

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7043201号

(P7043201)

(45)発行日 令和4年3月29日(2022.3.29)

(24)登録日 令和4年3月18日(2022.3.18)

(51)国際特許分類

F I

H 0 4 L 45/28 (2022.01)

H 0 4 L 45/28

H 0 4 L 47/125 (2022.01)

H 0 4 L 47/125

G 0 6 F 13/00 (2006.01)

G 0 6 F 13/00

請求項の数 20 (全24頁)

(21)出願番号	特願2017-156664(P2017-156664)	(73)特許権者	390019839
(22)出願日	平成29年8月14日(2017.8.14)		三星電子株式会社
(65)公開番号	特開2018-29337(P2018-29337A)		S a m s u n g E l e c t r o n i c s
(43)公開日	平成30年2月22日(2018.2.22)		C o . , L t d .
審査請求日	令和2年8月11日(2020.8.11)		大韓民国京畿道水原市靈通区三星路 1 2
(31)優先権主張番号	62/377487		9
(32)優先日	平成28年8月19日(2016.8.19)		1 2 9 , S a m s u n g - r o , Y e o
(33)優先権主張国・地域又は機関	米国(US)		n g t o n g - g u , S u w o n - s i
(31)優先権主張番号	15/344438		, G y e o n g g i - d o , R e p u b
(32)優先日	平成28年11月4日(2016.11.4)	(74)代理人	l i c o f K o r e a
(33)優先権主張国・地域又は機関	米国(US)		110000051
早期審査対象出願		(72)発明者	特許業務法人共生国際特許事務所
前置審査			グネスワラ アール . マリブディ
			アメリカ合衆国 9 4 5 5 5 カリフォル
			ニア州 フレモント ミモサ テラス 3 4
			最終頁に続く

(54)【発明の名称】 コンピューティング資源への高可用性アクセスを提供するコンピューティングシステム及び予備資源連結ファブリック

(57)【特許請求の範囲】

【請求項 1】

コンピューティング資源への高可用性アクセスを提供するコンピューティングシステムであって、

複数のインターフェイスと、

複数のコンピューティング資源セットと、

少なくとも3つからなる複数のスイッチと、を備え、

前記コンピューティング資源セットの各々は、複数のコンピューティング資源を含み、

前記複数のスイッチの各々は、ホストリンクを通じて前記インターフェイスの中の対応する1つに連結され、複数の資源連結を通じて前記複数のコンピューティング資源セットの中の対応する複数のコンピューティング資源セットのコンピューティング資源に直接連結され、前記複数のスイッチの中の1つが故障の場合、前記複数のスイッチの間の複数のクロス接続を通じて前記複数のスイッチの中の残りのスイッチにデータトラフィックが分散されるように構成され、

前記コンピューティングシステムは、N個（但し、Nは自然数）のインターフェイスを含み、

前記コンピューティング資源セットのコンピューティング資源の中の1つが全処理量を達成するのに要求されるトラフィック帯域幅の量を帯域幅ユニットB（但し、Bは正数）で表す場合、前記複数の資源連結の各々は、少なくとも1×Bのトラフィック帯域幅を有し、各コンピューティング資源セットがK（但し、Kは自然数）個のコンピューティング資源

を有する場合、前記複数のインターフェイスの各々の各ホストリンクは、 $K \times B$ のトラフィック帯域幅を支援し、

前記複数のクロス接続の各クロス接続は、少なくとも $K \times B / (N - 1)$ のトラフィック帯域幅を有することを特徴とするコンピューティングシステム。

【請求項 2】

前記少なくとも 3 つからなる複数のスイッチは、

第 1 スwitch と、

第 2 スwitch と、

第 1 クロス接続を通じて前記第 1 スwitch に直接連結され、第 2 クロス接続を通じて前記第 2 スwitch に直接連結される第 3 スwitch と、を含むことを特徴とする請求項 1 に記載のコンピューティングシステム。

10

【請求項 3】

第 1 ホストリンクを通じて前記第 1 スwitch に連結され、前記複数のインターフェイスの中の第 1 インターフェイスを含む第 1 サーバーノードと、

第 2 ホストリンクを通じて前記第 2 スwitch に連結され、前記複数のインターフェイスの中の第 2 インターフェイスを含む第 2 サーバーノードと、を更に含むことを特徴とする請求項 2 に記載のコンピューティングシステム。

