comprises sizing a depth of the source-oriented queue to the maximum number of source access credits.

15. A storage control system comprising:

a plurality of source-oriented queues that each provide access credits to network-remote sources requesting access to storage resources controlled by a physical storage controller portion of a converged network-storage controller, wherein each of the network-remote sources is a distinct instance of the converged network-storage controller;

a plurality of target-oriented queues, wherein each target-oriented queue controls access to a local, physical storage resource by limiting a count of target access credits permitted for each local physical storage resource; and

a multiplexer for mapping the plurality of source-oriented queues to the plurality of target-oriented queues, wherein a maximum number of access credits permitted for each of the plurality of source-oriented queues is limited by the multiplexer to no more than a total number of target access credits available from the plurality of target-oriented queues with which each source queue of the plurality of source-oriented queues is multiplexed.

**16.** The system of claim **15**, wherein access bandwidth and access latency are guaranteed independent of a number of a local, physical storage resources.

**17.** The system of claim **15**, wherein access bandwidth and access latency are guaranteed independent of a number of converged network-storage controllers.

**18.** A method of guaranteeing predictable access latency in a network-distributed storage system, comprising:

    multiplexing a plurality of source-oriented queues to a plurality of target-oriented queues; and

    limiting a maximum size of each of the plurality of source-oriented queues to no more than a combined size of the plurality of target-oriented queues with which the plurality of source-oriented queues are multiplexed.

**19.** The method of claim **18**, further comprising allocating credits to a source coupled to a multiplexed source-oriented queue in response a data transfer request from the source.

**20.** The method of claim **18**, further comprising limiting a count of credits allocated to a source coupled to at least one of the plurality of source-oriented queues to the maximum size of each of the plurality of source-oriented queues for a credit-based flow control of data transfer.

\* \* \* \* \*