

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2013-97788

(P2013-97788A)

(43) 公開日 平成25年5月20日(2013.5.20)

(51) Int.Cl.

G06F 13/10 (2006.01)

F I

G06F 13/10 340A

テーマコード (参考)

審査請求 未請求 請求項の数 19 O L (全 11 頁)

(21) 出願番号 特願2012-229045 (P2012-229045)  
 (22) 出願日 平成24年10月16日 (2012.10.16)  
 (31) 優先権主張番号 13/289,617  
 (32) 優先日 平成23年11月4日 (2011.11.4)  
 (33) 優先権主張国 米国 (US)

(特許庁注：以下のものは登録商標)

1. イーサネット

(71) 出願人 591007686  
 エルエスアイ コーポレーション  
 アメリカ合衆国カリフォルニア州95035, ミルピタス, バーバー・レーン 1621  
 (74) 代理人 100140109  
 弁理士 小野 新次郎  
 (74) 代理人 100075270  
 弁理士 小林 泰  
 (74) 代理人 100080137  
 弁理士 千葉 昭男  
 (74) 代理人 100096013  
 弁理士 富田 博行  
 (74) 代理人 100096068  
 弁理士 大塚 住江

最終頁に続く

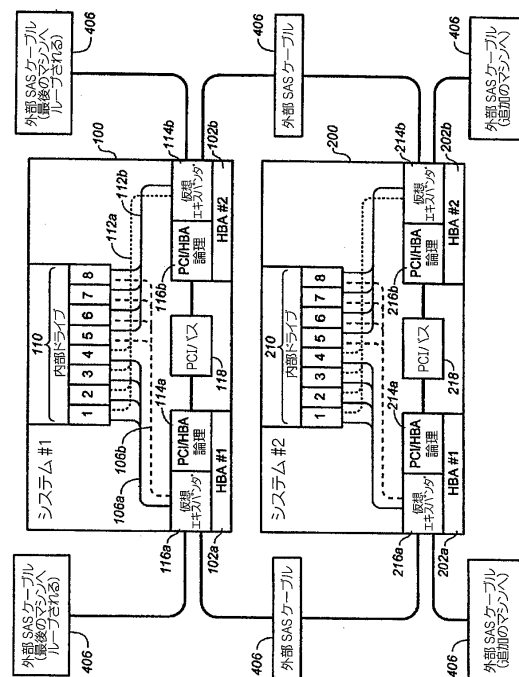
(54) 【発明の名称】 仮想SASエキスパンダを介して共有されるサーバ直接接続のストレージシステム

(57) 【要約】 (修正有)

【課題】複数のサーバそれぞれに備えられたストレージ装置を簡単な構成で安価に共有化する。

【解決手段】データストレージシステムは複数のサーバ100、200を備えている。サーバ100は、データを記憶する複数のストレージディスク110とホストバスアダプタ102aとを備え、該ホストバスアダプタは、仮想エキスパンダ116aと論理コンポーネント114aとを提供するプロセッサを備えている。サーバ200も同様に構成されている。ホストバスアダプタ102aは、SAS接続406を介してサーバ200のホストバスアダプタ202aに接続され、複数のストレージディスク110、210のそれぞれは、サーバ100、200それぞれによりアクセス可能である。同様に、サーバ200に更に他のサーバを接続することも可能である。

【選択図】 図4



**【特許請求の範囲】****【請求項 1】**

データストレージシステムであって、

第 1 のサーバであって

データを記憶するように構成された第 1 の複数のストレージディスクと、

第 1 の仮想エキスパンダと第 1 の論理コンポーネントとを提供するように構成された第 1 のプロセッサを含んだ第 1 のホストバスアダプタと

を備えた第 1 のサーバと、

第 2 のサーバであって

データを記憶するように構成された第 2 の複数のストレージディスクと、

第 2 の仮想エキスパンダと第 2 の論理コンポーネントとを提供するように構成された第 2 のプロセッサを含んだ第 2 のホストバスアダプタと

を備えた第 2 のサーバと

を備え、

第 1 のサーバの第 1 のホストバスアダプタは、シリアル接続のスモールコンピュータシステムインターフェース (S A S) 接続を介して、第 2 のサーバの第 2 のホストバスアダプタへ接続され、第 1 の複数のストレージディスク及び第 2 の複数のストレージディスクのそれぞれは、第 1 のサーバ及び第 2 のサーバのそれぞれによりアクセス可能であることを特徴とするデータストレージシステム。

10

**【請求項 2】**

請求項 1 記載のシステムにおいて、S A S 接続は、S A S ケーブルであることを特徴とするシステム。

20

**【請求項 3】**

請求項 1 又は 2 記載のシステムにおいて、第 1 のサーバの第 1 のホストバスアダプタの第 1 の仮想エキスパンダは、S A S 接続を介して、第 2 のサーバの第 2 のホストバスアダプタの第 2 の仮想エキスパンダへ接続されることを特徴とするシステム。