【請求項 4】

第 3 ホストリンクを通じて前記第 3 スwitch に連結される第 3 サーバーノードを更に含むことを特徴とする請求項 3 に記載のコンピューティングシステム。

20

【請求項 5】

前記第 1 サーバーノードに直接連結される第 1 補助スswitch と、

前記第 2 サーバーノードに直接連結される第 2 補助スswitch と、

前記第 3 サーバーノードに直接連結される第 3 補助スswitch と、

前記第 1 補助スswitch を前記第 2 補助スswitch に直接連結する第 1 クロス接続と、

前記第 1 補助スswitch を前記第 3 補助スswitch に直接連結する第 2 クロス接続と、

前記第 2 補助スswitch を前記第 3 補助スswitch に直接連結する第 3 クロス接続と、を更に含むことを特徴とする請求項 4 に記載のコンピューティングシステム。

【請求項 6】

前記複数のコンピューティング資源セットの中の第 1 コンピューティング資源セットは、前記複数のコンピューティング資源として複数のデータ格納装置を含む第 1 データ格納装置セットを含み、

30

前記第 1 データ格納装置セットのデータ格納装置の各々は、前記複数の資源連結の中の第 1 資源連結を通じて前記第 1 スwitch に直接連結される第 1 ポートと、前記複数の資源連結の中の第 2 資源連結を通じて前記第 2 スwitch に直接連結される第 2 ポートとを有し、前記複数のコンピューティング資源セットの中の第 2 コンピューティング資源セットは、前記複数のコンピューティング資源として複数のデータ格納装置を含む第 2 データ格納装置セットを含み、

前記第 2 データ格納装置セットのデータ格納装置の各々は、前記複数の資源連結の中の第 3 資源連結を通じて前記第 1 スwitch に直接連結される第 1 ポートと、前記複数の資源連結の中の第 4 資源連結を通じて前記第 2 スwitch に直接連結される第 2 ポートとを有することを特徴とする請求項 3 に記載のコンピューティングシステム。

40

【請求項 7】

第 3 クロス接続を通じて前記第 1 スwitch に直接連結され、第 4 クロス接続を通じて前記第 2 スwitch に直接連結される第 4 スwitch を更に含むことを特徴とする請求項 6 に記載のコンピューティングシステム。

【請求項 8】

前記複数のコンピューティング資源として複数のデータ格納装置を含む第 3 データ格納装置セットを更に含み、

前記第 3 データ格納装置セットのデータ格納装置の各々は、第 5 資源連結を通じて前記第

50

3 スイッチに直接連結される第 1 ポートと、第 6 資源連結を通じて前記第 4 スイッチに直接連結される第 2 ポートとを有することを特徴とする請求項 7 に記載のコンピューティングシステム。

【請求項 9】

第 4 ホストリンクを通じて前記第 4 スイッチに直接連結される第 4 サーバーノードを更に含むことを特徴とする請求項 8 に記載のコンピューティングシステム。

【請求項 10】

前記複数のコンピューティング資源として複数のデータ格納装置を含む第 4 データ格納装置セットを更に含み、

前記第 4 データ格納装置セットのデータ格納装置の各々は、第 7 資源連結を通じて前記第 3 スイッチに直接連結される第 1 ポートと、第 8 資源連結を通じて前記第 4 スイッチに直接連結される第 2 ポートとを有することを特徴とする請求項 9 に記載のコンピューティングシステム。

10

【請求項 11】

前記複数のスイッチの各々は、他のスイッチの故障を検出し、前記複数のスイッチの各々に対応する各クロス接続及び前記ホストリンクを通じて前記第 1 データ格納装置セット及び前記第 2 データ格納装置セットへのアクセスを提供するように構成されることを特徴とする請求項 6 に記載のコンピューティングシステム。

【請求項 12】

前記複数のスイッチは、P C I e (P e r i p h e r a l C o m p o n e n t I n t e r c o n n e c t E x p r e s s) スイッチであり、

20

前記ホストリンク及び前記複数の資源連結は、P C I e バスリンクであることを特徴とする請求項 1 に記載のコンピューティングシステム。

【請求項 13】

少なくとも 2 つの高可用性ペア (h i g h - a v a i l a b i l i t y p a i r s) を備えるコンピューティングシステムであって、

