

(51) International Patent Classification:  
**G06F 13/14** (2006.01) **G06F 12/00** (2006.01)95032 (US). **AMDAHL, Carlton, G** [US/US]; 48468  
Avalon Heights Terrace, Fremont, CA 94539 (US).(21) International Application Number:  
PCT/US2011/040996(74) Agent: **SMITH, Walstein**; P.O. Box 1668, Georgetown,  
TX 78627-1668 (US).(22) International Filing Date:  
17 June 2011 (17.06.2011)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:  
61/356,443 18 June 2010 (18.06.2010) US  
61/497,525 16 June 2011 (16.06.2011) US(71) Applicant (for all designated States except US): **SAND-  
FORCE, INC.** [US/US]; 691 S. Milpitas Blvd., Ste 100,  
Milpitas, CA 95035 (US).

(72) Inventors; and

(75) Inventors/Applicants (for US only): **CANEPA, Timo-  
thy, L** [US/US]; 418 University Ave, Los Gatos, CA(81) Designated States (unless otherwise indicated, for every  
kind of national protection available): AE, AG, AL, AM,  
AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ,  
CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO,  
DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT,  
HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP,  
KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD,  
ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI,  
NO, NZ, OM, PE, PG, PH, PL, PT, RO, RS, RU, SC, SD,  
SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR,  
TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.(84) Designated States (unless otherwise indicated, for every  
kind of regional protection available): ARIPO (BW, GH,  
GM, KE, LR, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG,  
ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ,  
TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK,  
EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU,

[Continued on next page]

(54) Title: SCALABLE STORAGE DEVICES

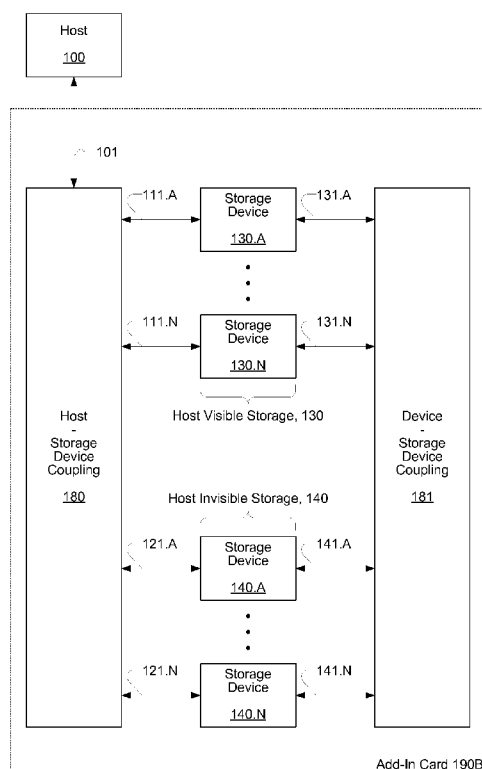


Fig. 1B

(57) Abstract: Techniques using scalable storage devices represent a plurality of host-accessible storage devices as a single logical interface, conceptually aggregating storage implemented by the devices. A primary agent of the devices accepts storage requests from the host using a host-interface protocol, processing the requests internally and/or forwarding the requests as sub-requests to secondary agents of the storage devices using a peer-to-peer protocol. The secondary agents accept and process the sub-requests, and report sub-status information for each of the sub-requests to the primary agent and/or the host. The primary agent optionally accumulates the sub-statuses into an overall status for providing to the host. Peer-to-peer communication between the agents is optionally used to communicate redundancy information during host accesses and/or failure recoveries. Various failure recovery techniques reallocate storage, reassign agents, recover data via redundancy information, or any combination thereof.



LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, **Published:**

SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ,  
GW, ML, MR, NE, SN, TD, TG).

— *without international search report and to be republished  
upon receipt of that report (Rule 48.2(g))*

**Declarations under Rule 4.17:**

— *of inventorship (Rule 4.17(iv))*

**SCALABLE STORAGE DEVICES****CROSS REFERENCE TO RELATED APPLICATIONS**

**[0001]** Priority benefit claims for this application are made in the accompanying Application Data Sheet, Request, or Transmittal (as appropriate, if any). To the extent permitted by the type of the instant application, this application incorporates by reference for all purposes the following applications, all owned by the owner of the instant application:

U.S. Provisional Application Serial No.: 61/356,443 (Docket No. SF-10-05), filed 18 June 2010, first named inventor Timothy Lawrence Canepa, and entitled **SCALABLE STORAGE DEVICES**; and

U.S. Provisional Application Serial No.: 61/497,525 (Docket No. SF-10-05B), filed 16 June 2011, first named inventor Timothy Lawrence Canepa, and entitled **SCALABLE STORAGE DEVICES**.

**BACKGROUND**

**[0002]** Field: Advancements in accessing of storage devices are needed to provide improvements in performance, efficiency, and utility of use.

**[0003]** Related Art: Unless expressly identified as being publicly or well known, mention herein of techniques and concepts, including for context, definitions, or comparison purposes, should not be construed as an admission that such techniques and concepts are previously publicly known or otherwise part of the prior art. All references cited herein (if any), including patents, patent applications, and publications, are hereby incorporated by reference in their entireties, whether specifically incorporated or not, for all purposes.

## SYNOPSIS

**[0004]** The invention may be implemented in numerous ways, including as a process, an article of manufacture, an apparatus, a system, a composition of matter, and a computer readable medium such as a computer readable storage medium (e.g. media in an optical and/or magnetic mass storage device such as a disk, or an integrated circuit having non-volatile storage such as flash storage) or a computer network wherein program instructions are sent over optical or electronic communication links. In this specification, these implementations, or any other form that the invention may take, may be referred to as techniques. The Detailed Description provides an exposition of one or more embodiments of the invention that enable improvements in performance, efficiency, and utility of use in the field identified above. The Detailed Description includes an Introduction to facilitate the more rapid understanding of the remainder of the Detailed Description. The Introduction includes Example Embodiments of one or more of systems, methods, articles of manufacture, and computer readable media in accordance with the concepts described herein. As is discussed in more detail in the Conclusions, the invention encompasses all possible modifications and variations within the scope of the issued claims.

## Brief Description of Drawings

**[0005]** Fig. 1A illustrates selected structural details of an embodiment of a technique for scalable storage devices, including a host, host visible storage having one or more storage devices operable as respective primary agents, and host invisible storage having one or more storage devices operable as respective secondary agents.

**[0006]** Fig. 1B illustrates selected structural details of another embodiment of a technique for scalable storage devices, including a host, host visible storage having one or more storage devices operable as respective primary agents, and host invisible storage having one or more storage devices operable as respective secondary agents.

**[0007]** Fig. 2 illustrates selected processing details of an embodiment of a technique for scalable storage devices, including actions performed by a primary agent and one or more secondary agents.

**[0008]** Fig. 3 illustrates selected details of host, primary agent, and secondary agent addressing in an embodiment of a technique for scalable storage devices.

**[0009]** Fig. 4 illustrates selected structural details of an embodiment of a scalable storage device enabled to operate as a primary agent.

**[0010]** Fig. 5 illustrates selected structural details of an embodiment of a scalable storage device enabled to operate as a secondary agent.

## List of Reference Symbols in Drawings

[0011]

Ref. Symbol	Element Name
<b>100</b>	Host
<b>101</b>	coupling
<b>110</b>	Host Visible Storage
<b>110.A</b>	Storage Device
<b>110.N</b>	Storage Device
<b>111.A</b>	coupling
<b>111.N</b>	coupling
<b>120</b>	Host Invisible Storage
<b>120.A</b>	Storage Device
<b>120.N</b>	Storage Device
<b>121.A</b>	coupling
<b>121.N</b>	coupling
<b>130</b>	Host Visible Storage
<b>130.A</b>	Storage Device
<b>130.N</b>	Storage Device
<b>131.A</b>	coupling
<b>131.N</b>	coupling
<b>140</b>	Host Invisible Storage
<b>140.A</b>	Storage Device
<b>140.N</b>	Storage Device
<b>141.A</b>	coupling
<b>141.N</b>	coupling
<b>151</b>	dashed-arrow
<b>152</b>	dashed-arrow
<b>153</b>	dashed-arrow
<b>154</b>	dashed-arrow
<b>180</b>	Host - Storage Device Coupling
<b>181</b>	Device - Storage Device Coupling
<b>190A</b>	Add-In Card
<b>190B</b>	Add-In Card

<b>Ref. Symbol</b>	<b>Element Name</b>
<b>201</b>	Start
<b>202</b>	Accept Req from Host
<b>203</b>	Local?
<b>203N</b>	No
<b>203Y</b>	Yes
<b>204</b>	Xlate to Local LBAs
<b>205</b>	Process Locally
<b>206</b>	Provide Status to Host
<b>207</b>	Forward to Secondary(s)
<b>207R</b>	Sub-Req
<b>208</b>	Wait for Completion
<b>209A</b>	Accept Sub-Status from Secondary(s)
<b>209</b>	Primary Actions
<b>211</b>	Accept Sub-Req from Primary
<b>212</b>	Xlate to Local LBAs
<b>213</b>	Process Locally
<b>214</b>	Provide Sub-Status to Primary
<b>214S</b>	Sub-Status
<b>219</b>	Secondary Actions
<b>299</b>	End
<b>310</b>	Host Address Space
<b>311</b>	Host LBA Range 1
<b>312</b>	Host LBA Range 2
<b>313</b>	Host LBA Range 3
<b>314</b>	Host LBA Range 4
<b>315</b>	Host LBA Range 5
<b>320</b>	Primary Address Space
<b>321</b>	Primary LBA Range 1
<b>322</b>	Primary LBA Range 2
<b>323</b>	Primary LBA Range 3
<b>330</b>	Secondary Address Space
<b>331A</b>	Secondary A LBA Range 1
<b>331B</b>	Secondary B LBA Range 1
<b>332A</b>	Secondary A LBA Range 2

<b>Ref. Symbol</b>	<b>Element Name</b>
<b>332B</b>	Secondary B LBA Range 2
<b>333A</b>	Secondary A LBA Range 3
<b>401</b>	PCIe Intfc
<b>402</b>	LBA, Length Xlate
<b>403</b>	Storage Intfc
<b>404</b>	Mass Storage
<b>404F</b>	Flash
<b>404M</b>	Magnetic
<b>404O</b>	Optical
<b>405</b>	Sub-Request Generation & Sub-Status Accumulation
<b>501</b>	Sub-Request Accepting & Sub-Status Generation



## DETAILED DESCRIPTION

**[0012]** A detailed description of one or more embodiments of the invention is provided below along with accompanying figures illustrating selected details of the invention. The invention is described in connection with the embodiments. The embodiments herein are understood to be merely exemplary, the invention is expressly not limited to or by any or all of the embodiments herein, and the invention encompasses numerous alternatives, modifications, and equivalents. To avoid monotony in the exposition, a variety of word labels (including but not limited to: first, last, certain, various, further, other, particular, select, some, and notable) may be applied to separate sets of embodiments; as used herein such labels are expressly not meant to convey quality, or any form of preference or prejudice, but merely to conveniently distinguish among the separate sets. The order of some operations of disclosed processes is alterable within the scope of the invention. Wherever multiple embodiments serve to describe variations in process, method, and/or program instruction features, other embodiments are contemplated that in accordance with a predetermined or a dynamically determined criterion perform static and/or dynamic selection of one of a plurality of modes of operation corresponding respectively to a plurality of the multiple embodiments. Numerous specific details are set forth in the following description to provide a thorough understanding of the invention. The details are provided for the purpose of example and the invention may be practiced according to the claims without some or all of the details. For the purpose of clarity, technical material that is known in the technical fields related to the invention has not been described in detail so that the invention is not unnecessarily obscured.

## INTRODUCTION

**[0013]** This introduction is included only to facilitate the more rapid understanding of the Detailed Description; the invention is not limited to the concepts presented in the introduction (including explicit examples, if any), as the paragraphs of any introduction are necessarily an abridged view of the entire subject and are not meant to be an exhaustive or restrictive description. For example, the introduction that follows provides overview information limited by space and organization to only certain embodiments. There are many other embodiments, including those to which claims will ultimately be drawn, discussed throughout the balance of the specification.

1 Acronyms

2

3 **[0014]** At least some of the various shorthand abbreviations, or acronyms, defined here  
 4 refer to certain elements used herein.

<u>Acronym</u>	<u>Description</u>
AHCI	Advanced Host Controller Interface
CF	Compact Flash
DIF	Data Integrity Field
DIX	Data Integrity eXtension
DMA	Direct Memory Access
eNVMHCI	Enterprise Non-Volatile Memory Host Controller Interface
eSATA	external Serial Advanced Technology Attachment
IC	Integrated Circuit
IDE	Integrated Drive Electronics
IO	Input/Output
JBOD	Just a Bunch Of Disks
LBA	Logical Block Address
MMC	MultiMediaCard
MsgD	Message Request with Data payload
PC	Personal Computer
PCIe	Peripheral Component Interconnect express (PCI express)
PDA	Personal Digital Assistant
RAID	Redundant Array of Inexpensive/Independent Disks
RMW	Read-Modify-Write
ROM	Read Only Memory
SAS	Serial Attached Small Computer System Interface (Serial SCSI)
SATA	Serial Advanced Technology Attachment (Serial ATA)
SCSI	Small Computer System Interface
SD	Secure Digital
SSD	Solid-State Disk/Drive
USB	Universal Serial Bus

5

6

7 **[0015]** Techniques using scalable storage devices represent a plurality of host-  
 8 accessible storage devices as a single logical interface, conceptually aggregating storage

1 implemented by the devices. A primary agent of the devices accepts storage requests from the  
2 host using a host-interface protocol, processing the requests internally and/or forwarding the  
3 requests as sub-requests to secondary agents of the storage devices using a peer-to-peer protocol.  
4 The secondary agents accept and process the sub-requests, and report sub-status information for  
5 each of the sub-requests to the primary agent and/or the host. The primary agent optionally  
6 accumulates the sub-statuses into an overall status for providing to the host. The primary agent  
7 reports available storage to the host that includes storage implemented by the primary agent as  
8 well as storage implemented by the secondary agents and allocated to aggregating. The  
9 secondary agents report zero available storage to the host or alternatively storage remaining after  
10 the aggregation allocation.

11  
12 **[0016]** Peer-to-peer communication between the agents is optionally used to  
13 communicate redundancy information during host accesses and/or failure recoveries. Various  
14 failure recovery techniques reallocate storage, reassign primary/secondary/configurable agents,  
15 recover data via redundancy information, or any combination thereof.

16  
17 **[0017]** PCIe interfaces have an inherent singular relationship to a host in a system.  
18 Each PCIe interface represents a single portal to a function of a device, whether the function is  
19 storage, networking, or some other computing system capability. Each PCIe device in the  
20 system represents a respective distinct entity in the system. However, through peer-to-peer  
21 traffic, several individual PCIe devices are grouped together so that the devices appear as a  
22 single PCIe device to the host. By leveraging PCIe peer-to-peer capability, a scalable storage  
23 architecture is enabled by ganging and/or grouping together a plurality of PCIe storage devices  
24 via peer-to-peer communication enabling the ganged/grouped devices to appear as a single  
25 storage subsystem to the host. In some embodiments, requests are distributed among the PCIe  
26 storage devices based only on LBAs. The LBA-only distribution does not require that the  
27 storage of any particular device in the group be functional for the particular device to route  
28 requests, thus eliminating the storage as a single point of failure.

29  
30 **[0018]** In some embodiments and/or usage scenarios, the ganging/grouping enables  
31 increased performance and/or capacity while appearing to the host as a single logical interface.  
32 In some embodiments and/or usage scenarios, all single points of storage failure are eliminated  
33 via LBA-based request routing. In some embodiments and/or usage scenarios, the  
34 ganging/grouping distributes request load across a number of storage devices, thus scaling  
35 capacity and performance transparently to a host. The capacity and performance scaling are  
36 transparent to the host in the sense that the host does not need to control and requires no

1 observability of the ganging/grouping, other than capacity and latency changes. In some  
2 embodiments and/or usage scenarios, the ganging/grouping provides for a symmetric multi-  
3 processor storage complex.  
4

5 **[0019]** In some embodiments and/or usage scenarios, the ganging/grouping enables a  
6 host to access, as a single logical device, an aggregation of resources distributed over a plurality  
7 of physical devices. For example, the ganging/grouping enables a host to access, as a single  
8 logical storage device, an aggregation of a plurality of storage ranges distributed over a plurality  
9 of physical storage devices. Forwarding an LBA (and optionally a length) as provided by the  
10 host as part of the request, preserves the LBA (and optional length) information. Preserving the  
11 LBA (and optional length) information and translating on a target device, enables maintaining  
12 accurate integrity metadata (such as DIV/DIX referential integrity) and/or enables independent  
13 (e.g. parallel) operation of devices.  
14

15 **[0020]** In some embodiments and/or usage scenarios, the ganging/grouping enables  
16 transparent redundancy and/or recovery with respect to storage capabilities, such as various  
17 RAID, mirroring, and fail-over implementations. In various embodiments, host traffic is  
18 predominantly data traffic while peer-to-peer traffic is predominantly control traffic. In various  
19 embodiments, peer-to-peer traffic includes data (e.g. during a RAID recovery operation, or  
20 during transport of accumulated parity information).  
21

22 **[0021]** In various embodiments and/or usage scenarios, the ganging/grouping enables  
23 systems with a host coupled to SSD(s) where bandwidths of one or more PCIe links coupled to  
24 the host (each having one or more lanes) are distributed (transparently to the host) across a  
25 plurality of SSDs, a plurality of flash controllers (such as used in SSDs), and/or a plurality of  
26 flash chips. For example, bandwidth of one host-coupled PCIe link (e.g. having eight lanes) is  
27 distributed across four SSDs or alternatively eight flash controllers (each SSD or flash controller  
28 having one PCIe link with, e.g., fewer than eight lanes). For another example, bandwidth of four  
29 host-coupled PCIe links is distributed across 64 SSDs or alternatively 256 flash controllers. For  
30 another example, bandwidth of one host-coupled PCIe link is distributed across 32 flash chips or  
31 alternatively across 64 flash chips. For another example, bandwidth of four host-coupled PCIe  
32 links is distributed across 128 flash chips or alternatively across 256 flash chips.  
33

34 **[0022]** In various embodiments, a primary agent determines host to secondary agent  
35 storage mapping (e.g. which host LBA ranges correspond to which secondary agents). For  
36 example, host LBAs are striped across physical drives at a pre-determined size (e.g. 64KB). For

another example, host LBAs are striped according to a RAID implementation (e.g. RAID 0, RAID 1, or RAID 5). For yet another example, striping is changed dynamically according to workload (e.g. a first striping is used for a first workload and a second striping is used for a second workload).

**[0023]** In some embodiments, a group of PCIe storage devices (accessible from one or more hosts) are coupled via a bus topology enabling peer-to-peer communication with each other. The group of PCIe storage devices is configured so that only one PCIe device in the group receives requests for a particular host logical connection. In some embodiments, there is a plurality of connections to the group, enabling different host logical connections to supply requests that are received by different PCIe storage devices in the group.

**[0024]** Configuration of the group of PCIe storage devices depends on how the devices are exposed to the host(s). If all of the devices in the group are directly exposed to the host, e.g., via a transparent switch, then all devices but the ones enabling host logical connections report zero storage capacity back to the host(s), while the devices enabling host logical connections report a configured capacity of the group for a particular host logical connection. The configured capacity depends on how storage of the group is partitioned across host logical connections as well as how the storage is configured or organized (such as JBOD or RAID). If some of the devices are hidden from the host behind, e.g., a bridge or a non-transparent switch, then the zero capacity reporting is not performed by the devices not exposed to the host. A host logical connection only sends requests directly to a single PCIe device interface that then distributes and coordinates requests to the other devices in the group.

**[0025]** In further embodiments, two or more devices enabling respective host logical connections report storage capacity that is overlapped within storage capacities of the devices and/or other devices. For example, a first device enabling a first host logical connection reports an aggregate storage capacity equal to a total storage capacity of the first device and a total storage capacity of a second device. A third device enabling a second host logical connection reports the aggregate storage capacity as well. Host requests received via the first host logical connection or received via the second host logical connection reference same physical storage of the first device and of the second device (optionally and/or selectively using different LBAs).

**[0026]** Storage implemented by a storage device need not be entirely exposed (or not exposed) directly to a host. For example, a particular storage device implements storage that is partially allocated to a group of storage devices (e.g. not directly exposed to a host) and that is

1 partially allocated for direct exposure to the host. The particular storage device, rather than  
2 reporting zero capacity to the host, instead reports an amount of storage corresponding to the  
3 partial allocation for direct exposure to the host.

4  
5 **[0027]** In some embodiments, initial configuration of devices of the group is performed  
6 by software residing on the host(s) and/or in one or more option ROMs. An example of the  
7 initial configuration is which of the devices is to enable host logical connections and which are  
8 not. Another example of the initial configuration is how much storage to allocate to a group of  
9 devices versus for direct exposure to a host.

10  
11 **[0028]** In some embodiments, configuration of devices of the group is performed while  
12 a system is operating, such as when a “hot spare” is inserted or when a device is added to or  
13 removed from the group dynamically, e.g. in response to a failure. For example, when a hot  
14 spare device is inserted, a discovery process involving one or more of a device driver on a host  
15 and a primary agent recognizes the inserted device, and configures the inserted device  
16 accordingly (such as to replace a failed device). The insertion of the device is indicated via, for  
17 instance, an announcement from the inserted device and/or a “hotplug” event.

18  
19 **[0029]** When a host issues a storage read/write request (such as an eNVMHCI or  
20 AHCI-compliant request), the device receiving the request uses an LBA (and optionally length)  
21 of the request to forward the request to device(s) of the group that implement the storage  
22 requested. Depending on the length of the request and the organization of the group (e.g. JBOD  
23 or RAID), the request is routed to more than one device of the group, for example as a plurality  
24 of sub-requests.

25  
26 **[0030]** According to various embodiments, when a particular device in the group  
27 receives a sub-request, the particular device translates the LBA (and optionally any supplied  
28 length) to a local LBA (and length), processes the sub-request, and then transfers data for the  
29 sub-request directly to/from host memory. An example of host memory is an element enabled to  
30 retain data written by the host via execution of an instruction that stores data to memory at a  
31 particular address and enabled for the host to read the written data via execution of an instruction  
32 that reads data from memory at the particular address. The particular device optionally and/or  
33 selectively obtains a host memory address of where to transfer the data from one or more of: (a)  
34 from the (forwarded) sub-request, (b) from host resident translation structures (such as  
35 scatter/gather lists), and (c) by issuing proxy requests to a device acting as a control interface to

1 the host. In some embodiments, the issuing of proxy requests is only performed for IO  
2 virtualization.

3  
4 **[0031]** When the data transfer for each sub-request is complete, the transferring device  
5 reports sub-status to the device acting as the control interface to the host. The device acting as  
6 the control interface to the host then posts status to the host when all of the devices involved in a  
7 request have reported sub-statuses.

8  
9 **[0032]** Non-data commands (for example a SATA IDENTIFY DEVICE command) are  
10 processed similarly. When a non-data command is received by a particular device acting as a  
11 control interface to a host, the particular device interrogates the other devices in the group (if  
12 necessary), and then reports back summarized and/or aggregated results of the interrogation to  
13 the host.

14  
15 **[0033]** In various embodiments, a host sees an aggregation of a plurality of devices  
16 (such as devices each having an x1 or x2 PCIe connection) as a single logical device  
17 encompassing storage of all the devices (e.g. the devices are virtualized as the single logical  
18 device). In various embodiments, the host sees any or all of the devices as having portions of  
19 overall storage, enabling splitting control traffic in addition to data traffic.