**【請求項 4】**

請求項 1 ~ 3 いずれかに記載のシステムにおいて、第 1 のサーバ及び第 2 のサーバのそれぞれはさらに、別のホストバスアダプタを備えることを特徴とするシステム。

**【請求項 5】**

請求項 4 記載のシステムにおいて、第 1 のサーバの別のホストバスアダプタは、第 2 のサーバの別のホストバスアダプタへ接続されることを特徴とするシステム。

30

**【請求項 6】**

請求項 4 又は 5 記載のシステムにおいて、該システムはさらに、第 1 のサーバの第 1 のホストバスアダプタと第 1 のサーバの別のホストバスアダプタとのそれぞれに接続されるバスを備えることを特徴とするシステム。

**【請求項 7】**

請求項 1 ~ 6 いずれかに記載のシステムにおいて、該システムはさらに第 3 のサーバを備え、該第 3 のサーバは、

データを記憶するように構成された第 3 の複数のストレージディスクと、

第 3 の仮想エキスパンダと第 3 の論理コンポーネントとを提供するように構成された第 3 のプロセッサを含んだ第 3 のホストバスアダプタと

を備えており、

第 1 のサーバは第 3 のサーバへ接続され、かつ第 2 のサーバは第 3 のサーバへ接続されることを特徴とするシステム。

40

**【請求項 8】**

請求項 7 記載のシステムにおいて、該システムは、第 1 のサーバと第 2 のサーバとの間、第 1 のサーバと第 3 のサーバとの間、及び第 2 のサーバと第 3 のサーバとの間のうち 1 つに、フェールオーバー接続を含むことを特徴とするシステム。

**【請求項 9】**

50

請求項 1 ~ 8 いずれかに記載のシステムにおいて、第 1 の複数のストレージディスクと第 2 の複数のストレージディスクとは、リダンダント・アレイ・オブ・インディペンデント・ディスク ( R A I D ) 構成で構成されることを特徴とするシステム。

【請求項 1 0】

データストレージシステムであって、

第 1 のサーバであって

データを記憶するように構成された第 1 の複数のストレージディスクと、

第 1 のマルチコアプロセッサを備えた第 1 のホストバスアダプタと

を備え、第 1 のマルチコアプロセッサの 1 つのコアが第 1 の仮想エキスパンドを提供するように構成された第 1 のサーバと、

第 2 のサーバであって

データを記憶するように構成された第 2 の複数のストレージディスクと、

第 2 のマルチコアプロセッサを備えた第 2 のホストバスアダプタと

を備え、第 2 のマルチコアプロセッサの 1 つのコアが第 2 の仮想エキスパンドを提供するように構成された第 2 のサーバと

を備え、

第 1 のサーバの第 1 のホストバスアダプタは、シリアル接続スモールコンピュータシステムインターフェース ( S A S ) 接続を介して、第 2 のサーバの第 2 のホストバスアダプタへ接続され、第 1 の複数のストレージディスク及び第 2 の複数のストレージディスクのそれぞれは、第 1 のサーバ及び第 2 のサーバのそれぞれによりアクセス可能であることを特徴とするデータストレージシステム。

【請求項 1 1】

請求項 1 0 記載のシステムにおいて、S A S 接続は、S A S ケーブルであることを特徴とするシステム。

【請求項 1 2】

請求項 1 0 又は 1 1 記載のシステムにおいて、第 1 のサーバの第 1 のホストバスアダプタの第 1 の仮想エキスパンドは、S A S 接続を介して、第 2 のサーバの第 2 のホストバスアダプタの第 2 の仮想エキスパンドへ接続されることを特徴とするシステム。

【請求項 1 3】

請求項 1 0 ~ 1 2 いずれかに記載のシステムにおいて、第 1 のサーバ及び第 2 のサーバそれぞれはさらに、別のホストバスアダプタを備えることを特徴とするシステム。

【請求項 1 4】

請求項 1 3 記載のシステムにおいて、第 1 のサーバの別のホストバスアダプタは、第 2 のサーバの別のホストバスアダプタへ接続されることを特徴とするシステム。

【請求項 1 5】

請求項 1 3 又は 1 4 記載のシステムにおいて、該システムはさらに、第 1 のサーバの第 1 のホストバスアダプタと第 1 のサーバの 2 番目のホストバスアダプタとへそれぞれ接続されるバスを備えることを特徴とするシステム。