各高可用性ペアは、

第 1 スイッチと、

第 2 スイッチと、

第 1 ホストリンクを通じて前記第 1 スイッチに連結される第 1 インターフェイスと、

30

第 2 ホストリンクを通じて前記第 2 スイッチに連結される第 2 インターフェイスと、

第 1 資源連結を通じて前記第 1 スイッチに直接連結され、第 2 資源連結を通じて前記第 2 スイッチに直接連結される複数のコンピューティング資源を含む第 1 コンピューティング資源セットと、

第 3 資源連結を通じて前記第 1 スイッチに直接連結され、第 4 資源連結を通じて前記第 2 スイッチに直接連結される複数のコンピューティング資源を含む第 2 コンピューティング資源セットと、

複数の第 1 クロス接続と、

複数の第 2 クロス接続と、を含み、

前記第 1 クロス接続の各々は、前記第 1 スイッチを他の高可用性ペアの各々の各スイッチに直接連結し、

40

前記第 2 クロス接続の各々は、前記第 2 スイッチを前記他の高可用性ペアの各々の各スイッチに直接連結し、

前記コンピューティングシステムは、N 個 (但し、N は自然数) のサーバーノードを含み、前記第 1 及び第 2 コンピューティング資源セットのコンピューティング資源の中の 1 つが全処理量を達成するのに要求されるトラフィック帯域幅の量を帯域幅ユニット B (但し、B は正数) で表す場合、前記第 1 ~ 第 4 資源連結の各々は、少なくとも $1 \times B$ のトラフィック帯域幅を有し、

各コンピューティング資源セットは K (但し、K は自然数) 個のコンピューティング資源を有する場合、各サーバーノードは、 $K \times B$ のトラフィック帯域幅を支援し、

50

前記第 1 及び第 2 クロス接続の各々は、少なくとも $K \times B / (N - 1)$ のトラフィック帯域幅を有することを特徴とするコンピューティングシステム。

【請求項 14】

前記第 1 コンピューティング資源セットは、前記複数のコンピューティング資源として複数のデータ格納装置を含む第 1 データ格納装置セットを含み、

前記第 1 データ格納装置セットのデータ格納装置の各々は、前記第 1 資源連結を通じて前記第 1 スイッチに直接連結される第 1 ポートと、前記第 2 資源連結を通じて前記第 2 スイッチに直接連結される第 2 ポートとを有し、

前記第 2 コンピューティング資源セットは、前記複数のコンピューティング資源として複数のデータ格納装置を含む第 2 データ格納装置セットを含み、

前記第 2 データ格納装置セットのデータ格納装置の各々は、前記第 3 資源連結を通じて前記第 1 スイッチに直接連結される第 1 ポートと、前記第 4 資源連結を通じて前記第 2 スイッチに直接連結される第 2 ポートとを有することを特徴とする請求項 13 に記載のコンピューティングシステム。

【請求項 15】

前記第 1 スイッチは、前記第 2 スイッチの故障を検出し、前記第 1 ホストリンク及び前記第 1 クロス接続を通じて前記第 1 データ格納装置セット及び前記第 2 データ格納装置セットへのアクセスを提供するように構成されることを特徴とする請求項 14 に記載のコンピューティングシステム。

【請求項 16】

前記第 1 及び第 2 スイッチは、PCIe (Peripheral Component Interconnect Express) スイッチであり、

前記第 1 及び第 2 ホストリンク及び前記第 1 ~ 第 4 資源連結は、PCIe バスリンクであることを特徴とする請求項 13 に記載のコンピューティングシステム。

【請求項 17】

前記少なくとも 2 つの高可用性ペアは、

第 1 高可用性ペアと、

第 2 高可用性ペアと、

第 3 高可用性ペアと、

第 4 高可用性ペアと、を含むことを特徴とする請求項 13 に記載のコンピューティングシステム。

【請求項 18】

予備資源連結ファブリックであって、

第 1 スイッチと、

第 2 スイッチと、

第 1 クロス接続を通じて前記第 1 スイッチに直接連結され、第 2 クロス接続を通じて前記第 2 スイッチに直接連結される第 3 スイッチと、

複数のコンピューティング資源を含む第 1 コンピューティング資源セットと、

複数のコンピューティング資源を含む第 2 コンピューティング資源セットと、を備え、

前記第 1 コンピューティング資源セットのコンピューティング資源の各々は、第 1 資源連結を通じて前記第 1 スイッチに直接連結される第 1 ポートと、第 2 資源連結を通じて前記第 2 スイッチに直接連結される第 2 ポートとを有し、