20  
21 **[0034]** A request from the host via a submission queue (such as one submission queue  
22 of a plurality of simultaneously active submission queues for parallel streams of activity) is  
23 processed, e.g., by a particular one of the devices that the host is aware of. An individual request  
24 in the submission queue is optionally and/or selectively forwarded as a sub-request (via peer-to-  
25 peer requests from the particular device) to the devices that implement all or portions of storage  
26 referenced by the request. For example, the storage is striped among the devices by individual  
27 or groups of LBAs. In some embodiments and/or usage scenarios, striping in groups of LBAs  
28 enables performance benefits via aggregation of larger chunks to individual devices. In some  
29 embodiments having a plurality of submission queues, the host notifies a particular one of the  
30 devices of new submission queues, and the particular device distributes servicing of the new  
31 submission queues to others of the devices.

32  
33 **[0035]** The individual devices process the forwarded sub-requests, including data  
34 transfer to/from the host, independently. Some protocols (such as some PCIe compatible  
35 protocols) use “completions” that are provided via a respective completion queue for each  
36 submission queue. For each entry in a submission queue, a corresponding completion

1 notification is sent back to the respective completion queue. Any one or more of the devices are  
2 optionally and/or selectively enabled to “aggregate” completion information and update the  
3 completion queue for a given request. In various embodiments, one device is assigned per  
4 completion queue to enable atomicity. The assigned device is, in some embodiments, the same  
5 device that processes entries in the corresponding submission queue, while in other  
6 embodiments, the assigned device is different than the device that processes the entries in the  
7 corresponding submission queue.

#### 10 EXAMPLE EMBODIMENTS

12 **[0036]** In concluding the introduction to the detailed description, what follows is a  
13 collection of example embodiments, including at least some explicitly enumerated as “ECs”  
14 (Example Combinations), providing additional description of a variety of embodiment types in  
15 accordance with the concepts described herein; these examples are not meant to be mutually  
16 exclusive, exhaustive, or restrictive; and the invention is not limited to these example  
17 embodiments but rather encompasses all possible modifications and variations within the scope  
18 of the issued claims.



- 1     **[0037]**           EC1) A system comprising:  
2           means for accepting a request from a host to access storage at an address;  
3           means for determining, based at least in part on the address, which one or more of a  
4           plurality of storage devices the request corresponds to;  
5           means for determining, based at least in part on the address, a respective sub-request  
6           corresponding to each of the determined storage devices;  
7           means for sending the sub-requests to the corresponding determined storage devices;  
8           means for accepting respective sub-status from each of the determined storage devices;  
9           means for determining an overall status based at least in part on each of the respective  
10          sub-statuses;  
11          means for providing the overall status to the host;  
12          wherein the means for accepting the request is compatible with a host-interface protocol  
13               and the means for sending the sub-requests and the means for accepting the sub-  
14               statuses are compatible with a peer-to-peer protocol operable independently of  
15               the host; and  
16          wherein at least one of the storage devices implements storage that is accessible via one  
17               of the sub-requests and the at least one of the storage devices disables accesses  
18               via the host-interface protocol to the implemented storage.  
19
- 20     **[0038]**           EC2) The system of EC1, wherein two or more of the means for accepting a  
21     request, the means for sending the sub-requests, the means for accepting respective sub-status,  
22     and the means for providing the overall status are operable on a shared physical interface.  
23
- 24     **[0039]**           EC3) The system of EC2, wherein at least one of a Peripheral Component  
25     Interconnect express (PCIe) topology and an Infiniband topology comprises the shared physical  
26     interface.  
27
- 28     **[0040]**           EC4) The system of EC2, wherein at least one of a Peripheral Component  
29     Interconnect express (PCIe) interface and an Infiniband interface comprises the shared physical  
30     interface.  
31

1   **[0041]**           EC5) The system of EC1 further comprising a particular one of the storage  
2   devices that comprises the means for accepting the request, the means for determining the  
3   storage devices, the means for determining the sub-requests, the means for sending the sub-  
4   requests, the means for accepting the sub-statuses, the means for determining the overall status,  
5   and the means for providing the overall status.

6  
7   **[0042]**           EC6) The system of EC1, wherein two or more of the storage devices are  
8   enabled to communicate with the host via the host interface protocol.

9  
10   **[0043]**           EC7) The system of EC6, wherein the two or more of the storage devices are all  
11   of the storage devices.

12  
13   **[0044]**           EC8) The system of EC1, wherein the address comprises a starting location and  
14   a length.

15  
16   **[0045]**           EC9) The system of EC1, wherein the address is a Logical Block Address  
17   (LBA).

18  
19   **[0046]**           EC10) The system of EC1, wherein the means are collectively implemented in  
20   a single Integrated Circuit (IC).

21  
22   **[0047]**           EC11) The system of EC1, wherein the means are collectively implemented in  
23   a single add-in card.

24  
25   **[0048]**           EC12) The system of EC1, wherein the means are comprised in a Solid-State  
26   Disk (SSD).

1     **[0049]**           EC13) The system of EC1, wherein the host-interface protocol is compatible  
2     with one or more of  
3           a Universal Serial Bus (USB) interface standard,  
4           a Compact Flash (CF) interface standard,  
5           a MultiMediaCard (MMC) interface standard,  
6           a Secure Digital (SD) interface standard,  
7           a Memory Stick interface standard,  
8           an xD-picture card interface standard,  
9           an Integrated Drive Electronics (IDE) interface standard,  
10          a Serial Advanced Technology Attachment (SATA) interface standard,  
11          an external SATA (eSATA) interface standard,  
12          a Small Computer System Interface (SCSI) interface standard,  
13          a Serial Attached Small Computer System Interface (SAS) interface standard,  
14          a Fibre Channel interface standard,  
15          an Ethernet interface standard, and  
16          a Peripheral Component Interconnect express (PCIe) interface standard.

17  
18   **[0050]**           EC14) The system of EC1, further comprising all or any portions of the host.

19  
20   **[0051]**           EC15) The system of EC14, wherein the host comprises one or more of  
21          a computer,  
22          a workstation computer,  
23          a server computer,  
24          a storage server,  
25          a Personal Computer (PC),  
26          a laptop computer,  
27          a notebook computer,  
28          a netbook computer,  
29          a Personal Digital Assistant (PDA),  
30          a media player,  
31          a media recorder,  
32          a digital camera,  
33          a cellular handset,  
34          a cordless telephone handset, and  
35          an electronic game.

36

1     **[0052]**           EC16) The system of EC1, further comprising all or any portions of the storage  
2     devices.

4     **[0053]**           EC17) The system of EC1, further comprising flash memory comprised in one  
5     or more of the storage devices.

7     **[0054]**           EC18) A computer readable medium having a set of instructions stored therein  
8     that when executed by a processing element cause the processing element to perform operations  
9     comprising:

10           managing accepting a request from a host to access storage at an address;

11           managing determining, based at least in part on the address, which one or more of a  
12           plurality of storage devices the request corresponds to;

13           managing determining, based at least in part on the address, a respective sub-request  
14           corresponding to each of the determined storage devices;

15           managing sending the sub-requests to the corresponding determined storage devices;

16           managing accepting respective sub-status from each of the determined storage devices;

17           managing determining an overall status based at least in part on each of the respective  
18           sub-statuses;

19           managing providing the overall status to the host;

20           wherein the accepting the request is via a host-interface protocol and the sending the  
21           sub-requests and the accepting the sub-statuses are via a peer-to-peer protocol  
22           operable independently of the host; and

23           wherein at least one of the storage devices implements storage that is accessible via one  
24           of the sub-requests and the at least one of the storage devices disables accesses  
25           via the host-interface protocol to the implemented storage.

27     **[0055]**           EC19) The computer readable medium of EC18, wherein the host-interface and  
28     the peer-to-peer protocols are compatible with a Peripheral Component Interconnect express  
29     (PCIe) channel.

1     **[0056]**           EC20) A method comprising:  
2           accepting a request from a host to access storage at an address;  
3           determining, based at least in part on the address, which one or more of a plurality of  
4           storage devices the request corresponds to;  
5           determining, based at least in part on the address, a respective sub-request corresponding  
6           to each of the determined storage devices;  
7           sending the sub-requests to the corresponding determined storage devices;  
8           accepting respective sub-status from each of the determined storage devices;  
9           determining an overall status based at least in part on each of the respective sub-statuses;  
10          providing the overall status to the host;  
11          wherein the accepting the request is compatible with a host-interface protocol and the  
12           sending the sub-requests and the accepting the sub-statuses are compatible with  
13           a peer-to-peer protocol operable independently of the host; and  
14          wherein at least one of the storage devices implements storage that is accessible via one  
15           of the sub-requests and the at least one of the storage devices disables accesses  
16           via the host-interface protocol to the implemented storage.  
17

1     **[0057]**           EC21) A system comprising:  
2           host interface logic hardware enabled to accept a request from a host to access storage at  
3           an address;  
4           control logic hardware enabled  
5           to determine, based at least in part on the address, which one or more of a  
6           plurality of storage devices the request corresponds to and a respective  
7           sub-request corresponding to each of the determined storage devices,  
8           to send the sub-requests to the corresponding determined storage devices,  
9           to accept respective sub-status from each of the determined storage devices, and  
10          to determine an overall status based at least in part on each of the respective  
11          sub-statuses;  
12          wherein the host interface logic hardware is further enabled to provide the overall status  
13          to the host;  
14          wherein the accepting the request is compatible with a host-interface protocol and the  
15          sending the sub-requests and the accepting the sub-statuses are compatible with  
16          a peer-to-peer protocol operable independently of the host; and  
17          wherein at least one of the storage devices implements storage that is accessible via one  
18          of the sub-requests and the at least one of the storage devices disables accesses  
19          via the host-interface protocol to the implemented storage.  
20

1     **[0058]**           EC22) A method comprising:  
2           a first storage device reading a first request from a first one of a plurality of host request  
3           queues;  
4           a second storage device reading a second request from a second one of the host request  
5           queues, the second host request queue being distinct from the first host request  
6           queue, the reading by the second storage device being independent of the  
7           reading by the first storage device;  
8           the first storage device forwarding at least a portion of the first request to the second  
9           storage device;  
10          the second storage device forwarding at least a portion of the second request to the first  
11          storage device, the forwarding by the second storage device being independent  
12          of the forwarding by the first storage device;  
13          the first storage device accessing first storage of the first storage device in accordance  
14          with the at least the portion of the second request and sending a first sub-status,  
15          based at least in part on the accessing by the first storage device, to the second  
16          storage device; and  
17          the second storage device accessing second storage of the second storage device in  
18          accordance with the at least the portion of the first request and sending a second  
19          sub-status, based at least in part on the accessing by the second storage device,  
20          to the first storage device.  
21

1     **[0059]**           EC23) A system comprising:  
2           means for a first storage device to read a first request from a first one of a plurality of  
3           host request queues;  
4           means for a second storage device to read a second request from a second one of the  
5           host request queues, the second host request queue being distinct from the first  
6           host request queue, the reading by the second storage device being independent  
7           of the reading by the first storage device;  
8           means for the first storage device to forward at least a portion of the first request to the  
9           second storage device;  
10          means for the second storage device to forward at least a portion of the second request to  
11          the first storage device, the forwarding by the second storage device being  
12          independent of the forwarding by the first storage device;  
13          means for the first storage device to access first storage of the first storage device in  
14          accordance with the at least the portion of the second request and to send a first  
15          sub-status, based at least in part on the accessing by the first storage device, to  
16          the second storage device; and  
17          means for the second storage device to access second storage of the second storage  
18          device in accordance with the at least the portion of the first request and to send  
19          a second sub-status, based at least in part on the accessing by the second storage  
20          device, to the first storage device.  
21



1     **[0060]**           EC24) A system comprising:  
2           host interface logic hardware of a first storage device enabled to read a first request from  
3           a first one of a plurality of host request queues;  
4           host interface logic hardware of a second storage device enabled to read a second  
5           request from a second one of the host request queues, the second host request  
6           queue being distinct from the first host request queue, the reading by the second  
7           storage device being independent of the reading by the first storage device;  
8           request forwarding logic hardware of the first storage device enabled to forward at least  
9           a portion of the first request to the second storage device;  
10          request forwarding logic hardware of the second storage device enabled to forward at  
11          least a portion of the second request to the first storage device, the forwarding  
12          by the second storage device being independent of the forwarding by the first  
13          storage device;  
14          storage interface logic hardware of the first storage device enabled to access first storage  
15          of the first storage device in accordance with the at least the portion of the  
16          second request;  
17          sub-status generation logic hardware of the first storage device enabled to determine a  
18          first sub-status, based at least in part on the accessing by the first storage device;  
19          sub-status forwarding logic hardware of the first storage device enabled to forward the  
20          first sub-status to the second storage device;  
21          storage interface logic hardware of the second storage device enabled to access second  
22          storage of the second storage device in accordance with the at least the portion  
23          of the first request;  
24          sub-status generation logic hardware of the second storage device enabled to determine  
25          a second sub-status, based at least in part on the accessing by the second storage  
26          device; and  
27          sub-status forwarding logic hardware of the second storage device enabled to forward  
28          the second sub-status to the first storage device.  
29

1     **[0061]**           EC25) A method comprising:  
2           a first storage device reading a first request from a first one of a plurality of host request  
3           queues;  
4           a second storage device reading a second request from a second one of the host request  
5           queues, the second host request queue being distinct from the first host request  
6           queue, the reading by the second storage device being independent of the  
7           reading by the first storage device;  
8           the first storage device forwarding at least a portion of the first request to the second  
9           storage device;  
10          the second storage device forwarding at least a portion of the second request to the first  
11          storage device, the forwarding by the second storage device being independent  
12          of the forwarding by the first storage device;  
13          the first storage device accessing first storage in accordance with the at least the portion  
14          of the second request and writing a first status, based at least in part on the  
15          accessing by the first storage device, to a first one of a plurality of host status  
16          queues; and  
17          the second storage device accessing second storage in accordance with the at least the  
18          portion of the first request and writing a second status, based at least in part on  
19          the accessing by the second storage device, to a second one of the host status  
20          queues, the accessing by the second storage device being independent of the  
21          accessing by the first storage device, the second storage being distinct from the  
22          first storage.  
23

1     **[0062]**           EC26) A system comprising:  
2           means for a first storage device to read a first request from a first one of a plurality of  
3           host request queues;  
4           means for a second storage device to read a second request from a second one of the  
5           host request queues, the second host request queue being distinct from the first  
6           host request queue, the reading by the second storage device being independent  
7           of the reading by the first storage device;  
8           means for the first storage device to forward at least a portion of the first request to the  
9           second storage device;  
10          means for the second storage device to forward at least a portion of the second request to  
11          the first storage device, the forwarding by the second storage device being  
12          independent of the forwarding by the first storage device;  
13          means for the first storage device to access first storage in accordance with the at least  
14          the portion of the second request and to write a first status, based at least in part  
15          on the accessing by the first storage device, to a first one of a plurality of host  
16          status queues; and  
17          means for the second storage device to access second storage in accordance with the at  
18          least the portion of the first request and to write a second status, based at least in  
19          part on the accessing by the second storage device, to a second one of the host  
20          status queues, the accessing by the second storage device being independent of  
21          the accessing by the first storage device, the second storage being distinct from  
22          the first storage.  
23

1     **[0063]**           EC27) A system comprising:  
2           host interface logic hardware of a first storage device enabled to read a first request from  
3           a first one of a plurality of host request queues;  
4           host interface logic hardware of a second storage device enabled to read a second  
5           request from a second one of the host request queues, the second host request  
6           queue being distinct from the first host request queue, the reading by the second  
7           storage device being independent of the reading by the first storage device;  
8           request forwarding logic hardware of the first storage device enabled to forward at least  
9           a portion of the first request to the second storage device;  
10          request forwarding logic hardware of the second storage device enabled to forward at  
11          least a portion of the second request to the first storage device, the forwarding  
12          by the second storage device being independent of the forwarding by the first  
13          storage device;  
14          storage interface logic hardware of the first storage device enabled to access first storage  
15          in accordance with the at least the portion of the second request;  
16          status generation logic hardware of the first storage device enabled to determine a first  
17          status, based at least in part on the accessing by the first storage device, and  
18          further enabled to provide the first status to the host interface logic hardware of  
19          the first storage device to be written to a first one of a plurality of host status  
20          queues;  
21          storage interface logic hardware of the second storage device enabled to access second  
22          storage in accordance with the at least the portion of the first request; and  
23          status generation logic hardware of the second storage device enabled to determine a  
24          second status, based at least in part on the accessing by the second storage  
25          device, and further enabled to provide the second status to the host interface  
26          logic hardware of the second storage device to be written to a second one of the  
27          host status queues, the accessing by the second storage device being  
28          independent of the accessing by the first storage device, the second storage  
29          being distinct from the first storage.  
30

- 1     **[0064]**           EC28) A system comprising:  
2           means for enabling a host to access, as a single logical interface, an aggregation of a  
3           plurality of storage ranges distributed over a plurality of physical storage  
4           devices, one of the physical storage devices enabled to operate as a primary  
5           agent and one or more of the physical storage devices enabled to operate as  
6           secondary agents;  
7           means for communicating peer-to-peer traffic between the primary agent and at least  
8           one of the secondary agents;  
9           wherein the primary agent is enabled to manage at least a portion of the secondary  
10          agents for redundant data storage; and  
11          wherein the peer-to-peer traffic comprises redundancy traffic to implement the  
12          redundant data storage.
- 13
- 14   **[0065]**           EC29) The system of EC28, wherein at least a portion of the redundancy traffic  
15   comprises accumulated parity information.
- 16
- 17   **[0066]**           EC30) The system of EC29, wherein the at least the portion is forwarded from  
18   the primary agent and/or one of the secondary agents.
- 19
- 20   **[0067]**           EC31) The system of EC29, wherein the at least the portion is forwarded to a  
21   one of the secondary agents that implements accumulated parity storage corresponding to the  
22   accumulated parity information.
- 23
- 24   **[0068]**           EC32) The system of EC28, wherein at least a portion of the redundancy traffic  
25   is provided by the primary agent to one or more of the secondary agents.
- 26
- 27   **[0069]**           EC33) The system of EC28, wherein at least a portion of the redundancy traffic  
28   is provided by one or more of the secondary agents to the primary agent.
- 29
- 30   **[0070]**           EC34) The system of EC28, wherein at least a portion of the redundancy traffic  
31   is between a plurality of the secondary agents.
- 32
- 33   **[0071]**           EC35) The system of EC28, wherein the redundancy traffic comprises one or  
34   more of control information, untransformed data, and transformed data.
- 35

1   **[0072]**           EC36) The system of EC35, wherein the control information comprises an  
2   indication of how many redundancy updates are to occur for a particular storage address.

4   **[0073]**           EC37) The system of EC36, wherein at least one of the primary agent and the  
5   secondary agents is enabled to use the indication to cache information relating to the redundancy  
6   updates until the redundancy updates have occurred for the particular storage address.

8   **[0074]**           EC38) The system of EC35, wherein the untransformed data comprises one or  
9   more portions of write data from the host.

11   **[0075]**           EC39) The system of EC38, wherein the portions of write data are determined  
12   based at least in part on a mirroring operation.

14   **[0076]**           EC40) The system of EC35, wherein the transformed data comprises one or  
15   more of

16       parity data based at least in part on write data from the host,  
17       exclusive-or (XOR) data based at least in part on write data from the host, and  
18       redundancy information based at least in part on a Redundant Array of  
19       Inexpensive/Independent Disks (RAID) implementation.

21   **[0077]**           EC41) The system of EC28, wherein the redundancy traffic comprises Read-  
22   Modify-Write (RMW) operation control and/or data redundancy traffic.

24   **[0078]**           EC42) The system of EC41, wherein the RMW operation data redundancy  
25   traffic is a copy of all or any portions of write data from the host.

27   **[0079]**           EC43) The system of EC41, wherein the RMW operation data redundancy  
28   traffic is transformed data based at least in part on a portion of write data from the host.

30   **[0080]**           EC44) The system of EC41, wherein the RMW operation control redundancy  
31   traffic is an indication of how many redundancy updates are to occur for a particular address.

33   **[0081]**           EC45) The system of EC28, wherein the redundancy traffic comprises data  
34   recovery operation control and/or data redundancy traffic.

1   **[0082]**           EC46) The system of EC28, wherein the redundant data storage is in  
2   accordance with one or more Redundant Array of Inexpensive/Independent Disks (RAID)  
3   techniques.

4  
5   **[0083]**           EC47) The system of EC28, wherein the primary agent is enabled to accept a  
6   request from the host to access storage, to forward the request as one or more sub-requests to all  
7   or any portions of the secondary agents, and to accept sub-statuses associated with the sub-  
8   requests from the secondary agents that the sub-requests were forwarded to.

9  
10   **[0084]**           EC48) The system of EC47, wherein the sub-requests are determined at least in  
11   part based on a redundancy technique that the redundant data storage is in accordance with.

12  
13   **[0085]**           EC49) The system of EC48, wherein the redundancy technique is according to  
14   a Redundant Array of Inexpensive/Independent Disks (RAID) implementation.

15  
16   **[0086]**           EC50) The system of EC47, wherein at least a portion of the sub-requests are  
17   based in part on a striping across the physical devices the sub-requests are forwarded to.

18  
19   **[0087]**           EC51) The system of EC50, wherein the striping is according to a Redundant  
20   Array of Inexpensive/Independent Disks (RAID) implementation.

21  
22   **[0088]**           EC52) The system of EC50, wherein the striping is changed dynamically  
23   according to workload.

24  
25   **[0089]**           EC53) The system of EC28, wherein each of the secondary agents is enabled to  
26   accept one or more sub-requests from the primary agent, translate host-context addressing  
27   information of the accepted sub-requests to access local storage and to provide respective sub-  
28   statuses to the primary agent based at least in part on the local accessing.

29  
30   **[0090]**           EC54) The system of EC28, wherein each of the secondary agents is enabled to  
31   accept one or more sub-requests from the primary agent, translate host-context addressing  
32   information of the accepted sub-requests to access local storage, access the local storage, and  
33   transfer data with the host, the data being read/write data of the access.

1     **[0091]**           EC55) The system of EC28, further comprising means for communicating host  
2     traffic between the host and the primary agent.

3

4     **[0092]**           EC56) The system of EC55, wherein the means for communicating peer-to-  
5     peer traffic is via at least one logical channel that is distinct from logical channels of the means  
6     for communicating host traffic.

7

8     **[0093]**           EC57) The system of EC55, wherein the means for communicating peer-to-  
9     peer traffic is via at least one physical channel that is distinct from physical channels of the  
10    means for communicating host traffic.

11

12    **[0094]**           EC58) The system of EC55, wherein one or more of the means for  
13    communicating peer-to-peer traffic and the means for communicating host traffic are compatible  
14    with a Peripheral Component Interconnect express (PCIe) standard.

15

16    **[0095]**           EC59) The system of EC55, wherein one or more of the means for  
17    communicating peer-to-peer traffic and the means for communicating host traffic are compatible  
18    with an Infiniband standard.

19

20    **[0096]**           EC60) The system of EC55, wherein the means for communicating peer-to-  
21    peer traffic comprises all or any portions of a non-transparent switch.

22

23    **[0097]**           EC61) The system of EC55, wherein the means for communicating host traffic  
24    comprises all or any portions of a transparent switch.

25



1     **[0098]**           EC62) A method comprising:  
2           enabling a host to access, as a single logical interface, an aggregation of a plurality of  
3           storage ranges distributed over a plurality of physical storage devices, one of the  
4           physical storage devices enabled to operate as a primary agent and one or more  
5           of the physical storage devices enabled to operate as secondary agents;  
6           communicating peer-to-peer traffic between the primary agent and at least one of the  
7           secondary agents;  
8           wherein the primary agent is enabled to manage at least a portion of the secondary  
9           agents for redundant data storage; and  
10          wherein the peer-to-peer traffic comprises redundancy traffic to implement the  
11          redundant data storage.