【請求項 1 6】

請求項 1 0 ~ 1 5 いずれかに記載のシステムにおいて、該システムはさらに第 3 のサーバを備え、第 3 のサーバは、

データを記憶するように構成された第 3 の複数のストレージディスクと、

第 3 のマルチコアプロセッサを備えた第 3 のホストバスアダプタと

を備え、

第 3 のマルチコアプロセッサの 1 つのコアは第 3 の仮想エキスパンドを提供するように構成されており、

第 1 のサーバは第 3 のサーバへ接続され、第 2 のサーバは第 3 のサーバへ接続されることを特徴とするシステム。

【請求項 1 7】

請求項 1 6 記載のシステムにおいて、該システムは、第 1 のサーバと第 2 のサーバとの間

10

20

30

40

50

、第 1 のサーバと第 3 のサーバとの間、又は第 2 のサーバと第 3 のサーバとの間のうち 1 つに、フェールオーバー接続を含むことを特徴とするシステム。

【請求項 18】

請求項 10 ~ 17 いずれかに記載のシステムにおいて、第 1 の複数のストレージディスクと第 2 の複数のストレージディスクとは、リダンダント・アレイ・オブ・インディペンデント・ディスク (RAID) 構成で構成されることを特徴とするシステム。

【請求項 19】

少なくとも 4 つのサーバを備えたデータストレージシステムであって、少なくとも 4 つのサーバのそれぞれは、

データを記憶するように構成された複数のストレージディスクと、

10

第 1 の仮想エキスパンドを提供するように構成された第 1 のプロセッサを備えた第 1 のホストバスアダプタと、

第 2 の仮想エキスパンドを提供するように構成された第 2 のプロセッサを備えた第 2 のホストバスアダプタと

を備え、

少なくとも 4 つのサーバのそれぞれが第 1 の接続構成を備え、該第 1 の接続構成は、少なくとも 4 つのサーバのうちの 1 つの第 1 の仮想エキスパンドを、少なくとも 4 つのサーバのうちの別の 2 つのサーバの異なる第 1 の仮想エキスパンドへ接続し、

少なくとも 4 つのサーバのそれぞれが第 2 の接続構成を備え、該第 2 の接続構成は、少なくとも 4 つのサーバのうちの 1 つの第 2 の仮想エキスパンドを、少なくとも 4 つのサーバのうちの別の 2 つのサーバの異なる第 2 の仮想エキスパンドへ接続し、

20

少なくとも 4 つのサーバのうち少なくとも 1 つのサーバの第 1 の接続構成は、どのサーバが第 1 の接続構成及び第 2 の接続構成に関連付けられるかにより、少なくとも 4 つのサーバのうち少なくとも 1 つのサーバの第 2 の接続構成とは相違している

ことを特徴とするデータストレージシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データストレージシステム (データ記憶システム) の技術分野に関し、より詳細には、仮想 SAS (シリアル接続 SCSI (スモールコンピュータシステムインターフェース)) エキスパンドを介して共有されるサーバ直接接続ストレージに関する。

30

【背景技術】

【0002】

クラウドコンピューティング技術は、オンデマンドのネットワークアクセスを可能にするモデルを提供することにより増加傾向にあり、モデルは設定可能なコンピューティング資源 (例えば、ネットワーク、サーバ、ストレージ (記憶装置)、アプリケーション、及びサービス) の共有プールへアクセスし、最小限の管理努力やサービス提供者とのやり取りで迅速な提供及び回収が可能なものである。クラウドコンピューティングは冗長性のためにクラスタリング (集団化) を使用するのが一般的であり、これはさまざまなストレージ構成により達成される。そのうち 4 つを本明細書に示すが、それぞれが問題のある特徴を含んでいる。

40

【0003】

4 つの構成は、(1) 各ノードが共有 SAN (ストレージエリアネットワーク) ファブリックへ接続されており、低レイテンシーのブロックインターフェースをストレージへ提供するもの、(2) 各ノードがイーサネットネットワークへ接続されており、共有ストレージへのファイルアクセスを使用するもの、(3) 外部 JBOD (just a bunch of disks「単純ディスク束」)、及び(4) 直接接続ドライブ (内部ドライブ) である。

【発明の開示】

【発明が解決しようとする課題】

【0004】

50

構成(1)及び(2)は、ノードを共有ストレージに接続してクラスタを形成するために、ファイバ又はイーサネットスイッチ等の追加的な外部要素を必要とすることがある。このような外部要素は、単一障害点をもたらすので望ましくない。結果として、冗長性コンポーネントは、高い利用可能性の構成を提供することが必要になり、システムに追加費用がかかることになる。

【0005】

構成(3)は費用効率がよいのが通常であるが、この構成ではクラスタ内のノードの量がJ B O D上のコネクタの数に限定されるので、極めて制限的でスケラビリティが限定される。さらに、構成(1)～(3)では、ストレージシステムを外部筐体に配置する必要があることが通常であり、追加的な電力、スペース、及び維持費用がかかる。