前記第 2 コンピューティング資源セットのコンピューティング資源の各々は、第 3 資源連結を通じて前記第 1 スイッチに直接連結される第 1 ポートと、第 4 資源連結を通じて前記第 2 スイッチに直接連結される第 2 ポートとを有し、

前記予備資源連結ファブリックは、N 個 (但し、N は自然数) のインターフェイスを更に含み、前記第 1 ~ 第 3 スイッチの各々は、ホストリンクを通じて前記 N 個のインターフェイスの対応する 1 つに連結され、

前記第 1 及び第 2 コンピューティング資源セットのコンピューティング資源の中の 1 つが全処理量を達成するのに要求されるトラフィック帯域幅の量を帯域幅ユニット B (但し、

10

20

30

40

50

Bは正数)で表す場合、前記第1～第4資源連結の各々は、少なくとも $1 \times B$ のトラフィック帯域幅を有し、

各コンピューティング資源セットがK(但し、Kは自然数)個のコンピューティング資源を有する場合、前記N個のインターフェイスの各々の各ホストリンクは、 $K \times B$ のトラフィック帯域幅を支援し、

前記第1クロス接続及び前記第2クロス接続の各々は、少なくとも $K \times B / (N - 1)$ のトラフィック帯域幅を有することを特徴とする予備資源連結ファブリック。

【請求項19】

前記第1スイッチは、前記第2スイッチの故障を検出し、前記第1クロス接続を通じて前記第1コンピューティング資源セット及び前記第2コンピューティング資源セットへのアクセスを提供するように構成されることを特徴とする請求項18に記載の予備資源連結ファブリック。

10

【請求項20】

第3資源連結を通じて前記第1スイッチに直接連結され、第4資源連結を通じて前記第2スイッチに直接連結される第4スイッチを更に含むことを特徴とする請求項18に記載の予備資源連結ファブリック。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はコンピューティング資源への高可用性アクセスを提供するコンピュータシステム及び予備資源連結ファブリックに係り、より詳しくは、冗長性(redundancy)を提供するコンピュータシステム及び予備資源連結ファブリックに関する。

20

【背景技術】

【0002】

コンピュータシステムの分野において、データ及び他のコンピューティング資源への信頼性ある高性能アクセス(reliable and high performance access)はビジネス及び日常生活において重要である。高可用性(HA; high availability)の用語は1つの要素の故障(又は1つの障害箇所; single points of failure)を除去するか、又は減少させる(例えば、1つの要素の故障が全体システムの故障とならないようにシステムに冗長性を提供する)システムを説明するために頻繁に使用される。

30

【0003】

高可用性の長所を有する例示的なコンピュータシステムはLAN(local area network)のようなコンピュータネットワーク又はインターネットを通じてデータの検索又は格納のためのデータ資源を提供できるSAN(storage area network)のようなデータ格納システムを含む。

【発明の概要】

【発明が解決しようとする課題】

【0004】

本発明の目的は向上した性能を有するコンピューティング資源への高可用性及び負荷均性を提供するコンピューティングシステム及び予備資源連結ファブリックを提供することにある。

40

【課題を解決するための手段】

【0005】

本発明の目的は、適応的多重経路ファブリックの使用を通じてコンピューティング資源への高可用性及び負荷均等性を提供することにある。

【0006】

本発明の一実施形態によれば、コンピューティング資源への高可用性アクセスを提供するコンピューティングシステムは、複数のインターフェイス、複数のコンピューティング資源セット、及び少なくとも3つのスイッチを含み、前記複数のコンピューティング資源セ

50

ットの各々は複数のコンピューティング資源を含み、前記少なくとも3つのスイッチの各々はホストリンクを通じて前記インターフェイスの中の対応する1つと連結され、複数の資源連結を通じて前記複数のコンピューティング資源セットの中の対応する1つと連結され、前記スイッチの中の1つが故障の場合、前記スイッチの間の複数のクロス接続を通じて前記スイッチの中の残されているスイッチにデータトラフィックが分散されるように構成される。