12  
13    **[0099]**           EC63) A system comprising:  
14          host interface logic hardware;  
15          a plurality of physical storage devices enabled to selectively communicate with a host at  
16          least in part via the host interface logic hardware, at least one of the physical  
17          storage devices enabled to operate as a primary agent and one or more of the  
18          physical storage devices enabled to operate as respective secondary agents,  
19          peer-to-peer communication logic hardware enabling peer-to-peer traffic between the  
20          primary agent and at least one of the secondary agents;  
21          wherein the physical storage devices implement an aggregation of a plurality of storage  
22          ranges distributed over the physical storage devices;  
23          wherein the primary agent is enabled to manage at least a portion of the secondary  
24          agents to implement redundant data storage; and  
25          wherein the peer-to-peer traffic comprises redundancy traffic to implement the  
26          redundant data storage.

27  
28    **[0100]**           EC64) A system comprising:  
29          means for interfacing between a plurality of storage agents via a peer-to-peer protocol;  
30          means for detecting a failure of a failed one of the storage agents, the failure being the  
31          failed storage agent no longer implementing a particular portion of storage; and  
32          means for allocating storage subject to reallocation to implement the particular portion  
33          of storage, the storage subject to reallocation comprising any combination of  
34          storage implemented by non-failed ones of the storage agents.

35

- 1   **[0101]**           EC65) The system of EC64, wherein the storage agents comprise secondary  
2   agents.  
3
- 4   **[0102]**           EC66) The system of EC64, wherein the storage agents comprise primary  
5   agents.  
6
- 7   **[0103]**           EC67) The system of EC64, wherein the means are comprised in a primary  
8   agent, and the storage agents comprise secondary agents.  
9
- 10   **[0104]**           EC68) The system of EC64, wherein each of the storage agents implements at  
11   least a respective portion of physical storage.  
12
- 13   **[0105]**           EC69) The system of EC68, wherein the respective portions comprise one or  
14   more non-volatile memories.  
15
- 16   **[0106]**           EC70) The system of EC64, further comprising local storage and wherein the  
17   storage subject to allocation further comprises the local storage.  
18
- 19   **[0107]**           EC71) The system of EC64, wherein the failure is a partial failure, and the  
20   failed storage agent continues to implement at least some storage after the partial failure.  
21
- 22   **[0108]**           EC72) The system of EC71, wherein the storage subject to allocation further  
23   comprises the at least some storage.  
24
- 25   **[0109]**           EC73) The system of EC64, wherein the means for detecting is implemented at  
26   least in part via one of the storage agents.  
27
- 28   **[0110]**           EC74) The system of EC64, wherein the means for detecting is implemented at  
29   least in part via a device driver.  
30
- 31   **[0111]**           EC75) The system of EC64, further comprising means for recovering data via  
32   one or more redundancy techniques.  
33
- 34   **[0112]**           EC76) The system of EC75, wherein the redundancy techniques comprise a  
35   mirroring technique.  
36

1   **[0113]**           EC77) The system of EC75, wherein the redundancy techniques comprise a  
2   Redundant Array of Inexpensive/Independent Disks (RAID) technique.

3  
4   **[0114]**           EC78) A method comprising:  
5           interfacing between a plurality of storage agents via a peer-to-peer protocol;  
6           detecting a failure of a failed one of the storage agents, the failure being the failed  
7           storage agent no longer implementing a particular portion of storage; and  
8           allocating storage subject to reallocation to implement the particular portion of storage,  
9           the storage subject to reallocation comprising any combination of storage  
10          implemented by non-failed ones of the storage agents.

11  
12   **[0115]**           EC79) A system comprising:  
13          peer-to-peer communication logic hardware enabling peer-to-peer protocol  
14          communication between a plurality of storage agents;  
15          failure detection logic hardware enabled to detect a failure of a failed one of the storage  
16          agents, the failure being the failed storage agent no longer implementing a  
17          particular portion of storage; and  
18          a processor enabled to execute instructions that when executed by the processor cause  
19          the processor to perform operations comprising allocating storage subject to  
20          reallocation to implement the particular portion of storage, the storage subject to  
21          reallocation comprising any combination of storage implemented by non-failed  
22          ones of the storage agents.

23  
24   **[0116]**           EC80) A computer readable medium having a set of instructions stored therein  
25   that when executed by a processing element cause the processing element to perform operations  
26   comprising:

27          managing interfacing between a plurality of storage agents via a peer-to-peer protocol;  
28          managing detecting a failure of a failed one of the storage agents, the failure being the  
29          failed storage agent no longer implementing a particular portion of storage; and  
30          managing allocating storage subject to reallocation to implement the particular portion  
31          of storage, the storage subject to reallocation comprising any combination of  
32          storage implemented by non-failed ones of the storage agents.

- 1     **[0117]**           EC81) A system comprising:  
2           means for detecting a failure of a failed one of a plurality of storage agents, the failure  
3           being the failed storage agent no longer enabling a host to access, as a single  
4           logical interface, an aggregation of a plurality of storage ranges distributed over  
5           the storage agents;  
6           means for identifying a replacement storage agent to replace the failed storage agent;  
7           and  
8           means for configuring the replacement storage agent to provide the single logical  
9           interface.  
10
- 11   **[0118]**           EC82) The system of EC81, wherein the replacement storage agent is a hot  
12   spare.  
13
- 14   **[0119]**           EC83) The system of EC81, wherein the replacement storage agent is one of  
15   the storage agents other than the failed storage agent.  
16
- 17   **[0120]**           EC84) The system of EC83, wherein the replacement storage agent is a primary  
18   agent.  
19
- 20   **[0121]**           EC85) The system of EC83, wherein the replacement storage agent is a  
21   configurable agent that is selectively operable as a primary agent or a secondary agent.  
22
- 23   **[0122]**           EC86) The system of EC81, wherein the means for detecting is implemented at  
24   least in part via one of the storage agents.  
25
- 26   **[0123]**           EC87) The system of EC81, wherein the means for detecting is implemented at  
27   least in part via a device driver executing on the host.  
28
- 29   **[0124]**           EC88) The system of EC81, further comprising means for interfacing to the  
30   host via a host interface protocol that comprises reading a storage request from an entry in a  
31   submission queue, the means for interfacing being in accordance with the providing the single  
32   logical interface.  
33
- 34   **[0125]**           EC89) The system of EC88, further comprising means for determining one or  
35   more storage sub-requests based at least in part on address information of the storage request.  
36

1     **[0126]**           EC90) The system of EC89, further comprising means for forwarding the sub-  
2     requests to secondary ones of the storage agents.

3

4     **[0127]**           EC91) The system of EC90, further comprising means for receiving sub-status  
5     information corresponding respectively to the sub-requests and returning overall status to the  
6     host based at least in part on the received sub-status information.

7

8     **[0128]**           EC92) The system of EC81, further comprising means for communicating peer-  
9     to-peer traffic among the storage agents.

10

11    **[0129]**           EC93) A method comprising:  
12           detecting a failure of a failed one of a plurality of storage agents, the failure being the  
13           failed storage agent no longer enabling a host to access, as a single logical  
14           interface, an aggregation of a plurality of storage ranges distributed over the  
15           storage agents;  
16           identifying a replacement storage agent to replace the failed storage agent; and  
17           configuring the replacement storage agent to provide the single logical interface.

18

19    **[0130]**           EC94) A system comprising:  
20           failure detection logic hardware enabled to detect a failure of a failed one of a plurality  
21           of storage agents, the failure being the failed storage agent no longer enabling a  
22           host to access, as a single logical interface, an aggregation of a plurality of  
23           storage ranges distributed over the storage agents; and  
24           a processor enabled to execute instructions that when executed by the processor cause  
25           the processor to perform operations comprising  
26                    identifying a replacement storage agent to replace the failed storage  
27                    agent, and  
28                    configuring the replacement storage agent to provide the single logical  
29                    interface.

30

1   **[0131]**           EC95) A computer readable medium having a set of instructions stored therein  
2   that when executed by a processing element cause the processing element to perform operations  
3   comprising:

4           managing detecting a failure of a failed one of a plurality of storage agents, the failure  
5                   being the failed storage agent no longer enabling a host to access, as a single  
6                   logical interface, an aggregation of a plurality of storage ranges distributed over  
7                   the storage agents;

8           managing identifying a replacement storage agent to replace the failed storage agent;

9                   and

10          managing configuring the replacement storage agent to provide the single logical  
11          interface.

12  
13   **[0132]**           EC96) A system comprising:

14          storage interface means for interfacing to mass storage;

15          host interface means for interfacing to a host via a host-interface protocol, the host

16                  interface means comprising means for reading a storage request from an entry in  
17                  a submission queue;

18          primary agent interface means for interfacing to a primary agent via a peer-to-peer  
19          protocol;

20          address determination means for determining at least one address to access the mass

21                  storage via the storage interface means, the means for determining operable

22                  based at least in part on address information comprised in the storage request;

23          storage capacity reporting means for reporting back to the host a zero storage capacity,

24                  and further for reporting back to the primary agent a storage capacity in

25                  accordance with a storage capacity of the mass storage; and

26          wherein the host interface means and the primary agent interface means are

27                  implemented at least in part via a same physical channel.

28  
29   **[0133]**           EC97) The system of EC96, further comprising destination determination  
30   means for determining, based at least in part on the address information, at least one destination  
31   to forward accumulated parity data to.

32  
33   **[0134]**           EC98) The system of EC97, wherein the at least one destination is reachable  
34   via the peer-to-peer protocol.

- 1     **[0135]**           EC99) The system of EC98, wherein the at least one destination is the primary  
2     agent.  
3
- 4     **[0136]**           EC100) The system of EC98, wherein the at least one destination is a secondary  
5     agent.  
6
- 7     **[0137]**           EC101) The system of EC96, wherein the address determination means  
8     operates in accordance with one or more redundancy techniques.  
9
- 10    **[0138]**           EC102) The system of EC101, wherein the redundancy techniques comprise  
11    one or more mirroring techniques and/or one or more Redundant Array of  
12    Inexpensive/Independent Disks (RAID) techniques.  
13
- 14    **[0139]**           EC103) The system of EC96, wherein the host interface means further  
15    comprises means for transferring data between the host and the mass storage.  
16
- 17    **[0140]**           EC104) The system of EC96, wherein the same physical channel is compatible  
18    with a Peripheral Component Interconnect express (PCIe) standard.  
19
- 20    **[0141]**           EC105) The system of EC96, wherein the same physical channel is compatible  
21    with an Infiniband standard.  
22
- 23    **[0142]**           EC106) The system of EC96, wherein the mass storage comprises one or more  
24    non-volatile memories.  
25

- 1     **[0143]**           EC107) A method comprising:  
2           interfacing to mass storage;  
3           interfacing to a host via a host-interface protocol, the interfacing to the host comprising  
4                 reading a storage request from an entry in a submission queue;  
5           interfacing to a primary agent via a peer-to-peer protocol;  
6           determining at least one address to access the mass storage via the interfacing to mass  
7                 storage, the determining operable based at least in part on address information  
8                 comprised in the storage request;  
9           reporting back to the host a zero storage capacity, and reporting back to the primary  
10                 agent a storage capacity in accordance with a storage capacity of the mass  
11                 storage; and  
12           wherein the interfacing to the host and the interfacing to the primary agent are  
13                 implemented at least in part via a same physical channel.  
14
- 15    **[0144]**           EC108) A system comprising:  
16           a mass storage hardware interface enabled to interface with mass storage;  
17           host interface logic hardware enabling interfacing with the host, compatible with a host-  
18                 interface protocol, and enabled to read a storage request from an entry in a  
19                 submission queue;  
20           peer-to-peer communication logic hardware enabling interfacing with a primary agent  
21                 via a peer-to-peer protocol;  
22           access address determination logic hardware enabled to determine at least one address to  
23                 access the mass storage via the mass storage hardware interface, the determining  
24                 operable based at least in part on address information comprised in the storage  
25                 request;  
26           storage capacity reporting logic hardware enabled to report back to the host a zero  
27                 storage capacity, and further enabled to report back to the primary agent a  
28                 storage capacity in accordance with a storage capacity of the mass storage; and  
29           wherein the interfacing with the host and the interfacing with the primary agent are  
30                 implemented at least in part via a same physical channel.  
31



1   **[0145]**           EC109) A computer readable medium having a set of instructions stored therein  
2   that when executed by a processing element cause the processing element to perform operations  
3   comprising:

4           managing interfacing to mass storage;  
5           managing interfacing to a host via a host-interface protocol, the interfacing to the host  
6                comprising reading a storage request from an entry in a submission queue;  
7           managing interfacing to a primary agent via a peer-to-peer protocol;  
8           managing determining at least one address to access the mass storage via the interfacing  
9                to mass storage, the determining operable based at least in part on address  
10           information comprised in the storage request;  
11          managing reporting back to the host a zero storage capacity, and reporting back to the  
12           primary agent a storage capacity in accordance with a storage capacity of the  
13           mass storage; and  
14          wherein the interfacing to the host and the interfacing to the primary agent are  
15           implemented at least in part via a same physical channel.

16  
17   **[0146]**           EC110) A system comprising:  
18          storage interface means for interfacing to mass storage;  
19          host interface means for interfacing to a host via a host-interface protocol, the host  
20           interface means comprising means for transferring data between the host and the  
21           mass storage;  
22          primary agent interface means for interfacing to a primary agent via a peer-to-peer  
23           protocol, the primary agent interface means comprising means for receiving a  
24           forwarded storage request from the primary agent;  
25          address determination means for determining at least one address to access the mass  
26           storage, the means for determining operable based at least in part on address  
27           information comprised in the forwarded storage request;  
28          wherein the forwarded storage request corresponds to a request obtained from a  
29           submission queue, the request and the forwarded request referring to identical  
30           one or more address ranges; and  
31          wherein the host interface means and the primary agent interface means are  
32           implemented at least in part via a same physical channel.

33

1     **[0147]**           EC111) The system of EC110, further comprising destination determination  
2     means for determining, based at least in part on the address information, at least one destination  
3     to forward accumulated parity data to.

4  
5     **[0148]**           EC112) The system of EC111, wherein the at least one destination is reachable  
6     via the peer-to-peer protocol.

7  
8     **[0149]**           EC113) The system of EC112, wherein the at least one destination is the  
9     primary agent.

10  
11    **[0150]**           EC114) The system of EC112, wherein the at least one destination is a  
12    secondary agent.

13  
14    **[0151]**           EC115) The system of EC110, wherein the address determination means  
15    operates in accordance with one or more redundancy techniques.

16  
17    **[0152]**           EC116) The system of EC115, wherein the redundancy techniques comprise  
18    one or more mirroring techniques and/or one or more Redundant Array of  
19    Inexpensive/Independent Disks (RAID) techniques.

20  
21    **[0153]**           EC117) The system of EC110, further comprising storage capacity reporting  
22    means for reporting back to the host a zero storage capacity, and further for reporting back to the  
23    primary agent a storage capacity in accordance with a storage capacity of the mass storage.

24  
25    **[0154]**           EC118) The system of EC110, wherein the same physical channel is  
26    compatible with a Peripheral Component Interconnect express (PCIe) standard.

27  
28    **[0155]**           EC119) The system of EC110, wherein the same physical channel is  
29    compatible with an Infiniband standard.

30  
31    **[0156]**           EC120) The system of EC110, wherein the mass storage comprises one or more  
32    non-volatile memories.

- 1     **[0157]**           EC121) A method comprising:  
2           interfacing to mass storage;  
3           interfacing to a host via a host-interface protocol, the interfacing to the host comprising  
4                transferring data between the host and the mass storage;  
5           interfacing to a primary agent via a peer-to-peer protocol, the interfacing to the primary  
6                agent comprising receiving a forwarded storage request from the primary agent;  
7           determining at least one address to access the mass storage via the interfacing to mass  
8                storage, the determining operable based at least in part on address information  
9                comprised in the storage request;  
10          wherein the forwarded storage request corresponds to a request obtained from a  
11                submission queue, the request and the forwarded request referring to identical  
12                one or more address ranges; and  
13          wherein the interfacing to the host and the interfacing to the primary agent are  
14                implemented at least in part via a same physical channel.  
15
- 16     **[0158]**           EC122) A system comprising:  
17           a mass storage hardware interface enabled to interface with mass storage;  
18           host interface logic hardware enabling interfacing with the host, compatible with a host-  
19                interface protocol, and enabled to transfer data between the host and the mass  
20                storage at least in part via the mass storage hardware interface;  
21           peer-to-peer communication logic hardware enabling interfacing with a primary agent  
22                via a peer-to-peer protocol;  
23           storage request accepting logic hardware enabling receiving a forwarded storage request  
24                from the primary agent at least in part via the peer-to-peer communication logic  
25                hardware;  
26           access address determination logic hardware enabled to determine at least one address to  
27                access the mass storage via the mass storage hardware interface, the determining  
28                operable based at least in part on address information comprised in the storage  
29                request;  
30          wherein the forwarded storage request corresponds to a request obtained from a  
31                submission queue, the request and the forwarded request referring to identical  
32                one or more address ranges; and  
33          wherein the interfacing with the host and the interfacing with the primary agent are  
34                implemented at least in part via a same physical channel.  
35

1   **[0159]**           EC123) A computer readable medium having a set of instructions stored therein  
2   that when executed by a processing element cause the processing element to perform operations  
3   comprising:

4           managing interfacing to mass storage;  
5           managing interfacing to a host via a host-interface protocol, the interfacing to the host  
6           comprising transferring data between the host and the mass storage;  
7           managing interfacing to a primary agent via a peer-to-peer protocol, the interfacing to  
8           the primary agent comprising receiving a forwarded storage request from the  
9           primary agent;  
10          managing determining at least one address to access the mass storage via the interfacing  
11          to mass storage, the determining operable based at least in part on address  
12          information comprised in the storage request;  
13          wherein the forwarded storage request corresponds to a request obtained from a  
14          submission queue, the request and the forwarded request referring to identical  
15          one or more address ranges; and  
16          wherein the interfacing to the host and the interfacing to the primary agent are  
17          implemented at least in part via a same physical channel.

18  
19   **[0160]**           EC124) A system comprising:

20          host interface means for interfacing to a host via a host-interface protocol, the host  
21          interface means comprising means for reading a storage request from an entry in  
22          a submission queue;  
23          secondary agent interface means for interfacing to a secondary agent via a peer-to-peer  
24          protocol, the secondary agent interface means comprising means for forwarding  
25          a storage sub-request to the secondary agent;  
26          sub-request generation means for determining the sub-request based at least in part on  
27          address information comprised in the storage request;  
28          storage capacity reporting means for reporting back to the host a total storage capacity,  
29          the total storage capacity based at least in part on a storage capacity of the  
30          secondary agent; and  
31          wherein the host interface means and the secondary agent interface means are  
32          implemented at least in part via a same physical channel.

33  
34   **[0161]**           EC125) The system of EC124, further comprising storage interface means for  
35   interfacing to mass storage.  
36

1     **[0162]**           EC126) The system of EC125, wherein the mass storage comprises one or more  
2     non-volatile memories.

3

4     **[0163]**           EC127) A method comprising:  
5           interfacing with a host via a host-interface protocol, the host interfacing comprising  
6           reading a storage request from an entry in a submission queue;  
7           interfacing with a secondary agent via a peer-to-peer protocol, the secondary agent  
8           interfacing comprising forwarding a storage sub-request to the secondary agent;  
9           determining the sub-request based at least in part on address information comprised in  
10          the storage request;  
11          reporting back to the host a total storage capacity, the total storage capacity based at  
12          least in part on a storage capacity of the secondary agent; and  
13          wherein the interfacing with the host and the interfacing with the secondary agent are  
14          implemented at least in part via a same physical channel.

15

16    **[0164]**           EC128) A system comprising:  
17          host interface logic hardware enabling interfacing with the host, compatible with a host-  
18          interface protocol, and enabled to read a storage request from an entry in a  
19          submission queue;  
20          peer-to-peer communication logic hardware enabling interfacing with a secondary agent  
21          via a peer-to-peer protocol;  
22          request forwarding logic hardware enabled to forward a storage sub-request to the  
23          secondary agent at least in part via the peer-to-peer communication logic  
24          hardware;  
25          sub-request determining logic hardware enabled to determine the sub-request based at  
26          least in part on address information comprised in the storage request;  
27          storage capacity reporting logic hardware enabled to report back to the host a total  
28          storage capacity, the total storage capacity based at least in part on a storage  
29          capacity of the secondary agent; and  
30          wherein the interfacing with the host and the interfacing with the secondary agent are  
31          implemented at least in part via a same physical channel.

32

1   **[0165]**           EC129) A computer readable medium having a set of instructions stored therein  
2   that when executed by a processing element cause the processing element to perform operations  
3   comprising:

4           managing interfacing with a host via a host-interface protocol, the host interfacing  
5                comprising reading a storage request from an entry in a submission queue;  
6           managing interfacing with a secondary agent via a peer-to-peer protocol, the secondary  
7                agent interfacing comprising forwarding a storage sub-request to the secondary  
8                agent;  
9           managing determining the sub-request based at least in part on address information  
10                comprised in the storage request;  
11           managing reporting back to the host a total storage capacity, the total storage capacity  
12                based at least in part on a storage capacity of the secondary agent; and  
13           wherein the interfacing with the host and the interfacing with the secondary agent are  
14                implemented at least in part via a same physical channel.

15  
16   **[0166]**           EC130) A system comprising:  
17           storage interface means for interfacing to mass storage;  
18           host interface means for interfacing to a host via a host-interface protocol, the host  
19                interface means comprising means for transferring data between the host and the  
20                mass storage;  
21           primary agent interface means for interfacing to a primary agent via a peer-to-peer  
22                protocol, the primary agent interface means comprising means for receiving a  
23                forwarded storage sub-request from the primary agent;  
24           storage capacity reporting means for reporting back to the host a zero storage capacity,  
25                and further for reporting back to the primary agent a storage capacity in  
26                accordance with a storage capacity of the mass storage; and  
27           wherein the host interface means and the primary agent interface means are  
28                implemented at least in part via a same physical channel.

29  
30   **[0167]**           EC131) The system of EC130, wherein the mass storage comprises one or more  
31   non-volatile memories.  
32

1     **[0168]**           EC132) A method comprising:  
2           interfacing to mass storage;  
3           interfacing with a host via a host-interface protocol, the host interfacing comprising  
4           transferring data between the host and the mass storage;  
5           interfacing with a primary agent via a peer-to-peer protocol, the primary agent  
6           interfacing comprising receiving a forwarded storage sub-request from the  
7           primary agent;  
8           reporting back to the host a zero storage capacity, and further reporting back to the  
9           primary agent a storage capacity in accordance with a storage capacity of the  
10          mass storage; and  
11          wherein the interfacing with the host and the interfacing with the primary agent are  
12          implemented at least in part via a same physical channel.

13  
14    **[0169]**           EC133) A system comprising:  
15          a mass storage hardware interface enabled to interface with mass storage;  
16          host interface logic hardware enabling interfacing with the host, compatible with a host-  
17          interface protocol, and enabled to transfer data between the host and the mass  
18          storage at least in part via the mass storage hardware interface;  
19          peer-to-peer communication logic hardware enabling interfacing with a primary agent  
20          via a peer-to-peer protocol;  
21          sub-request receiving logic hardware enabled to receive a forwarded storage sub-request  
22          from the primary agent at least in part via the peer-to-peer communication logic  
23          hardware;  
24          storage capacity reporting logic hardware enabled to report back to the host a zero  
25          storage capacity, and further enabled to report back to the primary agent a  
26          storage capacity in accordance with a storage capacity of the mass storage; and  
27          wherein the interfacing with the host and the interfacing with the primary agent are  
28          implemented at least in part via a same physical channel.