10

【0006】

構成(4)は、一般に経済的であるが、接続ドライブのための共有ストレージがないため、高い利用可能性のクラスタリングを提供することができない。このように、これらの構成の費用及び複雑性には問題があり、高い利用可能性のクラスタリングのためのストレージ要求(例えば、冗長性及び共有アクセス)に対して、望ましい解決策を提供するものではない。

【課題を解決するための手段】

【0007】

本発明の一実施例では、データストレージシステムは第1のサーバを備え、第1のサーバは、データを記憶するように構成された第1の複数のストレージディスクと第1のホストバスアダプタとを備え、第1のホストバスアダプタは、第1の仮想エキスパンドと第1の論理コンポーネントとを提供するように構成された第1のプロセッサを備えたものであり、このデータストレージシステムはさらに、第2のサーバを備え、第2のサーバは、データを記憶するように構成された第2の複数のストレージディスクと第2のホストバスアダプタとを備え、第2のホストバスアダプタは、第2の仮想エキスパンドと第2の論理コンポーネントとを提供するように構成された第2のプロセッサを備えたものであり、第1のサーバの第1のホストバスアダプタは、S A S接続を介して第2のサーバの第2のホストバスアダプタへ接続され、第1の複数のストレージディスク及び第2の複数のストレージディスクのそれぞれは、第1のサーバ及び第2のサーバのそれぞれによりアクセス可能である。

20

30

【0008】

本発明の別の実施例では、データストレージシステムは、第1のサーバを備え、第1のサーバは、データを記憶するように構成された第1の複数のストレージディスクと、第1のマルチコアプロセッサを備えた第1のホストバスアダプタとを備え、第1のマルチコアプロセッサの1つのコアは第1の仮想エキスパンドを提供するように構成されたものであり、このデータストレージシステムはさらに、第2のサーバを備え、第2のサーバは、データを記憶するように構成された第2の複数のストレージディスクと、第2のマルチコアプロセッサを備えた第2のホストバスアダプタとを備え、第2のマルチコアプロセッサの1つのコアは第2の仮想エキスパンドを提供するように構成されたものであり、第1のサーバの第1のホストバスアダプタは、S A S接続を介して第2のサーバの第2のホストバスアダプタへ接続され、第1の複数のストレージディスク及び第2の複数のストレージディスクのそれぞれは、第1のサーバ及び第2のサーバのそれぞれによりアクセス可能である。

40

【0009】

さらに別の本発明の実施例では、データストレージシステムは少なくとも4つのサーバを備え、少なくとも4つのサーバのそれぞれが、データを記憶するように構成された複数のストレージディスクと、第1の仮想エキスパンドを提供するように構成された第1のプロセッサを備えた第1のホストバスアダプタと、第2の仮想エキスパンドを提供するように構成された第2のプロセッサを備えた第2のホストバスアダプタとを備え、少なくとも4つのサーバのそれぞれが第1の接続形態を備え、第1の接続形態は、少なくとも4つの

50

サーバのうちの1つの第1の仮想エキスパンドを、少なくとも4つのサーバのうち別の2つのサーバの、違う第1の仮想エキスパンドへ接続するものであり、少なくとも4つのサーバのそれぞれは、第2の接続形態を備え、第2の接続形態は、少なくとも4つのサーバのうちの1つの第2の仮想エキスパンドを、少なくとも4つのサーバのうち別の2つのサーバの、違う第2の仮想エキスパンドへ接続するものであり、少なくとも4つのサーバのうち少なくとも1つのサーバの第1の接続形態は、少なくとも4つのサーバのうち少なくとも1つのサーバの第2の接続形態と異なっており、これは、どのサーバが第1の接続形態及び第2の接続形態に関連付けられているかによって異なっている。

#### 【0010】

上記した概要及び下記の詳細な説明はともに、例示及び説明にすぎず特許請求される本開示を限定するものとは限らない。添付図面は、明細書に含まれその一部を構成するものであり、本開示の実施形態を示し、概要とともに本開示の原理を説明する役割を果たすものである。当業者が添付図面を参照すれば本発明の多くの効果をよりよく理解できるであろう。