【0007】

前記少なくとも3つのスイッチは第1スイッチ、第2スイッチ、及び第1クロス接続を通じて前記第1スイッチと直接連結され、第2クロス接続を通じて前記第2スイッチと直接連結される第3スイッチを含む。

10

【0008】

前記コンピューティングシステムは第1ホストリンクを通じて前記第1スイッチと連結され、前記複数のインターフェイスの中の第1インターフェイスを含む第1サーバーノード、及び第2ホストリンクを通じて前記第2スイッチと連結され、前記複数のインターフェイスの中の第2インターフェイスを含む第2サーバーノードをさらに含む。

【0009】

前記コンピューティングシステムは第3ホストリンクを通じて前記第3スイッチと連結される第3サーバーノードをさらに含む。

【0010】

前記コンピューティングシステムは前記第1サーバーノードと直接連結される第1補助スイッチ、前記第2サーバーノードと直接連結される第2補助スイッチ、前記第3サーバーノードと直接連結される第3補助スイッチ、前記第1補助スイッチを前記第2補助スイッチと直接連結する第1クロス接続、前記第1補助スイッチを前記第3補助スイッチと直接連結する第2クロス接続、前記第2補助スイッチを前記第3補助スイッチと直接連結する第3クロス接続をさらに含む。

20

【0011】

前記複数のコンピューティング資源セットの中の第1コンピューティング資源セットは第1データ格納装置セットを含み、前記第1データ格納装置セットのデータ格納装置の各々は前記複数の資源連結の中の第1資源連結を通じて前記第1スイッチと直接連結された第1ポート及び前記複数の資源連結の中の第2資源連結を通じて前記第2スイッチと直接連結された第2ポートを含み、前記複数のコンピューティング資源セットの中の第2コンピューティング資源セットは第2データ格納装置セットを含み、前記第2データ格納装置セットのデータ格納装置の各々は前記複数の資源連結の中の第3資源連結を通じて前記第1スイッチと直接連結された第3ポート及び前記複数の資源連結の中の第4資源連結を通じて前記第2スイッチと直接連結される第4ポートを含む。

30

【0012】

前記コンピューティングシステムは第3クロス接続を通じて前記第1スイッチと直接連結され、第4クロス接続を通じて前記第2スイッチと直接連結される第4スイッチをさらに含む。

【0013】

前記コンピューティングシステムは第3データ格納装置セットをさらに含み、前記第3データ格納装置セットのデータ格納装置の各々は第5資源連結を通じて前記第3スイッチと直接連結される第5ポート及び第6資源連結を通じて前記第4スイッチと直接連結される第6ポートを含む。

40

【0014】

前記コンピューティングシステムは第4ホストリンクを通じて前記第4スイッチと直接連結された第4サーバーノードをさらに含む。

【0015】

前記コンピューティングシステムは第4データ格納装置セットをさらに含み、前記第4データ格納装置セットのデータ格納装置の各々は第7資源連結を通じて前記第3スイッチと

50

直接連結される第 7 ポート及び第 8 資源連結を通じて前記第 4 スイッチと直接連結される第 8 ポートを含む。

【 0 0 1 6 】

前記スイッチの各々は他のスイッチの故障を検出し、前記スイッチに対応する前記クロス接続及び前記ホストリンクを通じて前記第 1 データ格納装置セットのデータ格納装置及び前記第 2 データ格納装置セットのデータ格納装置へのアクセスを提供するように構成される。

【 0 0 1 7 】

前記コンピューティングシステムは N 個（但し、 N は自然数）のインターフェイスを含み、帯域幅ユニット B （但し、 B は正数）は前記コンピューティング資源セットの中の 1 つが全帯域幅（full bandwidth）にて動作するのに要求されるトラフィック帯域幅の量を示し、前記複数の資源連結の各々は少なくとも $1 \times B$ のトラフィック帯域幅を含み、前記複数のコンピューティング資源セットの各々は K （但し、 K は自然数）個以下のコンピューティング資源を含み、前記複数のインターフェイスの各々のホストリンクの各々は $K \times B$ の正常のトラフィック帯域幅及び $K \times B$ 以上の最大のトラフィック帯域幅を支援し、前記複数のクロス接続の各々は少なくとも $K \times B / (N - 1)$ のトラフィック帯域幅を含む。