29

1   **[0170]**           EC134) A computer readable medium having a set of instructions stored therein  
2   that when executed by a processing element cause the processing element to perform operations  
3   comprising:

4           managing interfacing to mass storage;  
5           managing interfacing with a host via a host-interface protocol, the host interfacing  
6           comprising transferring data between the host and the mass storage;  
7           managing interfacing with a primary agent via a peer-to-peer protocol, the primary agent  
8           interfacing comprising receiving a forwarded storage sub-request from the  
9           primary agent;  
10          managing reporting back to the host a zero storage capacity, and further reporting back  
11          to the primary agent a storage capacity in accordance with a storage capacity of  
12          the mass storage; and  
13          wherein the interfacing with the host and the interfacing with the primary agent are  
14          implemented at least in part via a same physical channel.

15  
16   **[0171]**           EC135) A system comprising:  
17          storage interface means for interfacing to mass storage;  
18          host interface means for interfacing to a host via a host-interface protocol, the host  
19          interface means comprising means for transferring data between the host and the  
20          mass storage;  
21          primary agent interface means for interfacing to a primary agent via a peer-to-peer  
22          protocol, the primary agent interface means comprising means for receiving a  
23          sub-request from the primary agent and means for sending a sub-status to the  
24          primary agent, the sub-status based at least in part on results of accessing the  
25          mass storage in accordance with the sub-request;  
26          storage capacity reporting means for reporting back to the host a zero storage capacity,  
27          and further for reporting back to the primary agent a storage capacity in  
28          accordance with a storage capacity of the mass storage; and  
29          wherein the host interface means and the primary agent interface means are  
30          implemented at least in part via a same physical channel.

31  
32   **[0172]**           EC136) The system of EC135, wherein the primary agent interface means  
33   further comprises means for communicating redundancy information with the primary agent.  
34



1     **[0173]**           EC137) The system of EC136, wherein the redundancy information is  
2     compatible with one or more mirroring techniques and/or one or more Redundant Array of  
3     Inexpensive/Independent Disks (RAID) techniques.

4  
5     **[0174]**           EC138) The system of EC135, wherein the same physical channel is  
6     compatible with a Peripheral Component Interconnect express (PCIe) standard.

7  
8     **[0175]**           EC139) The system of EC135, wherein the same physical channel is  
9     compatible with an Infiniband standard.

10  
11    **[0176]**           EC140) The system of EC135, wherein the mass storage comprises one or more  
12    non-volatile memories.

13  
14    **[0177]**           EC141) A method comprising:  
15           interfacing to mass storage;  
16           interfacing with a host via a host-interface protocol, the interfacing with the host  
17           comprising transferring data between the host and the mass storage;  
18           interfacing with a primary agent via a peer-to-peer protocol, the interfacing with the  
19           primary agent comprising receiving a sub-request from the primary agent and  
20           sending a sub-status to the primary agent, the sub-status based at least in part on  
21           results of accessing the mass storage in accordance with the sub-request;  
22           reporting back to the host a zero storage capacity, and further reporting back to the  
23           primary agent a storage capacity in accordance with a storage capacity of the  
24           mass storage; and  
25           wherein the interfacing with the host and the interfacing with the primary agent are  
26           implemented at least in part via a same physical channel.

27

1     **[0178]**           EC142) A system comprising:  
2           a mass storage hardware interface enabled to interface with mass storage;  
3           host interface logic hardware enabling interfacing with the host, compatible with a host-  
4           interface protocol, and enabled to transfer data between the host and the mass  
5           storage at least in part via the mass storage hardware interface;  
6           peer-to-peer communication logic hardware enabling interfacing with a primary agent  
7           via a peer-to-peer protocol;  
8           sub-request receiving logic hardware enabled to receive a forwarded storage sub-request  
9           from the primary agent at least in part via the peer-to-peer communication logic  
10          hardware;  
11          sub-status determination logic hardware enabled to determine a sub-status, based at least  
12          in part on accessing the mass storage in accordance with the sub-request and via  
13          the mass storage hardware interface;  
14          sub-status forwarding logic hardware enabled to forward the sub-status to the primary  
15          agent at least in part via the peer-to-peer communication logic hardware;  
16          storage capacity reporting logic hardware enabled to report back to the host a zero  
17          storage capacity, and further enabled to report back to the primary agent a  
18          storage capacity in accordance with a storage capacity of the mass storage; and  
19          wherein the interfacing with the host and the interfacing with the primary agent are  
20          implemented at least in part via a same physical channel.  
21

1   **[0179]**           EC143) A computer readable medium having a set of instructions stored therein  
2   that when executed by a processing element cause the processing element to perform operations  
3   comprising:  
4       managing interfacing to mass storage;  
5       managing interfacing with a host via a host-interface protocol, the interfacing with the  
6       host comprising transferring data between the host and the mass storage;  
7       managing interfacing with a primary agent via a peer-to-peer protocol, the interfacing  
8       with the primary agent comprising receiving a sub-request from the primary  
9       agent and sending a sub-status to the primary agent, the sub-status based at least  
10      in part on results of accessing the mass storage in accordance with the sub-  
11      request;  
12      managing reporting back to the host a zero storage capacity, and further managing  
13      reporting back to the primary agent a storage capacity in accordance with a  
14      storage capacity of the mass storage; and  
15      wherein the interfacing with the host and the interfacing with the primary agent are  
16      implemented at least in part via a same physical channel.  
17

- 1     **[0180]**           EC144) A method comprising:  
2           accepting a request from a host to access storage at an address;  
3           determining, based at least in part on the address, which one or more of a plurality of  
4           storage devices the request corresponds to;  
5           determining, based at least in part on the address, a respective sub-request corresponding  
6           to each of the determined storage devices;  
7           sending the sub-requests to the corresponding determined storage devices;  
8           accepting respective sub-status from each of the determined storage devices;  
9           determining an overall status based at least in part on each of the respective sub-statuses;  
10          providing the overall status to the host;  
11          wherein the accepting the request, the determining the storage devices, the determining  
12           and the sending the sub-requests, the accepting the sub-statuses, and the  
13           determining and the providing the overall status are performed by a particular  
14           one of the storage devices;  
15          wherein the accepting the request is via a host-interface protocol and the sending the  
16           sub-requests and the accepting the sub-statuses are via a peer-to-peer protocol  
17           operable independently of the host; and  
18          wherein at least one of the storage devices implements storage that is accessible via one  
19           of the sub-requests and the at least one of the storage devices disables accesses  
20           via the host-interface protocol to the implemented storage.  
21
- 22     **[0181]**           EC145) The method of EC144, wherein the host-interface and the peer-to-peer  
23     protocols are compatible with a Peripheral Component Interconnect express (PCIe) channel.  
24
- 25     **[0182]**           EC146) The method of EC144, wherein the accepting the request comprises  
26     reading an entry in a submission queue retained in host memory.  
27
- 28     **[0183]**           EC147) The method of EC144, wherein the providing the overall status  
29     comprises writing an entry in a completion queue retained in host memory.  
30
- 31     **[0184]**           EC148) The method of EC144, further comprising independently performing,  
32     by each of the determined storage devices, a data transfer portion of the corresponding sub-  
33     request.  
34

1     **[0185]**           EC149) The method of EC148, wherein the independently performing  
2 comprises accessing host memory.

3  
4     **[0186]**           EC150) A method comprising:  
5           a first storage device reading a first request from a first one of a plurality of host request  
6           queues;  
7           a second storage device reading a second request from a second one of the host request  
8           queues, the second host request queue being distinct from the first host request  
9           queue, the reading by the second storage device being independent of the  
10          reading by the first storage device;  
11          the first storage device forwarding at least a portion of the first request to the second  
12          storage device;  
13          the second storage device forwarding at least a portion of the second request to the first  
14          storage device, the forwarding by the second storage device being independent  
15          of the forwarding by the first storage device;  
16          the first storage device accessing first storage in accordance with the at least the portion  
17          of the second request and returning a first status to the second storage device,  
18          based at least in part on the accessing by the first storage device; and  
19          the second storage device accessing second storage in accordance with the at least the  
20          portion of the first request and returning a second status to the first storage  
21          device, based at least in part on the accessing by the second storage device, the  
22          accessing by the second storage device being independent of the accessing by  
23          the first storage device, the second storage being distinct from the first storage.

24  
25     **[0187]**           EC151) The method of EC150, wherein the host request queues are accessible  
26 by a host, and further comprising the second storage device returning a third status to the host  
27 based at least in part on the second request and the first storage device returning a fourth status  
28 to the host based at least in part on the first request.

29  
30     **[0188]**           EC152) The method of EC150, wherein the at least the portion of the second  
31 request is at least a first portion, and further comprising the second storage device accessing  
32 third storage in accordance with at least a second portion of the second request, the first and the  
33 second portions of the second request being distinct portions of the second request.

34

1     **[0189]**         EC153) The method of EC150, wherein the host request queues are accessible  
2     by a host and the at least the portion of the second request is at least a first portion, and further  
3     comprising the second storage device accessing third storage in accordance with at least a  
4     second portion of the second request and the second storage device returning a third status to the  
5     host based at least in part on the accessing of the third storage, the first and the second portions  
6     of the second request being distinct portions of the second request.

7  
8     **[0190]**         EC154) A method comprising:  
9         providing to a host a single logical interface to a plurality of physical storage devices,  
10         one of the physical storage devices operating as a primary agent and one or  
11         more of the physical storage devices operating as secondary agents;  
12         wherein the providing comprises  
13         the primary agent accepting a request from the host to access storage, forwarding the  
14         request as one or more sub-requests to all or any portions of the secondary  
15         agents, accepting sub-statuses associated with the sub-requests from the  
16         secondary agents that the sub-requests were forwarded to, and providing an  
17         overall status based at least on the sub-statuses to the host, and  
18         the all or any portions of the secondary agents accepting the sub-requests, translating  
19         host-context addressing information of the sub-requests to local-context  
20         addressing information, accessing local storage based at least in part on the  
21         local-context addressing information, and providing the sub-statuses to the  
22         primary agent;  
23         wherein the primary agent, when queried by the host regarding storage capacity, reports  
24         the storage capacity based at least in part on storage implemented by the  
25         primary agent and storage implemented by the secondary agents; and  
26         wherein each of the secondary agents, when queried by the host regarding storage  
27         capacity, report a respective storage capacity based at least in part on the storage  
28         capacity reported by the primary agent.

29  
30     **[0191]**         EC155) The method of EC154, wherein the storage capacity reported by the  
31     primary agent is increased in accordance with at least a portion of the storage implemented by a  
32     particular one of the secondary agents; and further comprising the particular secondary agent  
33     reporting to the host a storage capacity of the particular secondary agent as decreased in  
34     accordance with the at least a portion of the storage implemented by the particular second agent.

35

1     **[0192]**           EC156) The method of EC154, wherein a particular one of the secondary  
2     agents, when queried by the host regarding storage capacity, reports back zero storage capacity  
3     if all of the storage implemented by the particular secondary agent is accounted for in the  
4     storage capacity reported by the primary agent.

5  
6     **[0193]**           EC157) The method of EC154, wherein the primary agent is enabled to  
7     communicate with the host via a host-interface protocol.

8  
9     **[0194]**           EC158) The method of EC154, wherein the primary agent and the secondary  
10    agents are enabled to communicate with each other via a peer-to-peer protocol.

11  
12    **[0195]**           EC159) The method of EC154, wherein the primary agent is enabled to  
13    communicate with the host and at least one of the secondary agents via a same physical link.

14  
15    **[0196]**           EC160) A system comprising:  
16                    host interface logic enabled to interface to a host via a host-interface protocol and to  
17                    interface to a secondary agent via a peer-to-peer protocol, the interfacing to the  
18                    host comprising reading a storage request from an entry in a submission queue,  
19                    the interfacing to the secondary agent comprising forwarding a storage sub-  
20                    request to the secondary agent, the interfacing to the host and the interfacing to  
21                    the secondary agent being via a same physical channel;  
22                    sub-request generation logic enabled to determine the sub-request based at least in part  
23                    on address information comprised in the storage request; and  
24                    storage capacity reporting logic enabled to report back to the host a total storage  
25                    capacity, the total storage capacity based at least in part on a storage capacity of  
26                    the secondary agent.

27  
28    **[0197]**           EC161) The system of EC160, wherein the physical channel is compatible with  
29    at least one of a Peripheral Component Interconnect express (PCIe) channel and an Infiniband  
30    channel.

31  
32    **[0198]**           EC162) The system of EC160, further comprising one or more of the host, the  
33    secondary agent, and a switch coupling the host to the host interface logic and/or the secondary  
34    agent.

1     **[0199]**           EC163) The system of EC160, further comprising storage interface logic and  
2     mass storage, the storage interface logic enabled to interface to the mass storage, and wherein  
3     the total storage capacity is further in accordance with a storage capacity of the mass storage.  
4

5     **[0200]**           EC164) The system of EC160, wherein the interfacing to the host further  
6     comprises writing an overall status to an entry in a completion queue, and the interfacing to the  
7     secondary agent further comprises forwarding a sub-status from the secondary agent to the host  
8     as at least part of the overall status via the writing of the overall status to the completion queue;  
9     and further comprising sub-status accumulation logic enabled to determine the overall status  
10    based at least in part on the sub-status.  
11

12    **[0201]**           EC165) The system of EC160, wherein system is comprised in a Solid-State  
13    Disk (SSD).  
14

15    **[0202]**           EC166) A system comprising:  
16           host interface logic enabled to interface to a host via a host-interface protocol and to  
17           interface to a secondary agent via a peer-to-peer protocol, the interfacing to the  
18           host comprising reading a storage request from an entry in a submission queue,  
19           the interfacing to the secondary agent comprising forwarding a storage sub-  
20           request to the secondary agent, the interfacing to the host and the interfacing to  
21           the secondary agent being via a same physical channel;  
22           sub-request generation logic enabled to determine the sub-request based at least in part  
23           on address information comprised in the storage request; and  
24           storage capacity reporting logic enabled to report back to the host a total storage  
25           capacity, the total storage capacity based at least in part on a storage capacity of  
26           the secondary agent.  
27

28    **[0203]**           EC167) The system of EC166, wherein the physical channel is compatible with  
29    at least one of a Peripheral Component Interconnect express (PCIe) channel and an Infiniband  
30    channel.  
31

32    **[0204]**           EC168) The system of EC166, further comprising one or more of the host, the  
33    primary agent, a switch coupling the host to the host interface logic and/or the primary agent,  
34    and all or any portions of the mass storage.  
35



1     **[0205]**           EC169) The system of EC166, wherein the mass storage comprises a plurality  
2     of flash storage devices.

3

4     **[0206]**           EC170) The system of EC166, wherein the system is comprised in a solid-state  
5     disk.

6

7     **[0207]**           EC171) A system comprising:

8           a storage subsystem having physical components comprising a plurality of physical

9           storage devices and a physical switch-portion having a plurality of ports;

10          a host coupled via a dedicated point-to-point link to a dedicated port of the plurality of  
11          ports;

12          wherein each physical storage device comprises:

13                  at least one range of storage,

14                  at least one port enabled to be coupled via a respective point-to-point link to a

15                  respective port of the plurality of ports, and

16                  agent logic enabling the physical storage device to operate as at least one agent

17                  of one or more primary agents and one or more secondary agents of the

18                  storage subsystem, each primary agent enabled to accept host-initiated

19                  storage access requests, to generate sub-requests, and to accumulate

20                  sub-statuses, each secondary agent enabled to accept at least one of the

21                  sub-requests and to generate at least one of the sub-statuses; and

22          wherein the storage subsystem is enabled to operate one or more logical storage devices,

23                  each logical storage device enabling the host to access via a single logical

24                  interface an aggregation of the storage ranges distributed over the storage

25                  devices corresponding to one primary agent of the primary agents and at least

26                  one secondary agent of the secondary agents.

27

28     **[0208]**           EC172) The system of EC171, further comprising:

29           wherein the agent logic comprises configurability logic enabling the physical storage

30           device to operate at least sometimes as at least one of the primary agents and at

31           least sometimes as at least one of the secondary agents.

32

1     **[0209]**           EC173) The system of EC171, further comprising:  
2                   wherein the agent logic comprises concurrency logic enabling the physical storage  
3                   device to operate concurrently as at least one of the primary agents and at least  
4                   one of the secondary agents.

5  
6     **[0210]**           EC174) The system of EC171, further comprising:  
7                   wherein the agent logic comprises dedicated-agent logic enabling the physical storage  
8                   device to operate as a dedicated one of the primary agents and the secondary  
9                   agents.

10  
11    **[0211]**           EC175) The system of EC171, further comprising:  
12                   wherein a first logical storage device of the logical storage devices comprises at least a  
13                   first primary agent of the primary agents and one or more of the secondary  
14                   agents; and  
15                   wherein the physical switch-portion enables host-occluded peer-to-peer communications  
16                   between the first primary agent and each of the secondary agents.

17  
18    **[0212]**           EC176) The system of EC171, further comprising:  
19                   wherein the physical switch-portion having a plurality of ports is a first physical switch-  
20                   portion having a first plurality of ports and the physical components further  
21                   comprise a second physical switch-portion having a second plurality of ports;  
22                   and  
23                   wherein each physical storage device further comprises at least one port enabled to be  
24                   coupled via a respective point-to-point link to a respective port of the second  
25                   plurality of ports.

26  
27    **[0213]**           EC177) The system of EC176, further comprising:  
28                   wherein the second physical switch-portion enables storage redundancy techniques  
29                   implemented via peer-to-peer communication with inconsequential impact on  
30                   bandwidth available for communications via the first physical switch-portion.

31  
32    **[0214]**           EC178) The system of EC177, further comprising:  
33                   wherein the peer-to-peer communication is via the second physical switch-portion and  
34                   comprises one or more of control information, untransformed redundancy data,  
35                   and transformed redundancy data.

36

- 1     **[0215]**           EC179) The system of EC176, further comprising:  
2           wherein a unitary physical switch comprises the first physical switch-portion and the  
3           second physical switch-portion.  
4
- 5     **[0216]**           EC180) The system of EC171, further comprising:  
6           wherein the single logical interface enables one or more of mirroring, striping, RAID  
7           parity, and fail-over, to be implemented via host-occluded peer-to-peer  
8           communications between the primary and secondary agents.  
9
- 10    **[0217]**           EC181) The system of EC171, further comprising:  
11          wherein the primary agents comprise a first primary agent and a second primary agent,  
12          the secondary agents comprise a first secondary agent and a second secondary  
13          agent, the physical storage devices comprise a first physical storage device and a  
14          second physical storage device, and the logical storage devices comprise a first  
15          logical storage device and a second logical storage device;  
16          wherein with respect to a first subset of the requests, the first physical storage device  
17          operates as the first primary agent and as the first secondary agent;  
18          wherein with respect to a second subset of the requests, the second physical storage  
19          device operates as the second primary agent and as the second secondary agent;  
20          and  
21          wherein the first logical storage device comprises the first primary agent and the second  
22          secondary agent, the second logical storage device comprises the second  
23          primary agent and the first secondary agent, and the first logical storage device  
24          operates concurrently with the second logical storage device.  
25
- 26    **[0218]**           EC182) The system of EC181, further comprising:  
27          wherein the first subset of the requests are from a first host request queue and the second  
28          subset of the requests are from a second host request queue.  
29
- 30    **[0219]**           EC183) The system of EC181, further comprising:  
31          wherein the host is a first host and the first subset of the requests are from the first host,  
32          and the second subset of the requests are from a second host.  
33

1     **[0220]**           EC184) The system of EC181, further comprising:  
2                   wherein the first subset of the requests and the second subset of the requests are from the  
3                   same host request queue.  
4

5     **[0221]**           EC185) The system of EC171, further comprising:  
6                   wherein a first logical storage device of the logical storage devices comprises at least a  
7                   first primary agent of the primary agents and at least a first secondary agent of  
8                   the secondary agents;  
9                   wherein the aggregation of the storage ranges of the agents of the first logical storage  
10                  device is a first logical storage sub-space; and  
11                  wherein the agent logic comprises storage-recovery logic enabling the first primary  
12                  agent to identify free storage available to be provided by one or more of the  
13                  agents and, subsequent to a determination that a particular portion of the first  
14                  logical storage sub-space is no longer being provided by one of the agents, to  
15                  allocate from the free storage to restore the previously provided particular  
16                  portion.  
17

18    **[0222]**           EC186) The system of EC171, further comprising:  
19                   wherein a first logical storage device of the logical storage devices comprises at least a  
20                   first primary agent of the primary agents and at least a first secondary agent of  
21                   the secondary agents, the physical storage devices comprise a first physical  
22                   storage device and a second physical storage device, and wherein the first  
23                   physical device initially operates as the first primary agent;  
24                   a monitoring agent enabled to determine that the first primary agent has failed and to  
25                   transmit a primary-agent-replacement request to a selected one of the other  
26                   agents; and  
27                   wherein the agent logic of the second physical device comprises primary-agent-  
28                   replacement logic enabling the second physical storage device to operate as a  
29                   replacement for the first primary agent in response to receiving the primary-  
30                   agent-replacement request from the monitoring agent.  
31

- 1     **[0223]**           EC187) A method comprising:  
2           operating as a primary agent a first storage device and operating as one or more  
3           secondary agents at least a second storage device, each storage device of the  
4           primary and secondary agents having at least one respective storage range;  
5           the primary agent providing a logical interface enabling the host to access as a unitary  
6           logical device an aggregation of the storage ranges distributed over the storage  
7           devices of the primary and secondary agents, primary agent communications  
8           with the host via the logical interface comprising storage access requests  
9           accepted by the primary agent from the host and overall status sent by the  
10          primary agent to the host;  
11          the primary agent forwarding each request accepted as one or more sub-requests to all or  
12          any portions of the secondary agents, accepting sub-statuses associated with the  
13          sub-requests from the secondary agents that the sub-requests were forwarded to,  
14          and formulating each overall status sent based at least on the sub-statuses; and  
15          the all or any portions of the secondary agents accepting the sub-requests, translating  
16          host-context addressing information of the sub-requests to local-context  
17          addressing information, accessing local storage based at least in part on the  
18          local-context addressing information, and providing the sub-statuses to the  
19          primary agent.  
20
- 21     **[0224]**           EC188) The method of EC187, further comprising:  
22          the primary agent determining an aggregated storage capacity of the unitary logical  
23          device based at least in part on the storage ranges of the primary and secondary  
24          agents;  
25          in response to the host querying the unitary logical device regarding storage capacity,  
26          the primary agent reporting the aggregated storage capacity; and  
27          in response to the host querying the secondary agents regarding storage capacity, each of  
28          the secondary agents reporting a respective storage capacity that excludes  
29          representation of any storage range of the secondary agent that is represented in  
30          the aggregated storage capacity.  
31

- 1     **[0225]**           EC189) The method of EC187, further comprising:  
2                   in response to the host querying a particular one of the secondary agents regarding  
3                   storage capacity, the particular secondary agent reporting back a particular  
4                   storage capacity that excludes representation of those portions of storage  
5                   implemented by the particular secondary agent that are represented in an  
6                   aggregated storage capacity determined by the primary agent.  
7
- 8     **[0226]**           EC190) The method of EC187, further comprising:  
9                   performing at least some of the primary agent communications with the host in  
10                   accordance with a host-interface protocol.  
11
- 12    **[0227]**           EC191) The method of EC187, further comprising:  
13                   performing at least some of the primary agent communications with the host and at least  
14                   some primary agent communications with at least one of the secondary agents  
15                   via sharing a same physical link; and  
16                   wherein the primary agent communications with at least one of the secondary agents  
17                   comprise at least some of the sub-requests and at least some of the sub-statuses  
18                   exchanged between the primary agent and the secondary agents)  
19
- 20    **[0228]**           EC192) The method of EC187, further comprising:  
21                   performing at least some primary agent communications with the secondary agents in  
22                   accordance with a peer-to-peer protocol; and  
23                   wherein the primary agent communications with the secondary agents comprise at least  
24                   the sub-requests and the sub-statuses exchanged between the primary agent and  
25                   the secondary agents.  
26
- 27    **[0229]**           EC193) The method of EC192, wherein:  
28                   the primary agent communications with the secondary agents further comprise  
29                   exchanging redundancy information via the peer-to-peer protocol.  
30