#### 【図面の簡単な説明】

#### 【0011】

【図1】サーバの内部レイアウトの概略図である。

【図2】ホストバスアダプタの概略図である。

【図3A】カスケードDAS（直接接続ストレージ）クラスタの構成の概略図である。

【図3B】カスケードDASクラスタの別の構成の図である。

【図4】図3AのカスケードDASクラスタの一部の概略図である。

【図5】カスケードDASクラスタの一実施例の概略図である。

#### 【発明を実施するための形態】

#### 【0012】

以下に、本発明の好適な実施形態を添付図面を参照して詳細に説明する。

本発明は、サーバがノード（例えば、複数サーバ）のクラスタに属することができるようにする実施例を提供し、このノードのクラスタはストレージ（記憶装置）を共有し、スイッチ又は外部ストレージである外部要素を用いることがないものである。SAS技術は、各ノードにおける直接接続ディスクで、各ノードの間の接続で用いられるのが通常であり、それによりカスケードSASTポロジを通じてSAN環境をエミュレートする。近代のコンピューティングサーバは、SASを通じて組み込まれたディスクを含んでおり、あるサーバの内部ストレージを、接続された他のサーバ間で共有することができる。内部ストレージが共有される場合、大量のデータアクセスのために外部ストレージを必要とすることがなくなる。エキスパンドをエミュレートすることができるSAS HBA（ホストバスアダプタ）を使用することで、他の全てのノード及びそれに対応する付属ディスクへの双方向性トラフィックが可能になる。

#### 【0013】

図1は、ノードのクラスタに組み込まれたサーバ100の概略図を示す。サーバ100は、1又は複数のHBA（例えば、SAS HBA）を備えており、図1には2つのHBA、102a及び102bを示す。図2は、HBA102aのコンポーネントの概略図を示す。図示のように、HBA102aは、4つの外部コネクタのペアである104a及び104bと、4つの内部コネクタのペアである106a及び106bとで、合計16の物理層（phy）を備えている。HBA102aはさらに、HBA102aの動作を管理する2コアCPU108等のプロセッサを備えている。図1に示すように、4つの内部コネクタのペア106a及び106bは、HBA102aを、サーバ100上でストレージとして使用可能な複数のディスク110へ接続する。同様に、HBA102bは、HBA102bをサーバ100上の複数のディスク110へ接続するコネクタ112a及び112bを備えている。

#### 【0014】

HBA102a及びHBA102bの外部コネクタ（例えば104a及び104b）は

10

20

30

40

50

、サーバ 100 をクラスタの一部である他のサーバへ接続するように機能する。各サーバは、クラスタ内の他のサーバへ接続するための少なくとも 1 つの HBA を含み、1 サーバにつき 2 以上の HBA がある場合は冗長性が可能になる。例えば、各サーバ/ノードが、冗長性のために 2 つの他のノードへの SAS 接続（各サーバ/ノードの HBA を介して）を備えていてもよい。

#### 【0015】

図 3 A は、カスケード接続された DAS クラスタのための構成の概略図を示す。この構成は、5 つのサーバ/ノード 100、200、300、400、及び 500 を備え、サーバ/ノード 100 は最初のノードとし、サーバ/ノード 500 は最後のノードとされる。サーバ/ノード 100 は、コネクタ 104 a 及び 104 b を介してサーバ/ノード 200 へリンクされる。サーバ/ノード 200 はコネクタ 204 a 及び 204 b を介してサーバ/ノード 300 へリンクされる。サーバ/ノード 300 はコネクタ 304 a 及び 304 b を介してサーバ/ノード 400 へリンクされる。サーバ/ノード 400 はコネクタ 404 a 及び 404 b を介してサーバ/ノード 500 へリンクされる。最初のノード 100 及び最後のノード 500 も互いに接続されてもよいが、この接続はループ（例えば、無効な SAS トポロジ）を防止する等のために、無効化される。

図 3 A に示すように、サーバ/ノード 100 はコネクタ 504 a 及び 504 b を介してサーバ/ノード 500 へ接続されるが、コネクタ 504 a 及び 504 b はクラスタ内のノードが使用不可能になるまでディスエーブル（動作不能）状態にある。ノード又は接続がディスエーブル状態（例えば、ノード障害）になると、最初のノード及び最後のノードの間の無効化された接続（例えば、コネクタ 504 a 及び 504 b）はファームウェアにより直ちにイネーブル（有効化）状態にされて、使用可能な全てのノードへのアクセスが中断されないようにする。システムの各サーバ/ノードは、全てのノードへアクセス可能な、複数のディスク 110 等のローカル SAS（又は SATA（シリアル ATA タッチメント））ストレージを備えている。本明細書では、各ノードは 2 つの他のノードへの冗長接続を含み、すなわち、冗長性のために全ての端末装置へのデュアルパス（二重経路）が使用されるが、本開示の全ての実施形態で冗長接続が必要ということではない。