10

【 0 0 1 8 】

前記スイッチは PCIe（Peripheral Component Interconnect Express）スイッチであり、前記ホストリンク及び前記資源連結は PCIe バスリンクである。

20

【 0 0 1 9 】

本発明の一実施形態に係るコンピューティングシステムは少なくとも 2 つの高可用性ペア（high-availability pairs）を含み、前記高可用性ペアの各々は、第 1 スイッチ、第 2 スイッチ、第 1 ホストリンクを通じて前記第 1 スイッチと連結される第 1 インターフェイス、第 2 ホストリンクを通じて前記第 2 スイッチと連結される第 2 インターフェイス、第 1 資源連結を通じて前記第 1 スイッチと連結され、第 2 資源連結を通じて前記第 2 スイッチと連結される第 1 コンピューティング資源セット、第 3 資源連結を通じて前記第 1 スイッチと連結され、第 4 資源連結を通じて前記第 2 スイッチと連結される第 2 コンピューティング資源セット、複数の第 1 クロス接続、及び複数の第 2 クロス接続を含み、前記複数の第 1 クロス接続の各々は前記第 1 スイッチを他の高可用性ペアの各スイッチに直接連結し、前記第 2 クロス接続の各々は前記第 2 スイッチを前記他の高可用性ペアの各スイッチに直接連結する。

30

【 0 0 2 0 】

前記第 1 コンピューティング資源セットは第 1 データ格納装置セットを含み、前記第 1 データ格納装置セットのデータ格納装置の各々は前記第 1 資源連結を通じて前記第 1 スイッチと直接連結される第 1 ポート及び前記第 2 資源連結を通じて前記第 2 スイッチと直接連結される第 2 ポートを含み、前記第 2 コンピューティング資源セットは第 2 データ格納装置セットを含み、前記第 2 データ格納装置セットのデータ格納装置の各々は前記第 3 資源連結を通じて前記第 1 スイッチと直接連結される第 3 ポート及び前記第 4 資源連結を通じて前記第 2 スイッチと直接連結される第 4 ポートを含む。

40

【 0 0 2 1 】

前記第 1 スイッチは前記第 2 スイッチの故障を検出し、前記第 1 ホストリンク及び前記第 1 クロス接続を通じて前記第 1 データ格納装置セットのデータ格納装置及び前記第 2 データ格納装置セットのデータ格納装置へのアクセスを提供するように構成される。

【 0 0 2 2 】

前記コンピューティングシステムは、 N 個（但し、 N は自然数）のサーバーノードを含み、帯域幅ユニット B （但し、 B は正数）は前記コンピューティング資源の中の 1 つが全帯域幅（full bandwidth）にて動作するのに要求されるトラフィック帯域幅の量であり、資源連結の各々は少なくとも $1 \times B$ のトラフィック帯域幅を含み、コンピュ

50

ーティング資源セットの各々は K （但し、 K は自然数）個以下のコンピューティング資源を含み、前記サーバーノードの各々は $K \times B$ の正常（normal）のトラフィック帯域幅を支援し、 $K \times B$ 以上の最大（maximum）のトラフィック帯域幅を支援し、クロス接続の各々は少なくとも $K \times B / (N - 1)$ のトラフィック帯域幅を含む。

【0023】

前記スイッチはPCIe（Peripheral Component Interconnect Express）スイッチであり、前記ホストリンク及び前記資源連結はPCIeバスリンクである。

【0024】

前記少なくとも2つの高可用性ペアは第1高可用性ペア、第2高可用性ペア、第3高可用性ペア、及び第4高可用性ペアを含む。

10

【0025】

本発明の一実施形態に係る予備資源連結ファブリックは第1スイッチ、第2スイッチ、第1クロス接続を通じて前記第1スイッチと直接連結され、第2クロス接続を通じて前記第2スイッチと直接連結される第3スイッチ、第1コンピューティング資源セット、及び第2コンピューティング資源セットを含み、前記第1コンピューティング資源セットの各々は第1資源連結を通じて前記第1スイッチと直接連結された第1ポート及び第2資源連結を通じて前記第2スイッチと直接連結される第2ポートを含み、前記第2コンピューティング資源セットの各々は第3資源連結を通じて前記第1スイッチと直接連結される第3ポート及び第4資源連結を通じて前記第2スイッチと直接連結される第4ポートを含む。