- 1     **[0230]**         EC194) An apparatus comprising:  
2             means for operating as a primary agent a first storage device and operating as one or  
3             more secondary agents at least a second storage device, each storage device of  
4             the primary and secondary agents having at least one respective storage range;  
5             means for the primary agent providing a logical interface enabling the host to access as a  
6             unitary logical device an aggregation of the storage ranges distributed over the  
7             storage devices of the primary and secondary agents, primary agent  
8             communications with the host via the logical interface comprising storage  
9             access requests accepted by the primary agent from the host and overall status  
10            sent by the primary agent to the host;  
11            means for the primary agent forwarding each request accepted as one or more sub-  
12            requests to all or any portions of the secondary agents, accepting sub-statuses  
13            associated with the sub-requests from the secondary agents that the sub-requests  
14            were forwarded to, and formulating each overall status sent based at least on the  
15            sub-statuses; and  
16            means for the all or any portions of the secondary agents accepting the sub-requests,  
17            translating host-context addressing information of the sub-requests to local-  
18            context addressing information, accessing local storage based at least in part on  
19            the local-context addressing information, and providing the sub-statuses to the  
20            primary agent.
- 21
- 22     **[0231]**         EC195) The apparatus of EC194, further comprising:  
23             means for the primary agent determining an aggregated storage capacity of the unitary  
24             logical device based at least in part on the storage ranges of the primary and  
25             secondary agents;  
26             means for in response to the host querying the unitary logical device regarding storage  
27             capacity, the primary agent reporting the aggregated storage capacity; and  
28             means for in response to the host querying the secondary agents regarding storage  
29             capacity, each of the secondary agents reporting a respective storage capacity  
30             that excludes representation of any storage range of the secondary agent that is  
31             represented in the aggregated storage capacity.
- 32

- 1     **[0232]**           EC196) The apparatus of EC194, further comprising:  
2           means for in response to the host querying a particular one of the secondary agents  
3           regarding storage capacity, the particular secondary agent reporting back a  
4           particular storage capacity that excludes representation of those portions of  
5           storage implemented by the particular secondary agent that are represented in an  
6           aggregated storage capacity determined by the primary agent.  
7
- 8     **[0233]**           EC197) The apparatus of EC194, further comprising:  
9           means for performing at least some of the primary agent communications with the host  
10          in accordance with a host-interface protocol.  
11
- 12    **[0234]**           EC198) The apparatus of EC194, further comprising:  
13          means for performing at least some of the primary agent communications with the host  
14          and at least some primary agent communications with at least one of the  
15          secondary agents via sharing a same physical link; and  
16          wherein the primary agent communications with at least one of the secondary agents  
17          comprise at least some of the sub-requests and at least some of the sub-statuses  
18          exchanged between the primary agent and the secondary agents.  
19
- 20    **[0235]**           EC199) The apparatus of EC194, further comprising:  
21          means for performing at least some primary agent communications with the secondary  
22          agents in accordance with a peer-to-peer protocol; and  
23          wherein the primary agent communications with the secondary agents comprise at least  
24          the sub-requests and the sub-statuses exchanged between the primary agent and  
25          the secondary agents.  
26
- 27    **[0236]**           EC200) The apparatus of EC199, wherein:  
28          the primary agent communications with the secondary agents further comprise  
29          redundancy information exchanged via the peer-to-peer protocol.



## 1 SCALABLE STORAGE SYSTEM

2  
3 **[0237]** Fig. 1A illustrates selected structural details of an embodiment of a technique  
4 for scalable storage devices, including a host, host visible storage having one or more storage  
5 devices operable as respective primary agents, and host invisible storage having one or more  
6 storage devices operable as respective secondary agents. Host **100** is coupled to Host Visible  
7 Storage **110** and Host Invisible Storage **120** via Host - Storage Device Coupling **180**. The Host -  
8 Storage Device Coupling as well as the Host Visible and Host Invisible Storages are optionally  
9 implemented as a pluggable module (illustrated as Add-In Card **190A**) and/or coupling **101** is  
10 optionally implemented as a cable. In some embodiments, all or any portions of the add-in card  
11 are implemented as an SSD. The Host Visible Storage has one or more storage devices  
12 (illustrated as Storage Device **110.A** and Storage Device **110.N**). Similarly, the Host Invisible  
13 Storage has one or more storage devices (illustrated as Storage Device **120.A** and Storage  
14 Device **120.N**). In various embodiments, any one or more of the storage devices are physical  
15 storage devices, such as SSDs.

16  
17 **[0238]** Couplings **101**, **111.A**, **111.N**, **121.A**, and **121.N** enable request, status, and data  
18 transfers between the Host, the Host Visible Storage, and the Host Invisible Storage. One or  
19 more of the couplings enable the transfers via a host-interface protocol (such as the Host acting  
20 as a master and one of the elements of the Host Visible Storage operating as a slave). One or  
21 more of the couplings enable the transfers via a peer-to-peer protocol (such as one of the  
22 elements of the Host Visible Storage operating as a primary agent and one of the elements of the  
23 Host Invisible Storage or one of the other elements of the Host Visible Storage operating as a  
24 secondary agent). In various embodiments, one or more of the couplings are compatible with an  
25 interface standard (such as PCIe or Infiniband). In various embodiments, Host - Storage Device  
26 Coupling **180** is implemented via one or more PCIe and/or Infiniband switches. In some  
27 embodiments, Host - Storage Device Coupling **180** is integrated with Host **100** as all or any  
28 portions of a host computing complex.

29  
30 **[0239]** The Host Visible Storage elements implement storage, and during initialization,  
31 the elements are configured to enable Host access of at least some of the implemented storage,  
32 thus providing storage that is 'visible' to the Host such as via a host-interface protocol. The  
33 Host Invisible Storage elements also implement storage, but during initialization, the elements  
34 are configured to disable Host access of at least some of the implemented storage, thus providing  
35 storage that is 'invisible' to the Host. However, the storage that is invisible to the host is

1 indirectly accessible by the Host via elements of the Host Visible Storage, such as via a peer-to-  
2 peer protocol.

3  
4 **[0240]** Dashed-arrow **151** conceptually illustrates information transfer between Host  
5 **100** and Storage Device **110.A**, and is representative of information transfer between the Host  
6 and any of the elements of Host Visible Storage **110**. The information transfers include one or  
7 more of a request originating from the Host to access storage, a status response relating to the  
8 request, and a data transfer relating to the request. The request includes one or more of an  
9 address (such as an LBA) and a length (such as in bytes or LBA quanta). The information  
10 transfers from the host to the storage device are communicated via coupling **101**, Host - Storage  
11 Device Coupling **180**, and then coupling **111.A**, and in reverse for the information transfers from  
12 the storage device to the host.

13  
14 **[0241]** Dashed-arrow **152** conceptually illustrates information transfer between Storage  
15 Device **110.A** and Storage Device **120.A**, and is representative of peer-to-peer information  
16 transfer between any of the elements of Host Visible Storage **110** and any of the elements of  
17 Host Invisible Storage **120**. The information transfers include one or more of a sub-request from  
18 Storage Device **110.A** (acting as a primary agent) to any of the elements of Host Invisible  
19 Storage **120** (acting as a secondary agent), a sub-status response relating to the sub-request, and  
20 a data transfer relating to the sub-request. The sub-request includes one or more of an address  
21 (such as an LBA) and a length (such as in bytes or LBA quanta). The information transfers from  
22 the primary agent to the secondary agent are communicated via coupling **111.A**, Host - Storage  
23 Device Coupling **180**, and then coupling **121.A**, and in reverse for the information transfers from  
24 the secondary agent to the primary agent.

25  
26 **[0242]** Dashed-arrow **153** conceptually illustrates information transfer between Storage  
27 Device **110.A** and Storage Device **110.N**, and is representative of peer-to-peer information  
28 transfer between any of the elements of Host Visible Storage **110**. The information transfers  
29 include one or more of a sub-request from Storage Device **110.A** (acting as a primary agent) to  
30 any of the other elements of Host Visible Storage **110** (acting as a secondary agent), a sub-status  
31 response relating to the sub-request, and a data transfer relating to the sub-request. The  
32 information transfers from the primary agent to the secondary agent are communicated via  
33 coupling **111.A**, Host - Storage Device Coupling **180**, and then coupling **111.N**, and in reverse  
34 for the information transfers from the secondary agent to the primary agent.

1     **[0243]**           Dashed-arrow **154** conceptually illustrates information transfer between Host  
2     **100** and Storage Device **120.A**, and is representative of information transfer between the Host  
3     and any of the elements of Host Invisible Storage **120**. The information transfers include one or  
4     more of an overall status response relating to the request, and a data transfer relating to the  
5     request. Although not illustrated as such for clarity in the figure, the information transfers from  
6     the host to the storage device are communicated via coupling **101**, Host - Storage Device  
7     Coupling **180**, and then coupling **121.A**, and in reverse for the information transfers from the  
8     storage device to the host.

10    **[0244]**           In some embodiments, requests originating from a host are via one or more  
11    submission queues in host memory that are accessible via coupling **101** (e.g. by a primary agent  
12    reading entries in a submission queue via one or more DMA operations). In some embodiments,  
13    statuses provided to a host are via one or more completion queues in host memory that are  
14    accessible to primary and/or secondary agents (e.g. by a primary/secondary agent writing entries  
15    in a completion queue via one or more DMA operations). In some embodiments, requests  
16    and/or statuses are communicated at least in part via one or more doorbell and/or tail registers,  
17    optionally in conjunction with one or more submission and/or completion queues. In some  
18    embodiments, requests and/or statuses are communicated to primary agents at least in part via  
19    programmed IO operations initiated by a host. In some embodiments, submission and/or  
20    completion queue descriptors are communicated to primary agents at least in part via  
21    programmed IO operations initiated by a host.

23    **[0245]**           In some embodiments, accepting a request from a host includes reading one or  
24    more entries from a submission queue. E.g. a host writes a particular entry in a submission  
25    queue retained in host memory, modifies a tail register accordingly, and sets a corresponding  
26    doorbell register (such as a doorbell register implemented in and/or accessible by a primary  
27    agent). In response to the setting of the doorbell register, the primary agent accesses the tail  
28    register and reads the particular submission queue entry (such as via a DMA read from the host  
29    memory). In some embodiments, providing status to a host includes writing one or more entries  
30    in a completion queue. E.g. a primary agent writes a particular entry in a completion queue  
31    retained in host memory (such as via a DMA write to the host memory), modifies a tail register  
32    accordingly, and notifies the host of a completed operation (e.g. by setting a corresponding  
33    doorbell register or providing an interrupt). In response to the setting of the doorbell register,  
34    the host accesses the tail register and reads the particular completion queue entry.

1     **[0246]**           In some embodiments, sub-requests are sent to secondary agents using different  
2     routing, addressing, and/or encapsulation than a request the sub-requests are determined from,  
3     and the sub-requests are transported over a same coupling as the request. For example, in an  
4     exemplary system using PCIe couplings for transport of any or all of request, sub-request, and  
5     sub-status traffic, a request is provided to a primary agent as a completion from a read request,  
6     while a sub-request is forwarded to a secondary agent in an MsgD request. In some  
7     embodiments, a sub-request is repackaged to include information in addition to a request the  
8     sub-request is determined from. For example, a repackaged sub-request includes additional  
9     information such as flow control, augmented routing information, or other additional  
10    information about contents of the request. In some embodiments, augmented routing  
11    information enables sending the sub-request via one or more couplings to remote secondary  
12    agents.

13  
14   **[0247]**           In some embodiments, sub-statuses are sent to primary agents using different  
15    routing, addressing, and/or encapsulation than a request the sub-statuses correspond to, and the  
16    sub-statuses are transported over a same coupling as the request. For example, in an exemplary  
17    system using PCIe couplings for transport of any or all of request, sub-request, and sub-status  
18    traffic, a request is provided to a primary agent as a completion from a read request, while a sub-  
19    status is returned to the primary agent by a secondary agent in an MsgD request. In some  
20    embodiments, a sub-status is repackaged to include information in addition to the sub-status  
21    and/or a request the sub-status corresponds to. For example, a repackaged sub-status includes  
22    additional information such as flow control, augmented routing information, or other additional  
23    information about contents of the request. In some embodiments, augmented routing  
24    information enables sending the sub-status via one or more couplings from remote secondary  
25    agents to primary agents.

26  
27   **[0248]**           As a specific example, Host **100**, acting as a master, provides a request to access  
28    storage (conceptually illustrated by dashed-arrow **151**) via a host-interface protocol to Storage  
29    Device **110.A**. In response, Storage Device **110.A**, acting as a slave, responds to the request.  
30    Storage Device **110.A** accepts the request (also conceptually illustrated by dashed-arrow **151**)  
31    via the host-interface protocol, and then determines whether storage relating to the request is  
32    implemented by Storage Device **110.A** and/or by one or more of the elements of Host Invisible  
33    Storage **120** or any of the other elements of Host Visible Storage **110**. Storage Device **110.A**  
34    then processes the request internally (if at least some of the storage is implemented by Storage  
35    Device **110.A**) and/or, acting as a primary agent, forwards the request, as one or more peer-to-  
36    peer sub-requests (conceptually illustrated by dashed-arrow **152**) via the peer-to-peer protocol to

Storage Device **120.A**. Storage Device **120.A**, acting as a secondary, accepts the sub-request (also conceptually illustrated by dashed-arrow **152**) via the peer-to-peer protocol, and then processes the sub-request internally. Storage Device **120.A** then returns sub-status corresponding to the sub-request to Storage Device **110.A** (also conceptually illustrated by dashed-arrow **152**). Storage Device **110.A** then determines an overall status for the request and provides the overall status to the host (conceptually illustrated by dashed-arrow **151**). Alternatively, rather than Storage Device **120.A** returning the sub-status to Storage Device **110.A** (to then forward to Host **100**), Storage Device **120.A** returns the sub-status to Host **100** directly (conceptually illustrated by dashed-arrow **154**).

[0249] Continuing with the specific example, data for reads flows similarly to status. Storage Device **120.A** returns data corresponding to the sub-request to Storage Device **110.A** (conceptually illustrated by dashed-arrow **152**) and then Storage Device **110.A** provides the data to the host (conceptually illustrated by dashed-arrow **151**). Alternatively, rather than Storage Device **120.A** returning the data to Storage Device **110.A** (to then forward to Host **100**), Storage Device **120.A** returns the data to Host **100** directly (conceptually illustrated by dashed-arrow **154**). Status and data for a particular request are not restricted to a same flow. For example, for one request, status and data both flow from a secondary device (e.g. Storage Device **120.A**) via a primary agent (e.g. Storage Device **110.A**) to Host **100**. For another request, status flows from the secondary agent via the primary agent to the host and data flows directly from the secondary agent to the host. Data for writes flows from the host to the secondary agent via the primary agent (conceptually illustrated by dashed-arrows **151** and **152**), or alternatively directly from the host to the secondary agent (conceptually illustrated by dashed-arrow **154**).

[0250] Other examples include any one or more of the elements of Host Visible Storage **110** acting as one or more primary agents and any one or more of the elements of Host Invisible Storage **120** or any one or more of the other elements of Host Visible Storage **110** acting as one or more secondary agents.

[0251] Fig. 1B illustrates selected structural details of another embodiment of a technique for scalable storage devices, including a host, host visible storage having one or more storage devices operable as respective primary agents, and host invisible storage having one or more storage devices operable as respective secondary agents. The technique uses some elements that are identical in operation and structure to similarly identified elements of Fig. 1A (Host **100**, Host - Storage Device Coupling **180**, and couplings **101**, **111.A**, **111.N**, **121.A**, and **121.N**). Host Visible Storage **130** is similar to Host Visible Storage **110** of Fig. 1A, except that

the storage devices therein (Storage Device **130.A** and Storage Device **130.N**) are each enabled to communicate via an additional coupling (couplings **131.A** and **131.N**, respectively). Host Invisible Storage **140** is similar to Host Invisible Storage **120** of Fig. 1A, except that the storage devices therein (Storage Device **140.A** and Storage Device **140.N**) are enabled to communicate via an additional coupling (couplings **141.A** and **141.N**, respectively). Device - Storage Device Coupling **181** is coupled to couplings **131.A**, **131.N**, **141.A** and **141.N**, and enables additional communication bandwidth between any of the elements of Host Visible Storage **130** and Host Invisible Storage **140**, useful, e.g., for peer-to-peer communication. For example, Device - Storage Device Coupling **181** enables additional communication bandwidth between Storage Device **130.A** and any of the elements of Host Invisible Storage **140**. For another example, Device - Storage Device Coupling **181** enables additional communication bandwidth between Storage Device **130.A** and Storage Device **130.N**.

**[0252]** The Host - Storage Device Coupling and the Device - Storage Device Coupling as well as the Host Visible and Host Invisible Storages are optionally implemented as a pluggable module (illustrated as Add-In Card **190B**) and/or coupling **101** is optionally implemented as a cable. In some embodiments, all or any portions of the add-in card are implemented as an SSD. In various embodiments, any one or more of the storage devices are SSDs.

**[0253]** In some embodiments and/or usage scenarios, additional communication bandwidth provided by Device - Storage Device Coupling **181** enables storage redundancy techniques implemented via peer-to-peer communication (such as RAID 5), with little or no impact on bandwidth available for communication between the host and the storage devices. The peer-to-peer communication includes one or more of communication of untransformed data (such as during a RAID mirroring operation), transformed data (such as XOR data during RAID parity generation and/or writing), and control information.

**[0254]** In some embodiments, RAID redundancy is enabled by distributing RAID data in a striped manner across a plurality of storage devices, such as orthogonal to how (non-RAID) data is distributed. In some situations, when a host request is a read, information is read from two or more of the storage devices. When a host request is a write, maintaining RAID redundancy information optionally includes performing RMW operations or alternatively data copying operations with respect to two or more of the storage devices.

1     **[0255]**           When a primary agent or a secondary agent retrieves data as a result of a host  
2     write request, the agent communicates a peer-to-peer request (such as via Device - Storage  
3     Device Coupling **181**) to another agent (a primary agent or a secondary agent), along with  
4     particular RAID data. In some circumstances (such as mirror RAID), the particular RAID data  
5     is a copy of data for the write request. In some circumstances (such as RAID 5), the particular  
6     RAID data is computed based on the data for the write request, and previous data corresponding  
7     to an LBA of the write request (e.g. the previous data is used in a RMW operation such as an  
8     XOR with the data for the write request).

9  
10    **[0256]**           For example, the agent receiving the data for the write request computes a  
11    RAID delta by XORing the data for the write request with the previous data corresponding to the  
12    LBA of the write request. The agent sends the RAID delta to the one of the agents storing RAID  
13    redundancy associated with the LBA. The agent storing the RAID redundancy then performs an  
14    RMW to update the RAID redundancy using the RAID delta. If the agent storing the RAID  
15    redundancy has been informed that multiple updates to the RAID redundancy are expected, then  
16    the agent storing the RAID redundancy optionally and/or selectively combines the RMW  
17    operations for greater efficiency.

18  
19    **[0257]**           In some embodiments and/or usage scenarios, a primary agent informs another  
20    agent (another primary agent or a secondary agent) how many RAID updates there will be for  
21    particular LBAs. The informed agent is enabled, in various embodiments, to cache information  
22    relating to the RAID updates until all of the updates for the particular LBAs are performed. In  
23    various embodiments, the number of RAID updates is communicated as part of a sub-request,  
24    either implicitly or explicitly. When the RAID updates are completed, the informed agent sends  
25    a sub-status to the primary agent.

26  
27    **[0258]**           In various embodiments, a plurality of data transfers relating to a single host  
28    request is performed wholly or partially independently. For example, in some RAID RMW  
29    scenarios, transfer of write data is independent of transfer of previously accumulated parity data.  
30    In various situations, the write data is read from host memory by any of a primary agent and a  
31    plurality of secondary agents, the accumulated parity data is read from a storage device by any  
32    of the primary and secondary agents, and then new parity data is written, e.g., to the storage  
33    device the accumulated parity data was read from.

34  
35    **[0259]**           In various embodiments, one or more of couplings **131.A**, **131.N**, **141.A** and  
36    **141.N** are identical or substantially similar to couplings **111.A**, **111.N**, **121.A**, and **121.N** of Fig.

1 1A. In some embodiments, device-storage coupling Device - Storage Device Coupling **181** is  
2 identical or substantially similar in implementation to Host - Storage Device Coupling **180**. In  
3 some embodiments, Host - Storage Device Coupling **180** and Device - Storage Device Coupling  
4 **181** are combined into a single element. In some embodiments, any one or more of couplings  
5 **131.A**, **131.N**, **141.A** and **141.N** are coupled to additional ports of a variation of Host - Storage  
6 Device Coupling **180** instead of to Device - Storage Device Coupling **181**.

7  
8 **[0260]** In various embodiments, such as some embodiments related to either of Fig. 1A  
9 or Fig. 1B, host visible storage has one or more storage devices operable as respective primary  
10 agents and zero or more storage devices operable as secondary agents, while host invisible  
11 storage has zero or more storage devices operable as respective secondary agents.

12  
13 **[0261]** In various embodiments, a storage device corresponds to a physical storage  
14 device, and the physical storage device is enabled to implement, according to one or more  
15 operating modes, any combination of host visible storage and host invisible storage. In a first  
16 operating scenario, a particular physical device is operated to implement only host visible  
17 storage. In a second operating scenario, the particular physical device is operated to implement  
18 only host invisible storage. In a third operating scenario, the particular physical device is  
19 operated to implement a combination of host visible storage and host invisible storage. In some  
20 situations (such as in response to a change in workload and/or a failure), operation is  
21 dynamically changed from one of the operating scenarios to another.

22  
23 **[0262]** For example, with respect to Fig. 1A, operation begins as illustrated, with Host  
24 Visible Storage **110** having Storage Device **110.A** and Storage Device **110.N**, and Host Invisible  
25 Storage **120** having Storage Device **120.A** and Storage Device **120.N**. Continuing with the  
26 example, each of Storage Devices **110.A**, **110.N**, **120.A** and **120.N** corresponds to respective  
27 physical storage devices. In response to a change in workload, Storage Device **110.N** is  
28 dynamically (re)configured from operating as an element of the Host Visible Storage to  
29 operating as an element of the Host Invisible Storage. After the dynamic reconfiguration, Host  
30 Visible Storage **110** has Storage Device **110.A**, and Host Invisible Storage **120** has Storage  
31 Device **110.N**, Storage Device **120.A**, and Storage Device **120.N** (not illustrated).