#### 【0016】

図 3 B には、カスケード接続された DAS クラスタの別の構成が示されており、この構成では、クラスタは 2 つの異なるケーブル接続形式を含んでいる。例えば、コネクタ 104 a、204 a、304 a、404 a、及び 504 a の構成は図 3 A に関して述べた構成と同じであるが、図 3 B のコネクタ 104 b、204 b、304 b、404 b、及び 504 b の構成は、図 3 A のコネクタ 104 b、204 b、304 b、404 b、及び 504 b の構成と異なる。図 3 B に示す構成に異なるケーブル接続形式が含まれることにより、このケーブル接続形式は、サーバ/ノードの各 HBA が同じサーバ/ノードへ接続される場合よりも、レイテンシーを低減しシステム/ドライブの利用可能性を増大することができる。図 3 B のコネクタ 104 b 及び 504 a はフェールオーバー接続であり、フェールオーバー接続は、クラスタの各サーバ/ノードが動作可能なときにはディスエーブルされるが、クラスタ内のノード又は接続がもはや動作不能（例えば、ノード障害）であるときには有効化される。クラスタ内のノード又は接続がもはや動作不能であるときは、ファームウェアが直ちにコネクタ 104 b 及び / 又は 504 a を有効化し、クラスタ内の使用可能な全てのノードへのアクセスが中断されないようにする。

#### 【0017】

図 4 には、図 3 A のカスケード接続の DAS クラスタの一部の概略図が示されている。図 4 に示すように、各サーバ/ノードの各 HBA は、2 つの主要なコンポーネントを含み、それらは、（1）システム 100 上で HBA の動作及び複数の HBA 間の接続を提供するための PCI（周辺要素相互接続）論理及び HBA 論理と、（2）ドライブ及び HBA 論理コンポーネント間、並びに HBA 論理コンポーネント及び外部物理層間のトラフィックのルーティングを処理するための仮想エキスパンドである。例えば、サーバ/ノード 100 の HBA 102 a は、PCI/HBA 論理コンポーネント 114 a 及び仮想エキスパ

ンダ 1 1 6 a を含み、サーバ / ノード 1 0 0 の H B A 1 0 2 b は P C I / H B A 論理コンポーネント 1 1 4 b 及び仮想エキスパンダ 1 1 6 b を含んでいる。コネクタ 1 0 6 a 及び 1 0 6 b は、複数のドライブ 1 1 0 を H B A 1 0 2 a の仮想エキスパンダ 1 1 6 a へ接続し、コネクタ 1 1 2 a 及び 1 1 2 b は複数のドライブ 1 1 0 を H B A 1 0 2 b の仮想エキスパンダ 1 1 6 b へ接続する。同様の構成が、クラスタの一部である他のサーバ / ノードの構成にも現れる。例えば、サーバ / ノード 2 0 0 の H B A 2 0 2 a は P C I / H B A 論理コンポーネント 2 1 4 a 及び仮想エキスパンダ 2 1 6 a を含み、サーバ / ノード 2 0 0 の H B A 2 0 2 b は P C I / H B A 論理コンポーネント 2 1 4 b 及び仮想エキスパンダ 2 1 6 b を含み、複数のドライブ 2 1 0 と仮想エキスパンダ 2 1 6 a 及び 2 1 6 b との間に接続が行われる。

10

#### 【 0 0 1 8 】

各サーバ / ノードは、サーバ / ノードのコンポーネント間に通信を提供するバスも含んでいる。例えば、サーバ / ノード 1 0 0 は P C I バス 1 1 8 を含み、P C I バス 1 1 8 は H B A 1 0 2 a 及び H B A 1 0 2 b のそれぞれへ接続されており、サーバ / ノード 2 0 0 は、H B A 2 0 2 a 及び H B A 2 0 2 b のそれぞれへ接続された P C I バス 2 1 8 を含んでいる。さらに、各サーバ / ノードは、例えば図 3 A 及び図 3 B に関して述べたような 2 つの他のサーバ / ノードへ接続される。各サーバ / ノード間の接続は、S A S ケーブル 4 0 6 等の S A S コネクタを含み、各サーバ / ノード間に外部接続を提供する。図 4 に示すように、サーバ / ノード 1 0 0 は、2 つの外部 S A S ケーブル 4 0 6 を含み、これらはクラスタ内の最後のマシン（例えば、エンドノード）と接続するためにループ接続される。本明細書で述べるように、1 又は複数の S A S ケーブルが無効にされ、フェールオーバーのケーブルとして機能することにより、無効な S A S トポロジを防止する。