20

【0026】

前記第1スイッチは前記第2スイッチの故障を検出し、前記第1クロス接続を通じて前記第1コンピューティング資源セット及び前記第2コンピューティング資源セットへのアクセスを提供するように構成される。

【0027】

前記予備資源連結ファブリックは第3資源連結を通じて前記第1スイッチと直接連結され、第4資源連結を通じて前記第2スイッチと直接連結される第4スイッチをさらに含む。

【0028】

前記予備資源連結ファブリックは第5資源連結を通じて前記第3スイッチと直接連結され、第6資源連結を通じて前記第4スイッチと直接連結される第3コンピューティング資源セットをさらに含む。

30

【0029】

前記予備資源連結ファブリックは第7資源連結を通じて前記第3スイッチと直接連結され、第8資源連結を通じて前記第4スイッチと直接連結される第4コンピューティング資源セットをさらに含む。

【発明の効果】

【0030】

本発明に係るインターフェイス及びコンピューティング資源の間に位置した適応的ファブリック（adaptive fabric）は、故障状況においてコンピューティング資源へのアクセスを維持し、一部の実施形態において、故障状況において資源への最大限のパフォーマンスアクセス（full performance access）を維持する。したがって、向上した性能を有する高可用性及び負荷均等のためのコンピューティングシステム及び予備資源連結ファブリックが提供される。

40

【図面の簡単な説明】

【0031】

詳細な説明と共に、添付した図面は本発明の例示的な実施形態を示し、説明と共に本発明の原理を説明する。

【図1】本発明の一実施形態に係る適応的多重経路ファブリック（adaptive multipath fabric）を含むデータ格納システムを示すブロック図である。

【図2】本発明の一実施形態に係る適応的多重経路ファブリックを含むデータ格納システ

50

ムを示すブロック図である。

【図3】本発明の一実施形態に係るスイッチの故障及びデータ伝送トラフィックのリバランシング (rebalancing) を示すブロック図である。

【図4】本発明の一実施形態に係る、サーバーノードの故障及びデータ伝送のリバランシングを示すブロック図である。

【図5】本発明の一実施形態に係るサーバーノードを連結する補助ファブリックをさらに含むシステムのブロック図である。

【図6】本発明の一実施形態に係る2つのサーバーノード及び1つのデータ格納装置セットを含む小さい高可用性構成 (small high availability configuration) を示すブロック図である。

10

【図7】本発明の一実施形態に係る追加的なデータ格納装置セットを含む図6に図示された構成からの拡張を示す。

【図8】本発明の一実施形態に係る追加的なサーバーノードを含む図6に図示された構成からの拡張を示す。

【図9】本発明の一実施形態に係る、192GB/sの全体処理量のために8個のスイッチを通じて48個のデータ格納装置に連結された8個のサーバーノードを含むシステムを示すブロック図である。

【発明を実施するための形態】

【0032】

以下の詳細な説明において、説明を簡易にするために本発明の特定例示的な実施形態のみを説明する。当業者は、本発明が多様な他の形態に具現されるので、例示的な実施形態に限定されると理解されてはならない。詳細な説明の全体において、類似の参照番号は類似の構成要素を示す。

20

【0033】

本発明の実施形態は、インターフェイスのセットを通じてコンピューティング資源への高可用性アクセス (high availability access) を提供するシステム及び方法と一般的に連関される。インターフェイス及びコンピューティング資源の間に位置した適応的ファブリック (adaptive fabric) は故障状況においてもコンピューティング資源へのアクセスを維持し、実施形態において、故障状況においてコンピューティング資源への最大限のパフォーマンスアクセス (full performance access) を維持する。