32  
33 **[0263]** For another example, with respect to Fig. 1B, operation begins as illustrated,  
34 with Host Visible Storage **130** having Storage Device **130.A** and Storage Device **130.N**, and  
35 Host Invisible Storage **140** having Storage Device **140.A** and Storage Device **140.N**. Continuing  
36 with the example, each of Storage Devices **130.A**, **130.N**, **140.A**, and **140.N** corresponds to



1     respective physical storage devices. In response to a failure of Storage Device **130.A**, Storage  
2     Device **140.A** is dynamically (re)configured to operate as an element of the Host Visible Storage  
3     as well as continuing to operate as an element of the Host Invisible Storage. For instance,  
4     storage implemented by Storage Device **140.A** is allocated between being used as replacing all  
5     or any portions of storage implemented by Storage Device **130.A** before the failure, and  
6     continuing use as Host Invisible Storage. After the dynamic reconfiguration, Host Visible  
7     Storage **130** has a portion of storage implemented by Storage Device **140.A** and Storage Device  
8     **130.N**, and Host Invisible Storage **140** has Storage Device **140.A**, (less the portion used to  
9     replace Storage Device **130.A**) and Storage Device **140.N** (not illustrated).

10  
11     **[0264]**         In various embodiments, host visible storage is storage that is accessible to a  
12     host via a transparent switch, and/or host invisible storage is storage that is hidden from the host  
13     via a non-transparent switch. The host visible storage and the host invisible storage are enabled  
14     for peer-to-peer communication, e.g., via all or any portions of the transparent and/or the non-  
15     transparent switches. In some embodiments, a unified switch provides a transparent portion and  
16     a non-transparent portion corresponding respectively to the aforementioned transparent and non-  
17     transparent switches. For example, with respect to Fig. 1B, one or more of all or any portions of  
18     Host **100**, coupling **101**, Host - Storage Device Coupling **180**, and Device - Storage Device  
19     Coupling **181** implement a unified switch. The unified switch is operable as a transparent switch  
20     with respect to one or more of the elements of Host Visible Storage **130** and also simultaneously  
21     operable as a non-transparent switch with respect to one or more of the elements of Host  
22     Invisible Storage **140**.

## 23 24 25     SCALABLE STORAGE SYSTEM OPERATION

26  
27     **[0265]**         Fig. 2 illustrates selected processing details of an embodiment of a technique for  
28     scalable storage devices. Illustrated are actions performed by a storage device acting as a  
29     primary agent (such as one or more of the elements of Host Visible Storage **110** of Fig. 1A) as  
30     “Primary Actions” **209**, as well as actions performed by a storage device acting as a secondary  
31     agent (such as one or more of the elements of Host Invisible Storage **120** and/or any one or more  
32     of the other elements of Host Visible Storage **110** of Fig. 1A) as “Secondary Actions” **219**.

33  
34     **[0266]**         A scenario where a request to access storage is satisfied by a single responding  
35     agent is first described, and then scenarios where a single request is satisfied by more than one  
36     responding agent are described. A host (such as Host **100** of Fig. 1A) provides a request to

1 access storage to a storage device acting as a primary agent (represented conceptually by “Start”  
 2 **201**). The primary agent accepts the request from the host, such as by reading a submission  
 3 queue entry retained in host memory, optionally providing handshaking with the host to indicate  
 4 the request has been accepted without errors (“Accept Req from Host” **202**). The primary agent  
 5 then examines the request to determine if the request is to storage implemented by the primary  
 6 agent (“Local?” **203**). If so (“Yes”, **203Y**), then flow proceeds where the primary agent  
 7 processes the address information provided with the request to translate the address information  
 8 from a (host) context associated with the request to a (primary agent) context associated with  
 9 corresponding storage implemented by the primary agent (“Xlate to Local LBAs” **204**). For an  
 10 example, the request specifies LBA 100, and the primary agent translates LBA 100 to LBA 0,  
 11 corresponding to a beginning location of (local) storage implemented by the primary agent. The  
 12 primary agent then performs the request (“Process Locally” **205**), and provides status to the host  
 13 (“Provide Status to Host” **206**), such as by writing a completion queue entry retained in the host  
 14 memory, thus completing processing for the request (“End” **299**). In some embodiments, the  
 15 translation comprises a mapping between a function of the LBAs and local storage addresses,  
 16 similar to a page table.

17  
 18 **[0267]** If the request is not to storage implemented by the primary agent (“No”, **203N**),  
 19 then flow within the primary agent proceeds to send the request, as a sub-request, to a storage  
 20 device acting as a secondary agent (“Forward to Secondary(s)” **207**). The primary agent  
 21 determines a sub-request from the request, and communication of the sub-request from the  
 22 primary agent to the secondary agent is illustrated by dashed-line “Sub-Req” **207R**. Within the  
 23 primary agent, flow then proceeds to await completion of the sub-request by the secondary agent  
 24 (“Wait for Completion” **208**). When the secondary agent has completed the sub-request (either  
 25 with or without errors), the secondary agent supplies sub-status to the primary agent (illustrated  
 26 conceptually by dashed-line “Sub-Status” **214S**). The primary agent then accepts the sub-status  
 27 from the secondary agent, optionally providing handshaking with the secondary agent to indicate  
 28 the sub-status has been accepted without errors (“Accept Sub-Status from Secondary(s)” **209A**).  
 29 The primary agent then proceeds to provide the sub-status as status to the host (“Provide Status  
 30 to Host” **206**), and then processing for the request is complete (“End” **299**).

31  
 32 **[0268]** Within the secondary agent, flow begins by accepting the sub-request from the  
 33 primary agent, optionally providing handshaking with the primary agent to indicate the sub-  
 34 request has been accepted without errors (“Accept Sub-Req from Primary” **211**). The sub-  
 35 request includes address information (such as an LBA) as provided by the host in the request  
 36 that the sub-request corresponds to. The secondary agent processes the address information

1 provided with the sub-request to translate the address information from a (host) context  
2 associated with the request to a (secondary agent) context associated with corresponding storage  
3 implemented by the secondary agent (“Xlate to Local LBAs” **212**). For example, a request  
4 specifies LBA 200, the sub-request specifies LBA 200, and the secondary agent translates (host)  
5 LBA 200 to (local) LBA 0, corresponding to a beginning location of (local) storage implemented  
6 by the secondary agent. Continuing with the example, another request specifies LBA 201, a  
7 corresponding sub-request specifies LBA 201, and the secondary agent translates (host) LBA  
8 201 to (local) LBA 1, corresponding to a next location of the (local) storage implemented by the  
9 secondary agent. Other examples are described, following, in a description of Fig. 3. After  
10 translating the sub-request, flow then proceeds where the secondary agent performs the request  
11 (“Process Locally” **213**) and then returns sub-status to the primary agent (“Provide Sub-Status to  
12 Primary” **214**) as illustrated by dashed-line “Sub-Status” **214S**.

13

14 **[0269]** In a scenario where a single request to access storage is satisfied by more than  
15 one responding agent, conceptually the single request is broken into a plurality of sub-requests.  
16 Zero, one, or more of the sub-requests are satisfied by the primary agent, and zero, one, or more  
17 of the sub-requests are satisfied by zero, one, or more secondary agents. For example, a request  
18 to LBA 400 with a length of two LBAs is satisfied by a sub-request to a secondary agent for  
19 LBA 400 and another sub-request to another secondary agent for LBA 401. Each secondary  
20 agent independently translates the LBA of the respective sub-request to an LBA corresponding  
21 to local storage of the respective secondary agent.

22

23 **[0270]** The primary agent determines which (if any) secondary agents are to receive  
24 sub-requests, based on address and length information of the request. Thus in Fig. 2,  
25 determining if the request is to storage implemented by the primary agent (“Local?” **203**) is  
26 conceptually determined for each of a plurality of sub-requests. For some of the sub-requests  
27 the outcome is that the sub-requests are local, while for others of the sub-requests the outcome is  
28 that the other sub-requests are not local, and are forwarded along to one or more secondary  
29 agents as sub-requests (“Forward to Secondary(s)” **207** and “Sub-Req” **207R**). A single  
30 secondary agent receives, in some scenarios, a single sub-request, while in other scenarios the  
31 single secondary agent receives a plurality of sub-requests. The flow represented in “Secondary  
32 Actions” **219** is then performed for each of the plurality of sub-requests (in any combination of  
33 serial and/or parallel processing, in various embodiments). In some scenarios, each of a plurality  
34 of secondary agents receives one or more sub-requests. The flow represented in “Secondary  
35 Actions” **219** is then performed independently by each of the plurality of secondary agents.

36

1     **[0271]**           The primary agent determines how to map a single request into a plurality of  
2     sub-requests to one or more secondary agents as part of processing performed in “Forward to  
3     Secondary(s)” **207**, such as via one or more functions and/or mapping tables (implemented, e.g.  
4     in Storage Device **110.A** of Fig. 1A). When a plurality of sub-requests are sent to one or more  
5     secondary agents, the primary agent conceptually accumulates status information by waiting for  
6     completion (“Wait for Completion” **208**) and accepting corresponding sub-status (“Accept Sub-  
7     Status from Secondary(s)” **209A**) for each of the sub-requests independently of the other sub-  
8     requests. When sub-statuses for all of the sub-requests for a particular request have been  
9     received, the primary agent produces an overall status to provide to the host as part of processing  
10    performed in “Provide Status to Host” **206**.

11  
12   **[0272]**           In some embodiments, secondary agents receive an entirety of a request as  
13   identical sub-requests (or alternately read the request from an entry in a submission queue).  
14   Each of the secondary agents then determines which portion (if any) of the request the respective  
15   secondary agent is to perform, and optionally which agent to send accumulated parity data (if  
16   any) to. For example, a request specifies LBA 500 with a length of 200 LBAs. A primary agent  
17   sends identical sub-requests to first and second secondary agents (or alternatively a single sub-  
18   request directed to the first and the second secondary agents). The sub-request(s) specify LBA  
19   500 with a length of 200 LBAs. The first secondary agent (independently of the host and the  
20   second secondary agent) determines that LBAs 600-699 of the sub-request correspond to a  
21   portion of the sub-request that the first secondary agent is to perform. The second secondary  
22   agent (independently of the host and the first secondary agent) determines that LBAs 500-599 of  
23   the sub-request correspond to a portion of the sub-request that the second secondary agent is to  
24   perform. The first and the second secondary agents determine that accumulated parity  
25   information is an appropriate one of the agents, dependent on a particular redundancy scheme.

26  
27   **[0273]**           In some embodiments and/or usage scenarios, status is provided directly to the  
28   host by a secondary agent, rather than forwarding through a primary agent (conceptually  
29   illustrated by dashed-arrow **154** of Fig. 1A). For example, the direct providing of status occurs  
30   when sub-requests for an entirety of a request are forwarded to a single secondary agent. The  
31   single secondary agent then returns status directly to the host (for example by the single  
32   secondary agent writing an entry to a completion queue retained in host memory). Note that the  
33   action of the single secondary agent writing the entry is indistinguishable by the host from the  
34   primary agent writing the entry, and thus even though the single secondary agent writes the  
35   entry, satisfaction of the request by the single secondary agent is transparent to the host. In some  
36   embodiments, a secondary agent (rather than a primary agent) accumulates sub-statuses for a

1 plurality of sub-requests corresponding to a particular request and provides the accumulated sub-  
2 statuses as an overall status to the host.

3  
4 **[0274]** In some embodiments, one or more storage devices operate simultaneously as a  
5 primary agent and as a secondary agent, at least with respect to some requests from one or more  
6 hosts. For example, a system has two storage devices. The first of the storage devices operates  
7 as a primary agent with respect to the second of the storage devices operating as a secondary  
8 agent, and the second storage device operates as a primary agent with respect to the first storage  
9 device operating as a secondary agent. More specifically, the first storage device accepts a  
10 request from a particular one of the hosts and sends the request, as a sub-request, to the second  
11 storage device. The second storage device accepts another request from the particular host and  
12 forwards the other request, as a sub-request, to the first storage device. In some embodiments  
13 and/or usage scenarios, the requests are from a same host request queue. In some embodiments  
14 and/or usage scenarios (such as where the requests are from distinct hosts), the requests are from  
15 distinct host request queues.

16  
17 **[0275]** According to various embodiments, one or more of request, status, and data  
18 transfers are via DMA between, for example, memory accessible to (or part of) the host (such as  
19 host memory) and memory accessible to (or part of) one of the storage devices (such as a buffer  
20 implemented in a storage device). Thus, one or more of the aforementioned flows of request,  
21 status, and data transfers are by communication of one or more scatter/gather lists and  
22 subsequent DMA transfers according to the scatter/gather lists. In addition (or alternatively),  
23 one or more of the flows of request, status, and data transfers are by communication of  
24 completion and/or status queue information and subsequent DMA transfers to/from entries of the  
25 queues. In various embodiments, any of the request, the status, and the data transfers are  
26 performed directly and/or managed by one or more processors in one or more of the storage  
27 devices.

## 30 ADDRESSING

31  
32 **[0276]** Fig. 3 illustrates selected details of host, primary agent, and secondary agent  
33 addressing in an embodiment of a technique for scalable storage devices. The figure illustrates  
34 example windows of storage address space (Host Address Space **310**) as viewed by a host (such  
35 as Host **100** of Fig. 1A). The host storage address space has Host LBA Ranges 1-5 **311-315**.  
36 The figure further illustrates example windows of storage address spaces (Primary Address

1 Space **320**) as viewed by a primary agent (such as any one of the elements of Host Visible  
2 Storage **110** of Fig. 1A). The primary agent storage address space has Primary LBA Ranges 1-3  
3 **321-323**. The figure further illustrates example windows of storage address spaces (Secondary  
4 Address Spaces **330**) as viewed by two secondary agents (Secondary Agents A and B, such as  
5 any two of the elements of Host Invisible Storage **120** and/or any of the other elements of Host  
6 Visible Storage **110** of Fig. 1A). The Secondary Agent A storage address space has Secondary  
7 A LBA Ranges 1-3 **331A**, **332A** and **333A**. The Secondary Agent B storage address space has  
8 Secondary B LBA Ranges 1-2 **331B** and **332B**. While not drawn to any particular scale, in the  
9 figure the vertical height of the various address range elements varies, conceptually representing  
10 variation in how many locations are in each of the ranges.

11

12 **[0277]** In a first example, a single host address range (Host LBA Range 1 **311**)  
13 corresponds to a single address range of Secondary A (Secondary A LBA Range 1 **331A**). A  
14 primary agent (such as Storage Device **110.A** of Fig. 1A) determines that a host request to an  
15 address within Host LBA Range 1 **311** corresponds to storage that is not implemented by the  
16 primary agent (such as via processing associated with “Local?” **203** and “No”, **203N** of Fig. 2)  
17 and that the request corresponds to storage implemented by Secondary A. The primary agent  
18 then forwards a sub-request to Secondary A. Secondary A then translates the sub-request (such  
19 as processing associated with “Xlate to Local LBAs” **212** of Fig. 2), and performs a local access  
20 to storage accessible via Secondary A LBA Range 1 **331A** (such as processing via “Process  
21 Locally” **213** of Fig. 2).

22

23 **[0278]** In a second example, a single host address range (Host LBA Range 2 **312**)  
24 corresponds to a single address range of Secondary B (Secondary B LBA Range 1 **331B**). The  
25 primary agent determines that a request to an address within Host LBA Range 2 **312** corresponds  
26 to storage that is implemented by Secondary B, and then forwards a sub-request to Secondary B.  
27 Secondary B then translates the sub-request and performs a local access to storage accessible via  
28 Secondary B LBA Range 1 **331B**.

29

30 **[0279]** In a third example, a single host address range (Host LBA Range 3 **313**)  
31 corresponds to a single address range of the primary agent (Primary LBA Range 1 **321**). The  
32 primary agent determines that a request to an address within Host LBA Range 3 **313** corresponds  
33 to storage that is implemented by the primary agent, translates the request (such as via  
34 processing associated with “Xlate to Local LBAs” **204**), and performs a local access to storage  
35 accessible via Primary LBA Range 1 **321**.

36

1     **[0280]**           In a fourth example, a single host address range (Host LBA Range 4 **314**)  
2     corresponds to two address ranges of Secondary A (Secondary A LBA Ranges 2 **332A** and 3  
3     **333A**) and a single address range of Secondary B (Secondary B LBA Range 2 **332B**). The  
4     primary agent determines that a request to an entirety of LBAs of Host LBA Range 4 **314**  
5     (based, e.g. on a starting LBA and a number of contiguous LBAs to access) corresponds to  
6     storage that is implemented by Secondary A and Secondary B. The primary agent then forwards  
7     two sub-requests to Secondary A and a single sub-request to Secondary B. Secondary A  
8     translates the two sub-requests and performs local accesses to storage accessible via Secondary  
9     A LBA Ranges 2 **332A** and 3 **333A** (in series or wholly or partially in parallel). Secondary B  
10    translates the single sub-request and performs a local access to storage accessible via Secondary  
11    B LBA Range 2 **332B**.

12  
13    **[0281]**           In a fifth example, a single host address range (Host LBA Range 5 **315**)  
14    corresponds to two address ranges of the primary agent (Primary LBA Ranges 2 **322** and 3 **323**).  
15    The primary agent determines that a request to an entirety of LBAs of Host LBA Range 5 **315**  
16    (based, e.g. on a starting LBA and a number of contiguous LBAs to access) corresponds to  
17    storage that is implemented by the primary agent. The primary agent translates the request into  
18    two local LBA ranges, and performs local accesses to storage accessible via Primary LBA  
19    Ranges 2 **322** and 3 **323**.

20  
21    **[0282]**           In some embodiments, a primary agent is enabled to forward each sub-request  
22    of a host request to a particular one of one or more secondary agents separately. In other  
23    embodiments, a primary agent is enabled to forward two or more sub-requests of a host request  
24    to a particular one of a plurality of secondary agents as a single combined sub-request. For  
25    example, in various embodiments, a primary agent determines one or more secondary agents to  
26    process one or more sub-requests of a host request and forwards the host request with a same  
27    LBA range and a same length to the one or more secondary agents. The one or more secondary  
28    agents are enabled to then interpret the LBA range and the length to determine respective  
29    portions of the host request each of the one or more secondary agents is to process and respond  
30    to. In other embodiments, a primary agent processes an LBA range and a length of a host  
31    request and sends a processed version of the LBA range and the length to each of one or more  
32    secondary agents. For example, if data is striped among N agents (N-1 secondary agents and the  
33    primary agent) on 64KB boundaries, the primary agent divides the LBA range by N\*64KB to  
34    distribute sub-requests to the N-1 secondary agents. Each of the N-1 secondary agents receives  
35    a respective sub-request with a starting divided-down LBA, a portion of the length  
36    corresponding to striping of data among the N agents, and optionally and/or selectively a

1 remainder from dividing down the LBA; thus a first and/or a last of the secondary agents  
2 optionally process a transfer smaller than 64KB.

### 5 SCALABLE STORAGE DEVICES

7 **[0283]** Fig. 4 illustrates selected structural details of an embodiment of a scalable  
8 storage device enabled to operate as a primary agent, specifically Storage Device **110.A** of Fig.  
9 1A. The storage device includes an interface to a PCIe bus (PCIe Intfc **401**). The PCIe interface  
10 includes a DMA unit to enable optional and/or selective transfers directly between host memory  
11 and buffers of the PCIe interface, for example to read/write completion/status queue entries  
12 and/or to directly transfer data from/to host memory. The PCIe interface is coupled to a  
13 translation block (LBA, Length Xlate logic **402**). The translation block is enabled to translate  
14 host storage address information (e.g. LBA and length) into local storage address information,  
15 such as described with respect to “Xlate to Local LBAs” **204** of Fig. 2. The translation block is  
16 coupled to a mass storage interface (Storage Intfc logic **403**) that is coupled to one or more mass  
17 storage units (illustrated collectively as Mass Storage **404**). The mass storage units include any  
18 one or more of a flash storage unit (Flash **404F**), a magnetic disk storage unit (Magnetic **404M**),  
19 an optical disk storage unit (Optical **404O**), and/or any type of non-volatile storage unit. In  
20 various embodiments the flash storage unit comprises one or more flash memory chips and/or  
21 dice, such as NAND flash or NOR flash.

23 **[0284]** Storage Device **110.A** includes a control block enabled to perform various  
24 operations relating to operation as a primary agent (Sub-Request Generation & Sub-Status  
25 Accumulation **405**) that is optionally coupled to one or more of PCIe Intfc **401**, LBA, Length  
26 Xlate logic **402**, and Storage Intfc logic **403**. The operations include, for example, implementing  
27 and/or managing all or any portions of processing associated with any one or more of the  
28 elements of “Primary Actions” **209** of Fig. 2. In some embodiments, Sub-Request Generation &  
29 Sub-Status Accumulation **405** implements one or more functions and/or mapping tables that  
30 track how host addresses are allocated to various secondary devices (such as the mapping tables  
31 used by “Forward to Secondary(s)” **207** of Fig. 2). The operations further include responding to  
32 host enquiries regarding storage capacity as if storage capacity implemented by one or more  
33 secondary agents (such as Storage Device **120.A** of Fig. 5) were implemented by Storage Device  
34 **110.A**.



1   **[0285]**           Fig. 5 illustrates selected structural details of an embodiment of a scalable  
2   storage device enabled to operate as a secondary agent, such as Storage Device **120.A** of Fig.  
3   1A. The storage device includes some elements identical in operation and structure to similarly  
4   identified elements of Fig. 4 (PCIe Intfc **401**, LBA, Length Xlate logic **402**, Storage Intfc logic  
5   **403**, and Mass Storage **404** having elements Flash **404F**, Magnetic **404M**, and **404O**).

6  
7   **[0286]**           Storage Device **120.A** includes a control block enabled to perform various  
8   operations relating to operation as a secondary agent (Sub-Request Accepting & Sub-Status  
9   Generation **501**) that is optionally coupled to one or more of PCIe Intfc **401**, LBA, Length Xlate  
10   logic **402**, and Storage Intfc logic **403**. The operations include, for example, implementing  
11   and/or managing all or any portions of processing associated with any one or more of the  
12   elements of "Secondary Actions" **219** of Fig. 2. The operations further include responding to  
13   host enquiries regarding storage capacity as if there were no storage implemented by Storage  
14   Device **120.A**, or alternatively as if storage allocated for accessing by a host via a primary agent  
15   (rather than directly by the host) were not implemented by the storage device. Instead, the  
16   storage implemented by Storage Device **120.A** is reported by a primary agent, such as Storage  
17   Device **110.A** of Fig. 4.

18  
19   **[0287]**           In some embodiments (not illustrated), a configurable storage device includes  
20   elements corresponding to functionality of all of the elements illustrated in Storage Device  
21   **110.A** of Fig. 4 and further includes elements corresponding to functionality of Sub-Request  
22   Accepting & Sub-Status Generation **501** of Fig. 5. An example configurable storage device  
23   includes PCIe Intfc **401**, LBA, Length Xlate logic **402**, Storage Intfc logic **403**, and Mass  
24   Storage **404**, as well as Sub-Request Generation & Sub-Status Accumulation **405** and Sub-  
25   Request Accepting & Sub-Status Generation **501**. The configurable storage device is enabled to  
26   operate selectively as a primary agent or a secondary agent, for example in response to included  
27   mode state that is programmable. The mode state is programmable and/or alterable by various  
28   techniques, such as via the included PCIe interface by management software using special (e.g.  
29   vendor-specific) commands, via an initializing element such as non-volatile memory (e.g. flash  
30   memory or an option ROM), or via detection of an event by the configurable storage device  
31   and/or another storage device in communication with the configurable storage device. In some  
32   embodiments and/or usage scenarios, a configurable storage device enables a dynamic change in  
33   operating mode, such as a secondary agent dynamically replacing a failed primary agent, thus  
34   eliminating storage as a single point of failure.

## 1 FAILURE SCENARIOS

2

3 **[0288]** Several failure recovery techniques are enabled in various embodiments and/or  
4 usage scenarios where a primary agent as well as one or more secondary agents appear to a host  
5 as a single logical interface. The recovery techniques are either fully transparent to the host or  
6 are visible only to a device driver executing on the host.