20

#### 【 0 0 1 9 】

図 5 は、カスケード接続の D A S クラスタの一実施例の概略図を示している。通常、図 5 のカスケード接続の D A S クラスタの実施例は図 4 に示したものと異なり、この相違はサーバ / ノード間の接続に基づくものである。図 5 の実施例では、サーバ / ノードのシステムの片側に外部物理層があり、このサーバ / ノードの他方側は対応する外部物理層と結合形式が異なる。例えば、図示のように、サーバ / ノード 1 0 0 の H B A 1 0 2 a はコネクタ 5 0 2 a を介してサーバ / ノード 4 0 0 の H B A 4 0 2 a へ接続され、コネクタ 5 0 4 a を介してサーバ / ノード 2 0 0 の H B A 2 0 2 a へ接続されるが、H B A 1 0 2 b はコネクタ 5 0 2 b を介してサーバ / ノード 4 0 0 の H B A 4 0 2 b へ接続され、コネクタ 5 0 4 b を介してサーバ / ノード 3 0 0 の H B A 3 0 2 b へ接続される。サーバ / ノード 2 0 0 の H B A 2 0 2 a はコネクタ 5 0 4 a を介してサーバ / ノード 1 0 0 の H B A 1 0 2 a へ接続され、コネクタ 5 0 6 a を介してサーバ / ノード 3 0 0 の H B A 3 0 2 a へ接続されるが、サーバ / ノード 2 0 0 の H B A 2 0 2 b はコネクタ 5 0 6 b を介してサーバ / ノード 4 0 0 の H B A 4 0 2 b へ接続され、コネクタ 5 0 8 b を介してサーバ / ノード 3 0 0 の H B A 3 0 2 b へ接続される。そして、H B A 3 0 2 a はコネクタ 5 0 6 a を介してサーバ / ノード 2 0 0 の H B A 2 0 2 a へ接続され、コネクタ 5 0 8 a を介してサーバ / ノード 4 0 0 の H B A 4 0 2 a へ接続されるが、H B A 3 0 2 b はコネクタ 5 0 8 b を介してサーバ / ノード 2 0 0 の H B A 2 0 2 b へ接続され、コネクタ 5 0 4 b を介してサーバ / ノード 1 0 0 の H B A 1 0 2 b へ接続される。このような結合形式は、サーバ / ノードの各 H B A が同じサーバ / ノードへ接続される場合よりも、レイテンシーを低減し、システム / ドライブの使用可能性を増大することができる。

30

40

#### 【 0 0 2 0 】

コネクタ 5 0 2 a 及び 5 0 8 b はフェールオーバー接続であり、フェールオーバー接続は、クラスタの各サーバ / ノードが動作可能なときにはディスエーブル状態（無効）であるが、クラスタ内のノード又は接続がもはや動作不能（例えば、ノード障害）であるときには有効化される。クラスタ内のノード又は接続がもはや動作不能であるときは、ファームウェアが直ちにコネクタ 5 0 2 a 及び / 又は 5 0 8 b を有効にし、クラスタ内の使用可能な全てのノードへのアクセスが中断されないようにする。

50



## 【 0 0 2 1 】

データアクセス / 処理を迅速にするために、入力される I O ( 入力 / 出力 ) は効率的なルーティングアルゴリズムにより処理され、このアルゴリズムは、H B A のマルチコアプロセッサを使用し、マルチコアプロセッサは図 2 に示した 2 コア C P U 1 0 8 等である。この使用により、H B A の仮想エキスパンダ ( 例えば、仮想エキスパンダ 1 1 6 a ) のレイテンシーが低減される。例えば、H B A が 2 コアプロセッサを含む場合、第 2 のコアは仮想エキスパンダ専用にすることができる。

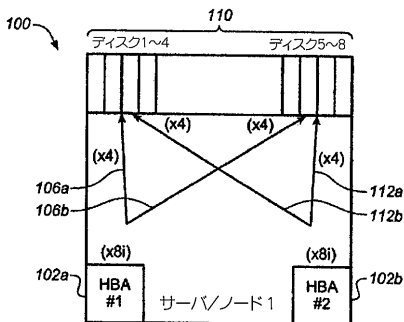
## 【 0 0 2 2 】

図 5 に示したカスケード接続の D A S クラスタの実施例は、R A I D ( リダンダント・アレイ・オブ・インディペンデント・ディスク ) の動作を実施できるように構成することもできる。例えば、クラスタのサーバ / ノードの複数のドライブ 1 1 0、2 1 0、3 1 0、4 1 0 を R A I D 構成 ( 図 5 に示したもの等 ) に配置して、ドライブ障害、システム障害、B H A 障害、又はケーブル障害の 1 又は複数を抑制すること等により、可用性が増大したクラスタを提供できるようにする。

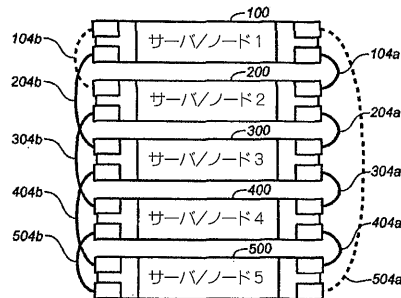
## 【 0 0 2 3 】

本発明の構成及び作用効果が上述の記載により理解されるであろう。また本発明の構成要素の形態、構成、及び配置においてさまざまな変更が可能であり、これは本発明の範囲及び技術思想から逸脱することや本発明の重要な効果のいずれも犠牲にすることなしに行うことができることが、明らかであろう。本明細書ですでに述べた形態は、例示的な実施形態に過ぎず、以下に記載の請求項はこのような変更を包括するものとする。

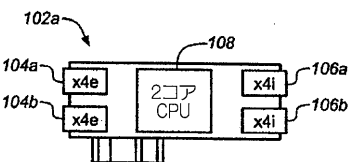
【 図 1 】



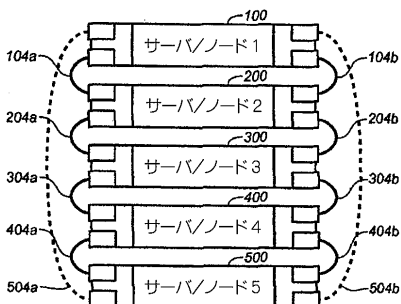
【 図 3 B 】



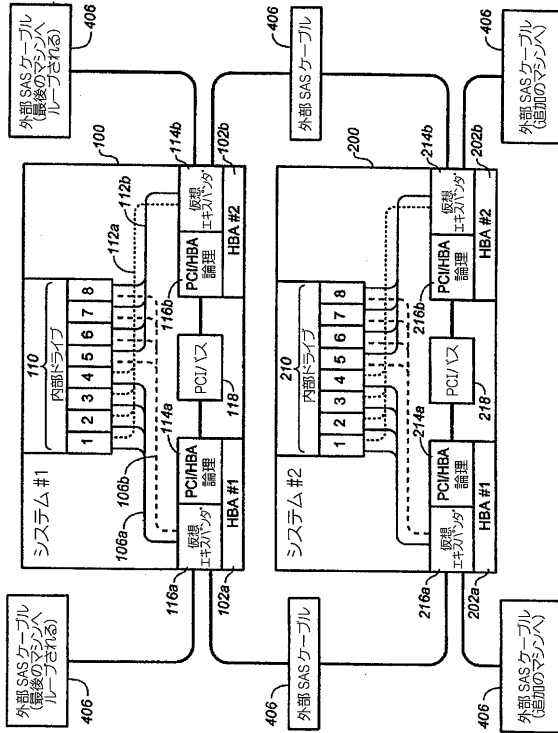
【 図 2 】



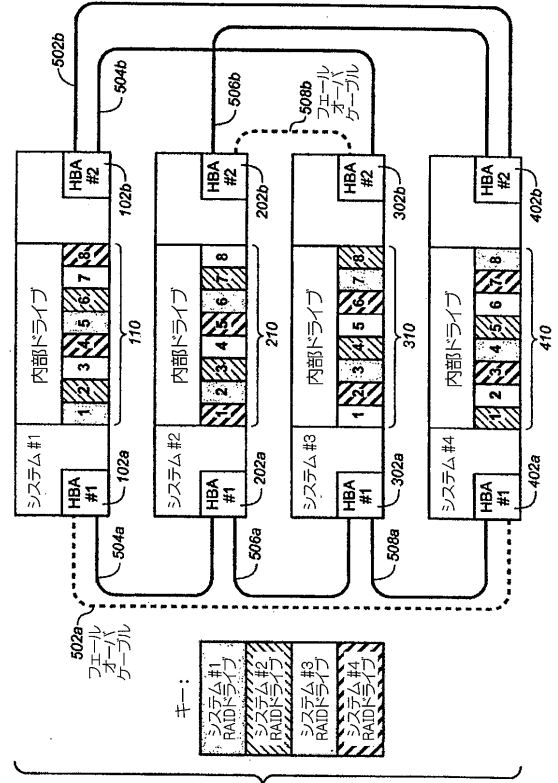
【 図 3 A 】



【図 4】



【図 5】



---

フロントページの続き

- (72)発明者 ルイス・ディー・ヴァーチャウトチック  
アメリカ合衆国カンザス州 6 7 2 2 0 , ウィチタ , イースト・ファルコン・コート 4 6 0 9
- (72)発明者 ジェイソン・エイ・アンレイン  
アメリカ合衆国カンザス州 6 7 2 2 0 , ウィチタ , イースト・ファルコン・コート 4 6 1 7
- (72)発明者 リード・エイ・カウフマン  
アメリカ合衆国カンザス州 6 7 0 0 2 , アンドーバー , ウェスト・セカンド・ストリート 6 1 0