7

8 **[0289]** A first failure recovery technique is applicable when a secondary agent fails  
9 wholly or partially, such that the failed secondary agent no longer implements a particular  
10 portion of storage. The primary agent becomes aware of the failure (e.g. the primary agent  
11 detects the failure or the failed secondary agent reports the failure to the primary agent). In  
12 response, the primary agent (re)allocates storage to implement the particular portion of storage,  
13 via allocation from storage of the primary agent and/or from zero or more of the secondary  
14 agents, optionally and/or selectively including the failed secondary agent. The primary agent  
15 and the secondary agents then operate according to any of the aforementioned embodiments  
16 and/or scenarios described with respect to Figs 1A, 1B, and 2-5.

17

18 **[0290]** In some situations, such as when a secondary agent fails such that the secondary  
19 agent is unable to contribute storage to the reallocation, the reallocation excludes the failed  
20 secondary agent. The secondary agent is unable to contribute storage for various reasons in  
21 various scenarios, such as having insufficient free storage that is operable, an inability to  
22 properly communicate with the primary agent and/or the host, physical removal, or total failure.  
23 In other situations, such as when a secondary agent partially fails and continues to implement at  
24 least some storage, the reallocation optionally includes storage of the partially failed second  
25 agent. In alternate related embodiments and/or usage scenarios, a device driver executing on the  
26 host becomes aware of the failure (rather than the primary agent) and informs the primary agent  
27 of the failure and/or directs the primary agent to perform a reallocation.

28

29 **[0291]** A second failure recovery technique is applicable when a primary agent fails  
30 wholly or partially, such that the primary agent no longer implements a particular portion of  
31 storage. The second failure recovery technique is similar to the first failure recovery technique,  
32 except that the reallocation is with respect to storage no longer implemented by the wholly or  
33 partially failing primary agent. As in the first failure recovery technique, the failure of the  
34 primary agent to properly implement storage is detectable by one or more of the primary agent, a  
35 secondary agent, or a device driver on a host.

36

1   **[0292]**           A third failure recovery technique is applicable when a primary agent fails to  
2 properly act as a primary agent, e.g. to communicate information (e.g. sub-requests and/or sub-  
3 statuses) with one or more secondary agents and/or a host. A monitoring agent becomes aware  
4 of the failure and in response identifies another primary or secondary agent to act as a  
5 replacement primary agent. In some embodiments, the replacement primary agent is a hot spare.  
6 An example of the monitoring agent is another primary (or secondary) agent of a storage device  
7 that includes the failed primary agent. Another example of the monitoring agent is a device  
8 driver on a host, such as a device driver enabling communication between an OS and/or  
9 applications executing on the host and the failed primary agent prior to the failure. If the failed  
10 primary agent failure also results in the failed primary agent no longer implementing a particular  
11 portion of storage, then the replacement primary agent reallocates storage to implement the  
12 particular portion of storage, such as in the second failure recovery technique. The reallocation  
13 is via allocation from storage of any one or more of: the failed primary agent (if any storage  
14 remains operable in the failed primary agent), the replacement primary agent, and/or zero or  
15 more of the secondary agents. The replacement primary agent acts as the failed primary agent  
16 did before the failure, and operates according to any of the aforementioned embodiments and/or  
17 scenarios described with respect to Figs 1A, 1B, and 2-5.

18  
19   **[0293]**           A fourth failure recovery technique is enabled in a system with one or more  
20 configurable agents, and is applicable when a primary agent fails to properly act as a primary  
21 agent. An example of the configurable agent is a configurable storage device operable either as  
22 a primary agent or a secondary agent, such as based on a configurable storage device that  
23 conceptually includes all elements of Fig. 4 and Fig. 5. Similar to the third failure recovery  
24 technique, in response to a monitoring agent becoming aware of failure of a primary agent to  
25 properly act as a primary agent, the monitoring agent identifies one of the configurable agents to  
26 act as a replacement primary agent. The monitoring agent requests that the identified  
27 configurable agent be configured (e.g. via programming mode information) to operate as a  
28 replacement primary agent. The identified configurable element dynamically responds to the  
29 configuration change request and begins acting as a primary agent to perform as a replacement  
30 primary agent. Similar to the third failure recovery technique, if the failed primary agent failure  
31 also results in the failed primary agent no longer implementing a particular portion of storage,  
32 then the replacement primary agent reallocates storage to implement the particular portion of  
33 storage, such as in the second failure recovery technique. The replacement primary agent acts as  
34 the failed primary agent did before the failure, and operates according to any of the  
35 aforementioned embodiments and/or scenarios described with respect to Figs 1A, 1B, and 2-5.

36

1     **[0294]**           In any one or more of the foregoing failure recovery techniques, after a failure,  
2     data is optionally recovered as possible via redundancy information. The data recovery is  
3     optionally via any combination of a host-side coupling (e.g. Host - Storage Device Coupling **180**  
4     of Figs. 1A or 1B) and a device-side coupling (e.g. Device - Storage Device Coupling **181** of  
5     Fig. 1B). As described elsewhere herein, the redundancy information is in accordance with, e.g.,  
6     various RAID and/or mirroring implementations.

7  
8     **[0295]**           In any one or more of the foregoing failure recovery techniques, the  
9     reallocations optionally reduce free space available, e.g., as represented to the host.

10  
11    **[0296]**           In any one or more of the foregoing failure recovery techniques, any one or  
12    more of the reallocations are performed and/or are managed at least in part by a processor  
13    executing instructions.

14  
15    **[0297]**           In any one or more of the foregoing failure recovery techniques, recognition of a  
16    failure of a primary or a secondary agent is via one or more of: a specific detection of a failure  
17    by an agent, a failure to receive a heartbeat indication by an agent, and a higher-level indication  
18    that a failure has occurred. The specific detection includes any one or more of a primary agent,  
19    a secondary agent, or a device driver determination that another agent has failed, e.g. via  
20    determining that a status returned is inappropriate, a protocol has been violated, or a request has  
21    timed out.

22  
23    **[0298]**           In some embodiments where a transparent switch enables communication  
24    between a host and host visible storage, a failure of an element of the host visible storage is  
25    visible, e.g., to a device driver executing on the host. The device driver optionally participates in  
26    recovery from the failure. In some embodiments having a non-transparent switch between a host  
27    and host invisible storage, a failure of an element of the host invisible storage is invisible to the  
28    host. Recovery from the failure is optionally invisible to the host.

29  
30    **[0299]**           In various embodiments, all or any portions of operations performed by any of  
31    the elements of Host Visible Storage **110**, Host Invisible Storage **120**, Host Visible Storage **130**,  
32    and Host Invisible Storage **140** of Figs. 1A and 1B, as well as units and sub-units within, are  
33    managed and/or implemented by one or more processors executing instructions (such as  
34    firmware) stored in an optionally tangible computer-readable medium. For example, Sub-  
35    Request Generation & Sub-Status Accumulation **405** and/or Sub-Request Accepting & Sub-  
36    Status Generation **501** are managed and/or implemented by a processor executing instructions

(such as firmware) stored in a computer-readable medium. Continuing with the example, the computer-readable medium is a portion of flash memory included in the scalable storage device (such as Flash **404F**). In various embodiments, all or any portions of operations performed by any of the elements of Host Visible Storage **110**, Host Invisible Storage **120**, Host Visible Storage **130**, and Host Invisible Storage **140** of Figs. 1A and 1B, as well as units and sub-units within, are managed and/or implemented by one or more state-machines. For example, all or any portions of operations performed by Sub-Request Generation & Sub-Status Accumulation **405** and/or Sub-Request Accepting & Sub-Status Generation **501** are managed and/or implemented by one or more state-machines implemented at least in part by logic gate circuitry.

**[0300]** In various embodiments, all or any portions of a scalable storage system having one or more scalable storage devices (enabled to operate as primary and/or secondary agents) is implemented in one or more drawers and/or shelves of an equipment rack, such as in a datacenter. The scalable storage system is optionally enabled to operate with one or more hosts, such as computing complexes implemented in other drawers and/or shelves of the equipment rack.

#### EXAMPLE IMPLEMENTATION TECHNIQUES

**[0301]** In some embodiments, various combinations of all or portions of operations performed by or blocks included in a scalable storage device (such as Storage Device **110.A** of Fig. 4 or Storage Device **120.A** of Fig. 5), and portions of a processor, microprocessor, system-on-a-chip, application-specific-integrated-circuit, hardware accelerator, or other circuitry providing all or portions of the aforementioned operations or blocks, are specified by a specification compatible with processing by a computer system. The specification is in accordance with various descriptions, such as hardware description languages, circuit descriptions, netlist descriptions, mask descriptions, or layout descriptions. Example descriptions include: Verilog, VHDL, SPICE, SPICE variants such as PSpice, IBIS, LEF, DEF, GDS-II, OASIS, or other descriptions. In various embodiments, the processing includes any combination of interpretation, compilation, simulation, and synthesis to produce, to verify, or to specify logic and/or circuitry suitable for inclusion on one or more integrated circuits. Each integrated circuit, according to various embodiments, is designable and/or manufacturable according to a variety of techniques. The techniques include a programmable technique (such as a field or mask programmable gate array integrated circuit), a semi-custom technique (such as a wholly or partially cell-based integrated circuit), and a full-custom technique (such as an

1 integrated circuit that is substantially specialized), any combination thereof, or any other  
2 technique compatible with design and/or manufacturing of integrated circuits.

3  
4 **[0302]** In some embodiments, various combinations of all or portions of operations as  
5 described by a computer readable medium having a set of instructions stored therein, are  
6 performed by execution and/or interpretation of one or more program instructions, by  
7 interpretation and/or compiling of one or more source and/or script language statements, or by  
8 execution of binary instructions produced by compiling, translating, and/or interpreting  
9 information expressed in programming and/or scripting language statements. The statements are  
10 compatible with any standard programming or scripting language (such as C, C++, Fortran,  
11 Pascal, Ada, Java, VBscript, and Shell). One or more of the program instructions, the language  
12 statements, or the binary instructions, are optionally stored on one or more computer readable  
13 storage medium elements. In various embodiments some, all, or various portions of the program  
14 instructions are realized as one or more functions, routines, sub-routines, in-line routines,  
15 procedures, macros, or portions thereof.

## CONCLUSION

**[0303]** Certain choices have been made in the description merely for convenience in preparing the text and drawings and unless there is an indication to the contrary the choices should not be construed per se as conveying additional information regarding structure or operation of the embodiments described. Examples of the choices include: the particular organization or assignment of the designations used for the figure numbering and the particular organization or assignment of the element identifiers (the callouts or numerical designators, e.g.) used to identify and reference the features and elements of the embodiments.

**[0304]** The words “includes” or “including” are specifically intended to be construed as abstractions describing logical sets of open-ended scope and are not meant to convey physical containment unless explicitly followed by the word “within.”

**[0305]** Although the foregoing embodiments have been described in some detail for purposes of clarity of description and understanding, the invention is not limited to the details provided. There are many embodiments of the invention. The disclosed embodiments are exemplary and not restrictive.

**[0306]** It will be understood that many variations in construction, arrangement, and use are possible consistent with the description, and are within the scope of the claims of the issued patent. For example, interconnect and function-unit bit-widths, clock speeds, and the type of technology used are variable according to various embodiments in each component block. The names given to interconnect and logic are merely exemplary, and should not be construed as limiting the concepts described. The order and arrangement of flowchart and flow diagram process, action, and function elements are variable according to various embodiments. Also, unless specifically stated to the contrary, value ranges specified, maximum and minimum values used, or other particular specifications (such as flash memory technology types; and the number of entries or stages in registers and buffers), are merely those of the described embodiments, are expected to track improvements and changes in implementation technology, and should not be construed as limitations.

**[0307]** Functionally equivalent techniques known in the art are employable instead of those described to implement various components, sub-systems, operations, functions, routines, sub-routines, in-line routines, procedures, macros, or portions thereof. It is also understood that many functional aspects of embodiments are realizable selectively in either hardware (i.e.,

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

11

12 **[0308]** The embodiments have been described with detail and environmental context  
13 well beyond that required for a minimal implementation of many aspects of the embodiments  
14 described. Those of ordinary skill in the art will recognize that some embodiments omit  
15 disclosed components or features without altering the basic cooperation among the remaining  
16 elements. It is thus understood that much of the details disclosed are not required to implement  
17 various aspects of the embodiments described. To the extent that the remaining elements are  
18 distinguishable from the prior art, components and features that are omitted are not limiting on  
19 the concepts described herein.

20

21 **[0309]** All such variations in design are insubstantial changes over the teachings  
22 conveyed by the described embodiments. It is also understood that the embodiments described  
23 herein have broad applicability to other computing and networking applications, and are not  
24 limited to the particular application or industry of the described embodiments. The invention is  
25 thus to be construed as including all possible modifications and variations encompassed within  
26 the scope of the claims of the issued patent.



## WHAT IS CLAIMED IS:

1     1. A system comprising:  
2         a storage subsystem having physical components comprising a plurality of physical  
3             storage devices and a physical switch-portion having a plurality of ports;  
4         a host coupled via a dedicated point-to-point link to a dedicated port of the plurality of  
5             ports;  
6         wherein each physical storage device comprises:  
7             at least one range of storage,  
8             at least one port enabled to be coupled via a respective point-to-point link to a  
9             respective port of the plurality of ports, and  
10            agent logic enabling the physical storage device to operate as at least one agent  
11            of one or more primary agents and one or more secondary agents of the  
12            storage subsystem, each primary agent enabled to accept host-initiated  
13            storage access requests, to generate sub-requests, and to accumulate  
14            sub-statuses, each secondary agent enabled to accept at least one of the  
15            sub-requests and to generate at least one of the sub-statuses; and  
16         wherein the storage subsystem is enabled to operate one or more logical storage devices,  
17             each logical storage device enabling the host to access via a single logical  
18             interface an aggregation of the storage ranges distributed over the storage  
19             devices corresponding to one primary agent of the primary agents and at least  
20             one secondary agent of the secondary agents.

1     2. The system of claim 1, further comprising:  
2         wherein the agent logic comprises configurability logic enabling the physical storage  
3             device to operate at least sometimes as at least one of the primary agents and at  
4             least sometimes as at least one of the secondary agents.

1     3. The system of claim 1, further comprising:  
2         wherein the agent logic comprises concurrency logic enabling the physical storage  
3             device to operate concurrently as at least one of the primary agents and at least  
4             one of the secondary agents.

- 1     4. The system of claim 1, further comprising:  
2             wherein the agent logic comprises dedicated-agent logic enabling the physical storage  
3             device to operate as a dedicated one of the primary agents and the secondary  
4             agents.
- 1     5. The system of claim 1, further comprising:  
2             wherein a first logical storage device of the logical storage devices comprises at least a  
3             first primary agent of the primary agents and one or more of the secondary  
4             agents; and  
5             wherein the physical switch-portion enables host-occluded peer-to-peer communications  
6             between the first primary agent and each of the secondary agents.
- 1     6. The system of claim 1, further comprising:  
2             wherein the physical switch-portion having a plurality of ports is a first physical switch-  
3             portion having a first plurality of ports and the physical components further  
4             comprise a second physical switch-portion having a second plurality of ports;  
5             and  
6             wherein each physical storage device further comprises at least one port enabled to be  
7             coupled via a respective point-to-point link to a respective port of the second  
8             plurality of ports.
- 1     7. The system of claim 6, further comprising:  
2             wherein the second physical switch-portion enables storage redundancy techniques  
3             implemented via peer-to-peer communication with inconsequential impact on  
4             bandwidth available for communications via the first physical switch-portion.
- 1     8. The system of claim 7, further comprising:  
2             wherein the peer-to-peer communication is via the second physical switch-portion and  
3             comprises one or more of control information, untransformed redundancy data,  
4             and transformed redundancy data.
- 1     9. The system of claim 6, further comprising:  
2             wherein a unitary physical switch comprises the first physical switch-portion and the  
3             second physical switch-portion.

1 10. The system of claim 1, further comprising:

2 wherein the single logical interface enables one or more of mirroring, striping, RAID  
3 parity, and fail-over, to be implemented via host-occluded peer-to-peer  
4 communications between the primary and secondary agents.

1 11. The system of claim 1, further comprising:

2 wherein the primary agents comprise a first primary agent and a second primary agent,  
3 the secondary agents comprise a first secondary agent and a second secondary  
4 agent, the physical storage devices comprise a first physical storage device and a  
5 second physical storage device, and the logical storage devices comprise a first  
6 logical storage device and a second logical storage device;

7 wherein with respect to a first subset of the requests, the first physical storage device  
8 operates as the first primary agent and as the first secondary agent;

9 wherein with respect to a second subset of the requests, the second physical storage  
10 device operates as the second primary agent and as the second secondary agent;  
11 and

12 wherein the first logical storage device comprises the first primary agent and the second  
13 secondary agent, the second logical storage device comprises the second  
14 primary agent and the first secondary agent, and the first logical storage device  
15 operates concurrently with the second logical storage device.

1 12. The system of claim 11, further comprising:

2 wherein the first subset of the requests are from a first host request queue and the second  
3 subset of the requests are from a second host request queue.

1 13. The system of claim 11, further comprising:

2 wherein the host is a first host and the first subset of the requests are from the first host,  
3 and the second subset of the requests are from a second host.

1 14. The system of claim 11, further comprising:

2 wherein the first subset of the requests and the second subset of the requests are from the  
3 same host request queue.

1     15. The system of claim 1, further comprising:  
2             wherein a first logical storage device of the logical storage devices comprises at least a  
3             first primary agent of the primary agents and at least a first secondary agent of  
4             the secondary agents;  
5             wherein the aggregation of the storage ranges of the agents of the first logical storage  
6             device is a first logical storage sub-space; and  
7             wherein the agent logic comprises storage-recovery logic enabling the first primary  
8             agent to identify free storage available to be provided by one or more of the  
9             agents and, subsequent to a determination that a particular portion of the first  
10            logical storage sub-space is no longer being provided by one of the agents, to  
11            allocate from the free storage to restore the previously provided particular  
12            portion.

1     16. The system of claim 1, further comprising:  
2             wherein a first logical storage device of the logical storage devices comprises at least a  
3             first primary agent of the primary agents and at least a first secondary agent of  
4             the secondary agents, the physical storage devices comprise a first physical  
5             storage device and a second physical storage device, and wherein the first  
6             physical device initially operates as the first primary agent;  
7             a monitoring agent enabled to determine that the first primary agent has failed and to  
8             transmit a primary-agent-replacement request to a selected one of the other  
9             agents; and  
10            wherein the agent logic of the second physical device comprises primary-agent-  
11            replacement logic enabling the second physical storage device to operate as a  
12            replacement for the first primary agent in response to receiving the primary-  
13            agent-replacement request from the monitoring agent.

- 1     17. A method comprising:  
2         operating as a primary agent a first storage device and operating as one or more  
3         secondary agents at least a second storage device, each storage device of the  
4         primary and secondary agents having at least one respective storage range;  
5         the primary agent providing a logical interface enabling the host to access as a unitary  
6         logical device an aggregation of the storage ranges distributed over the storage  
7         devices of the primary and secondary agents, primary agent communications  
8         with the host via the logical interface comprising storage access requests  
9         accepted by the primary agent from the host and overall status sent by the  
10        primary agent to the host;  
11        the primary agent forwarding each request accepted as one or more sub-requests to all or  
12        any portions of the secondary agents, accepting sub-statuses associated with the  
13        sub-requests from the secondary agents that the sub-requests were forwarded to,  
14        and formulating each overall status sent based at least on the sub-statuses; and  
15        the all or any portions of the secondary agents accepting the sub-requests, translating  
16        host-context addressing information of the sub-requests to local-context  
17        addressing information, accessing local storage based at least in part on the  
18        local-context addressing information, and providing the sub-statuses to the  
19        primary agent.
- 1     18. The method of claim 17, further comprising:  
2         the primary agent determining an aggregated storage capacity of the unitary logical  
3         device based at least in part on the storage ranges of the primary and secondary  
4         agents;  
5         in response to the host querying the unitary logical device regarding storage capacity,  
6         the primary agent reporting the aggregated storage capacity; and  
7         in response to the host querying the secondary agents regarding storage capacity, each of  
8         the secondary agents reporting a respective storage capacity that excludes  
9         representation of any storage range of the secondary agent that is represented in  
10        the aggregated storage capacity.

- 1     19. The method of claim 17, further comprising:  
2             in response to the host querying a particular one of the secondary agents regarding  
3             storage capacity, the particular secondary agent reporting back a particular  
4             storage capacity that excludes representation of those portions of storage  
5             implemented by the particular secondary agent that are represented in an  
6             aggregated storage capacity determined by the primary agent.
- 1     20. The method of claim 17, further comprising:  
2             performing at least some of the primary agent communications with the host in  
3             accordance with a host-interface protocol.
- 1     21. The method of claim 17, further comprising:  
2             performing at least some of the primary agent communications with the host and at least  
3             some primary agent communications with at least one of the secondary agents  
4             via sharing a same physical link; and  
5             wherein the primary agent communications with at least one of the secondary agents  
6             comprise at least some of the sub-requests and at least some of the sub-statuses  
7             exchanged between the primary agent and the secondary agents.
- 1     22. The method of claim 17, further comprising:  
2             performing at least some primary agent communications with the secondary agents in  
3             accordance with a peer-to-peer protocol; and  
4             wherein the primary agent communications with the secondary agents comprise at least  
5             the sub-requests and the sub-statuses exchanged between the primary agent and  
6             the secondary agents.
- 1     23. The method of claim 22, wherein:  
2             the primary agent communications with the secondary agents further comprise  
3             exchanging redundancy information via the peer-to-peer protocol.

1     24. An apparatus comprising:  
2             means for operating as a primary agent a first storage device and operating as one or  
3             more secondary agents at least a second storage device, each storage device of  
4             the primary and secondary agents having at least one respective storage range;  
5             means for the primary agent providing a logical interface enabling the host to access as a  
6             unitary logical device an aggregation of the storage ranges distributed over the  
7             storage devices of the primary and secondary agents, primary agent  
8             communications with the host via the logical interface comprising storage  
9             access requests accepted by the primary agent from the host and overall status  
10            sent by the primary agent to the host;  
11            means for the primary agent forwarding each request accepted as one or more sub-  
12            requests to all or any portions of the secondary agents, accepting sub-statuses  
13            associated with the sub-requests from the secondary agents that the sub-requests  
14            were forwarded to, and formulating each overall status sent based at least on the  
15            sub-statuses; and  
16            means for the all or any portions of the secondary agents accepting the sub-requests,  
17            translating host-context addressing information of the sub-requests to local-  
18            context addressing information, accessing local storage based at least in part on  
19            the local-context addressing information, and providing the sub-statuses to the  
20            primary agent.

1     25. The apparatus of claim 24, further comprising:  
2             means for the primary agent determining an aggregated storage capacity of the unitary  
3             logical device based at least in part on the storage ranges of the primary and  
4             secondary agents;  
5             means for in response to the host querying the unitary logical device regarding storage  
6             capacity, the primary agent reporting the aggregated storage capacity; and  
7             means for in response to the host querying the secondary agents regarding storage  
8             capacity, each of the secondary agents reporting a respective storage capacity  
9             that excludes representation of any storage range of the secondary agent that is  
10            represented in the aggregated storage capacity.

1 26. The apparatus of claim 24, further comprising:  
2 means for in response to the host querying a particular one of the secondary agents  
3 regarding storage capacity, the particular secondary agent reporting back a  
4 particular storage capacity that excludes representation of those portions of  
5 storage implemented by the particular secondary agent that are represented in an  
6 aggregated storage capacity determined by the primary agent.

1 27. The apparatus of claim 24, further comprising:  
2 means for performing at least some of the primary agent communications with the host  
3 in accordance with a host-interface protocol.

1 28. The apparatus of claim 24, further comprising:  
2 means for performing at least some of the primary agent communications with the host  
3 and at least some primary agent communications with at least one of the  
4 secondary agents via sharing a same physical link; and  
5 wherein the primary agent communications with at least one of the secondary agents  
6 comprise at least some of the sub-requests and at least some of the sub-statuses  
7 exchanged between the primary agent and the secondary agents.

1 29. The apparatus of claim 24, further comprising:  
2 means for performing at least some primary agent communications with the secondary  
3 agents in accordance with a peer-to-peer protocol; and  
4 wherein the primary agent communications with the secondary agents comprise at least  
5 the sub-requests and the sub-statuses exchanged between the primary agent and  
6 the secondary agents.

1 30. The apparatus of claim 29, wherein:  
2 the primary agent communications with the secondary agents further comprise  
3 redundancy information exchanged via the peer-to-peer protocol.



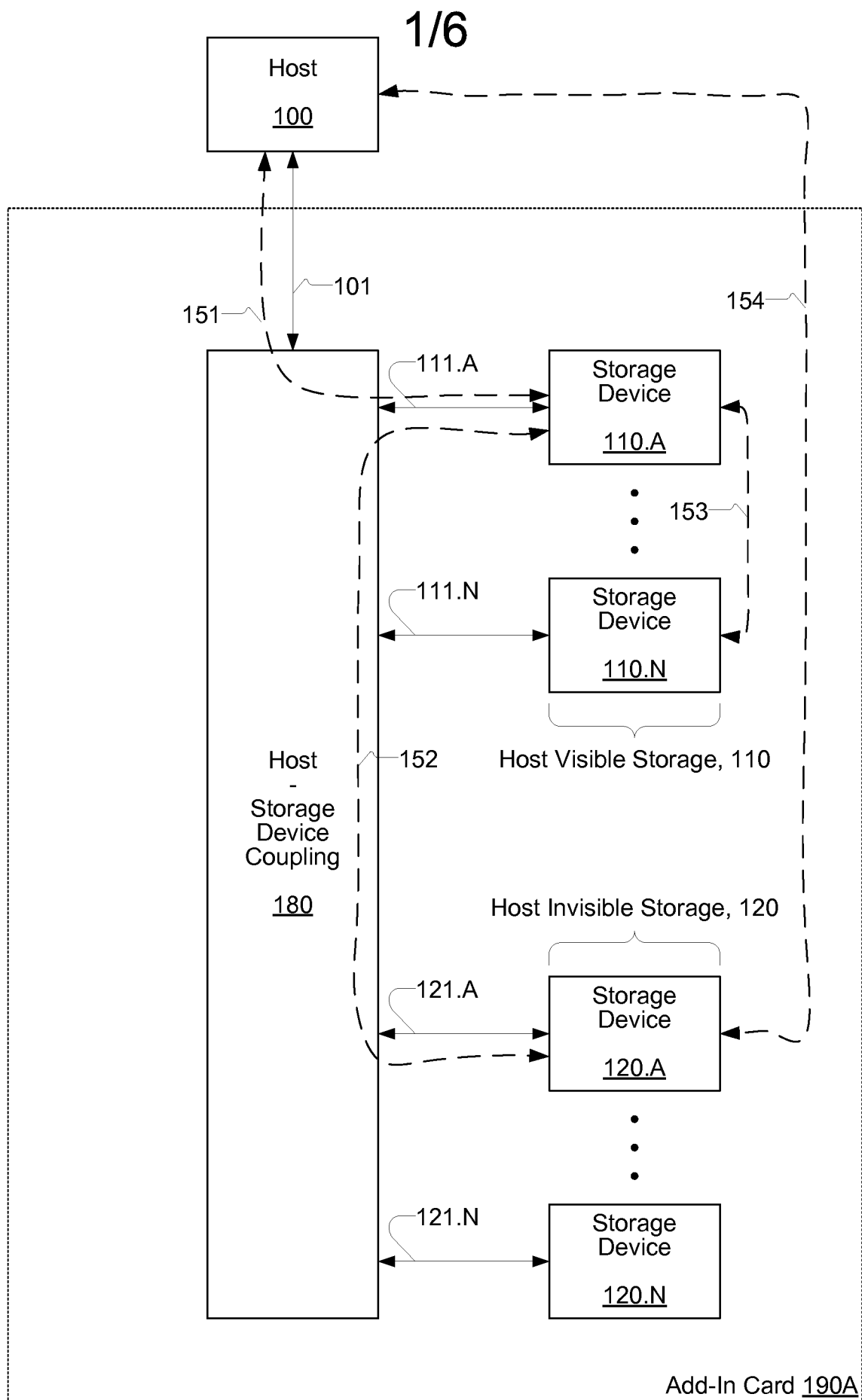


Fig. 1A

2/6

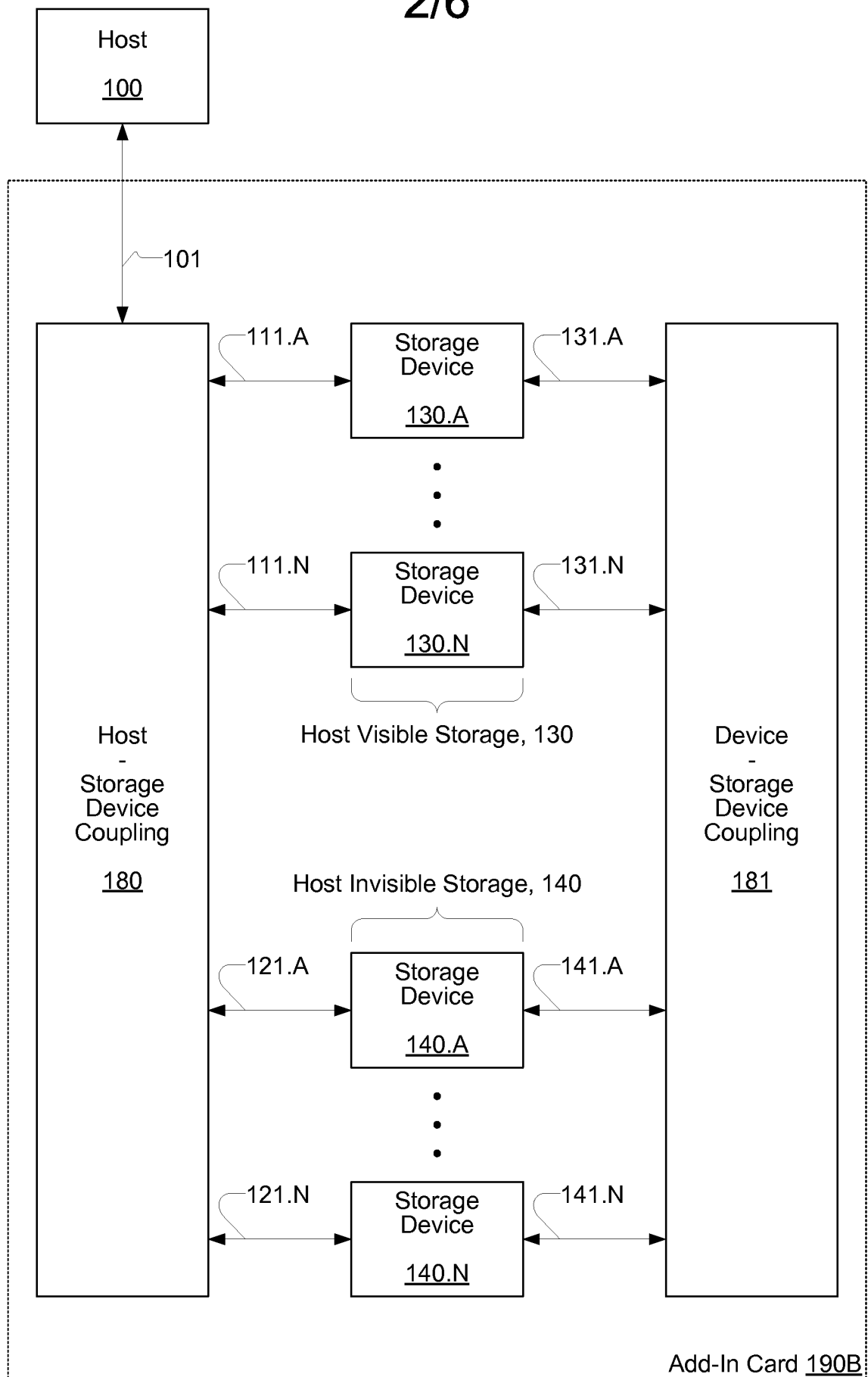


Fig. 1B

3/6

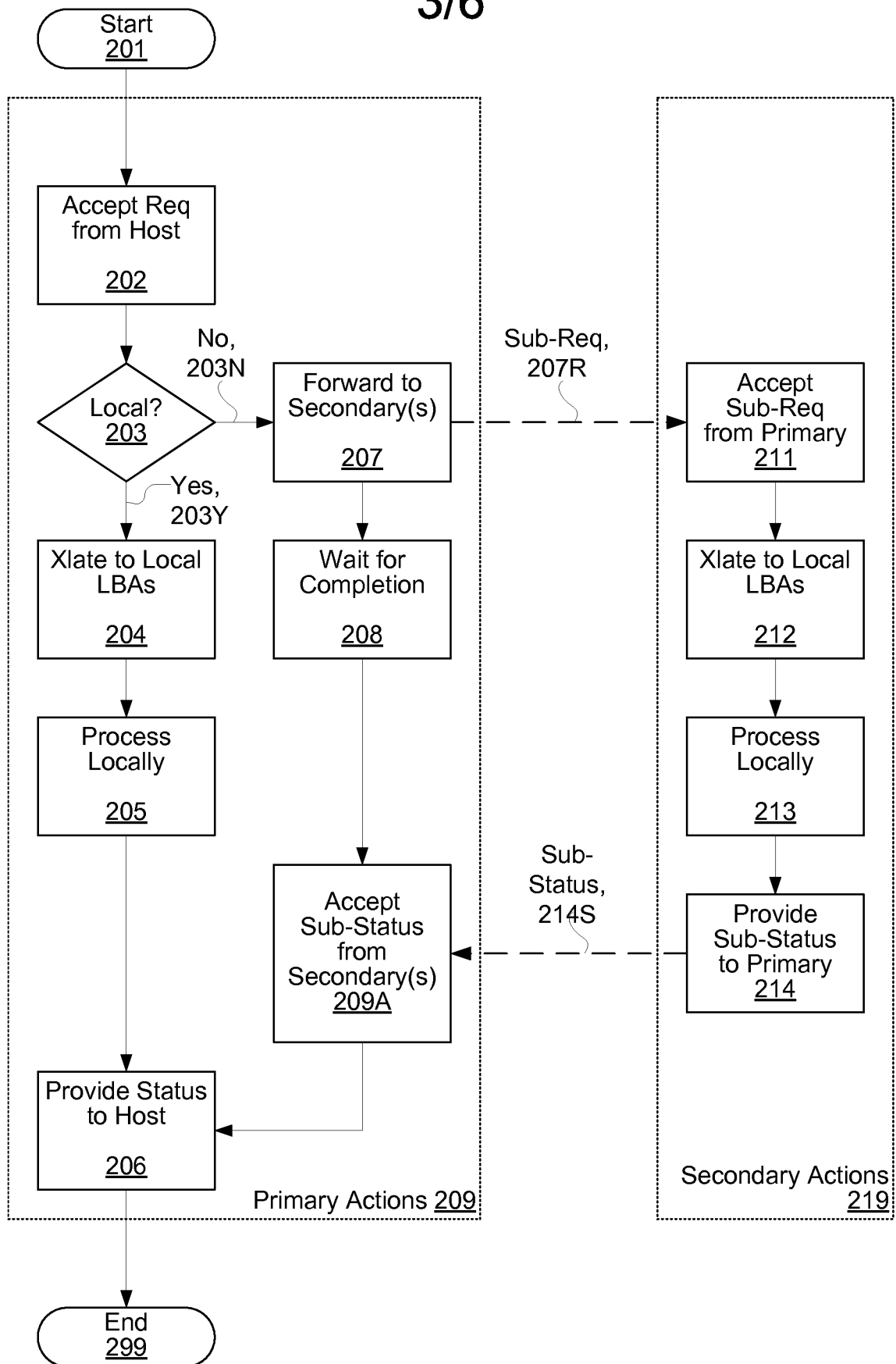


Fig. 2

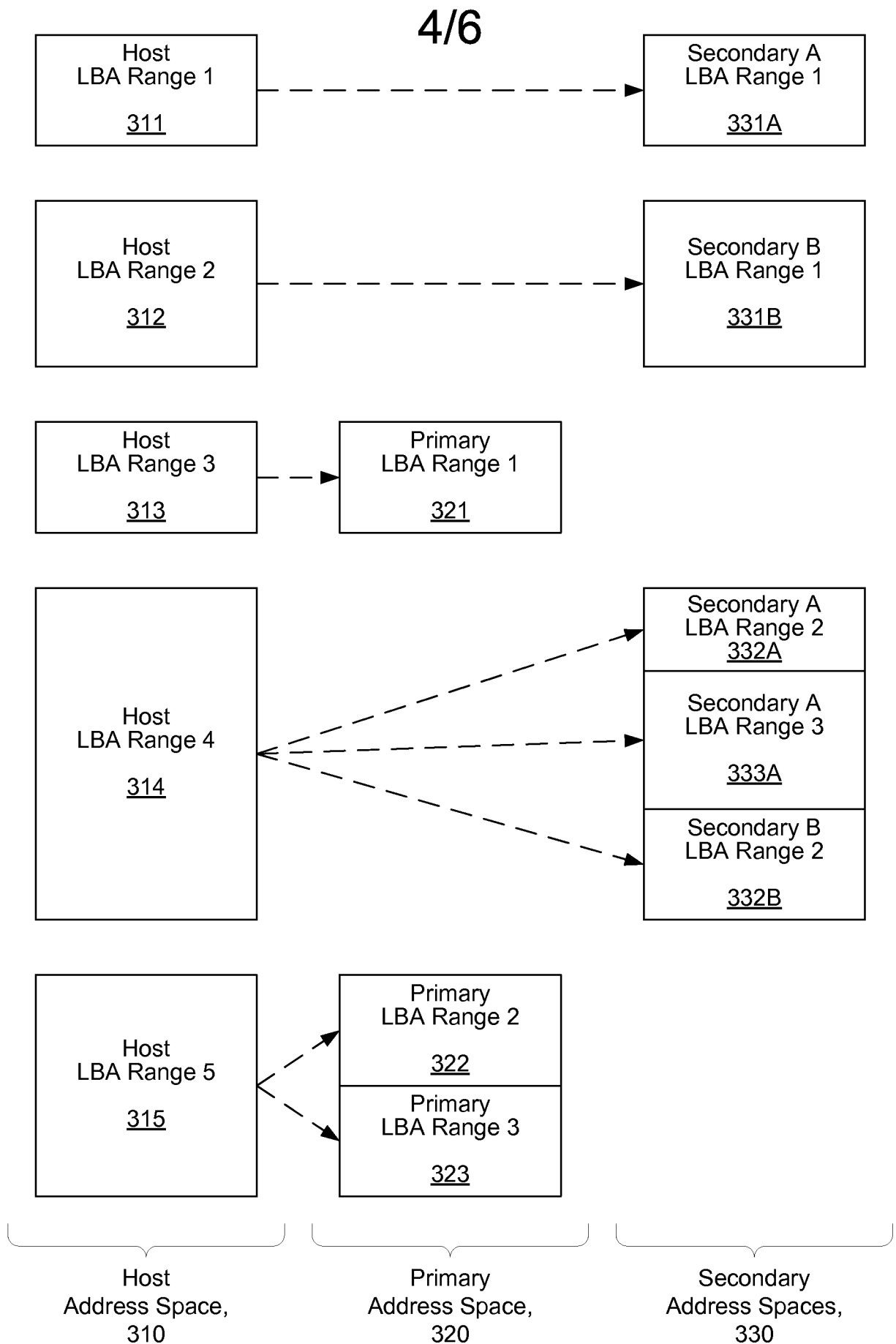


Fig. 3

5/6

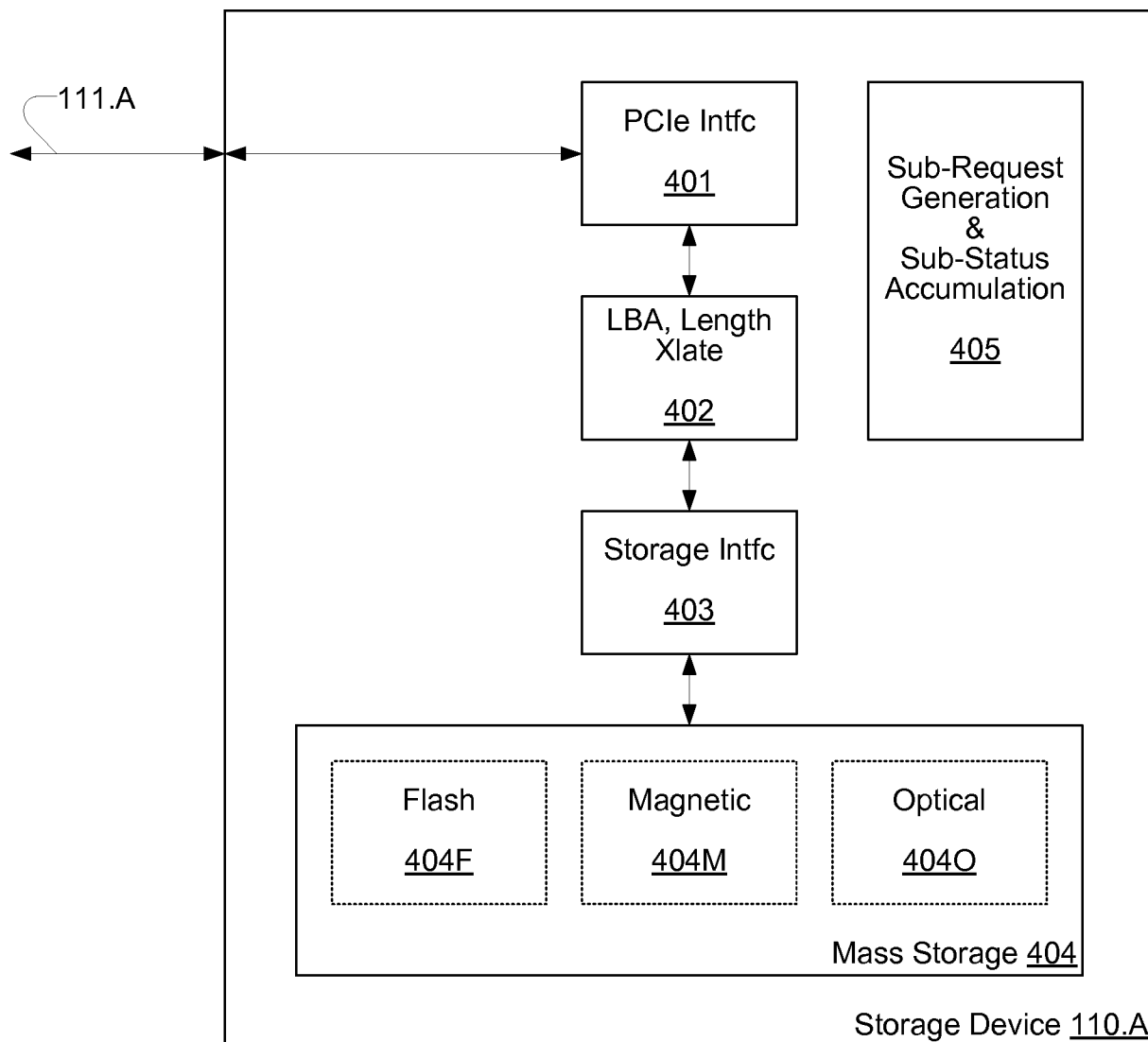


Fig. 4

6/6

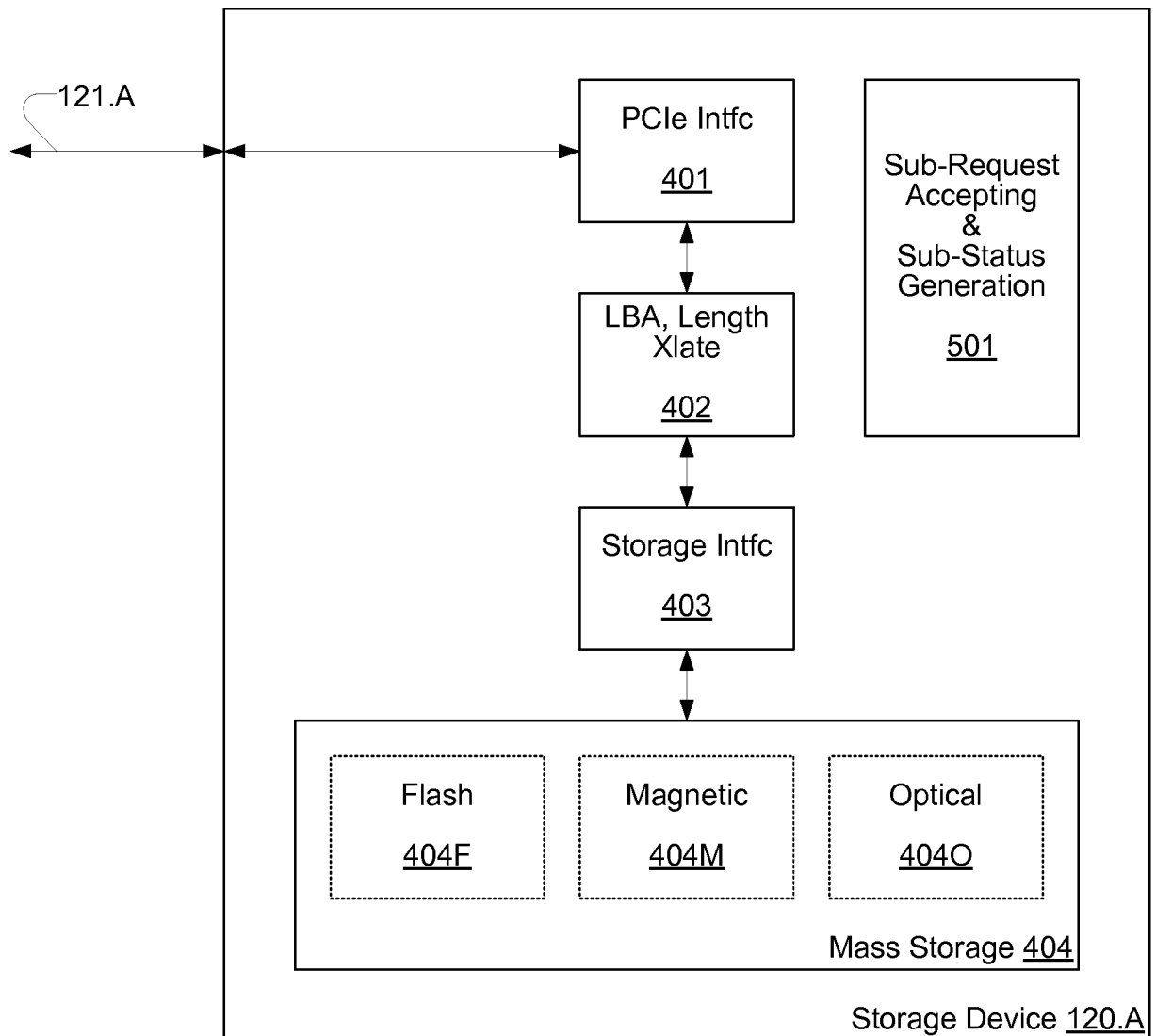


Fig. 5