30

【0034】

例えば、コンピューティング資源への高可用性アクセス (high-availability access) を提供するデータ格納コンピューティングシステムは複数のインターフェイスと、複数のコンピューティング資源セットと、少なくとも3つのスイッチと、を含む。コンピューティング資源セットの各々は複数のコンピューティング資源を含む。スイッチの各々はホストリンクを通じてインターフェイスの中の対応する1つと連結され、資源連結 (resource connection) を通じてコンピューティング資源セットの中の対応する1つと連結される。スイッチの各々はスイッチの中の1つが故障の場合、スイッチ間の複数のクロス接続 (cross-connections) を通じてデータトラフィックがスイッチの中の残るスイッチに分散されるように構成される。

40

【0035】

さらに具体的に、SAN (storage area network) システムのようなデータ格納システムは1つ以上のサーバーノードと連結されたデータ格納装置を含む。例えば、データ格納装置はPCIe (peripheral component interconnect express) バスのようなバスを通じてサーバーノードと連結されたSSD (solid state drive) 又はHDD (hard disk drive) である。例えば、各サーバーノードは中央処理ユニット、メモリ、及びデータ格納装置に格納されたデータへの遠隔アクセスを提供するネットワークインターフェイスを含む。この時、データ格納装置はサーバーノードにマッピングされる。しかし、1つの

50

サーバーノード (single server node) が故障の場合 (例えば、ネットワークインターフェイス、バス、又はCPUが故障の場合)、データ格納装置に格納されたデータへのアクセスが損失されるので、1つのサーバーノードのみを使用するネットワークストレージシステムは高可用性 (HA) を提供できない。

【0036】

このように、本発明の実施形態は、コンピューティング資源への高可用性アクセスを提供する適応的ファブリック (adaptive fabric) に係る。一実施形態において、コンピューティング資源の各々は複数のスイッチと連結され、スイッチの各々は適応的ファブリックのクロス接続を通じて少なくとも1つの他のスイッチと連結される。スイッチの各々はコンピューティング資源の利用者との通信のためのインターフェイス (例えば、ネットワークインターフェイス) と連結される。インターフェイスの故障又はスイッチの故障が発生した場合、適応的ファブリックは他のインターフェイスのクロス接続を通じてデータへの経路を再設定 (reroute) する。

10

【0037】

一実施形態において、インターフェイスはホストサーバーノードの構成要素である。この時、サーバーノードはプロセッサ (CPU) 及びメモリを含むコンピュータシステムである。サーバーノードは、サーバーノードのメモリに格納され、サーバーノードのプロセッサによって実行され、駆動されるアプリケーションを通じて利用者にコンピューティング資源と連結されたサービスへのアクセスを提供する。例えば、コンピューティング資源はデータ格納装置セットであり、この場合、アプリケーションはネットワークファイルサーバー、ウェブサーバー、データベースサーバー等である。他の例として、コンピューティング資源はローレイテンシキャッシュ (low latency caches) を提供する動的メモリ (dynamic memory) である。その他の例として、コンピューティング資源はグラフィック処理ユニット (GPU; graphical processing unit) であり、この場合、アプリケーションは、例えば3次元レンダリングエンジン、マシンラーニングトレーニングプラットフォーム (例えば、トレーニングニューラルネットワーク)、暗号通貨マイナー (cryptocurrency miner) (例えば、ビットコイン) 等である。

20

【0038】

本発明の実施形態は、データ格納装置のようなコンピューティング資源への十分な帯域幅 (伝送速度) の提供と関連する。サーバーノードにあまりにも多くのデータ格納装置が連結された場合、データ格納装置に最大限のパフォーマンス (full performance) を可能にする、サーバーノード及びデータ格納装置の間の可用である帯域幅が十分でないこともあり得る。さらに具体的に、1つのサーバーノードシステムにおいて、8個のSSDがPCIeスイッチと連結され、各SSDがPCIeスイッチへの4レーンリンク (X4) を飽和させ、サーバーノードがPCIeスイッチへの32レーンリンク (X32) を含む場合、8個のSSDだけでサーバーノードへの連結を飽和させるのに十分である。追加的なデータ格納装置がシステムに追加される場合、サーバーノード及びPCIeスイッチの間の連結がシステムにおいて隘路現象 (bottleneck) として作用するので、データ格納装置の全体を最大限のパフォーマンスにより動作させるのに帯域幅が十分ではないことがあり得る。一部の状況において、サーバーノード及びネットワークアダプターの間の連結がシステムにおける隘路現象に類似して作用する。

30

40

【0039】

一部のデータ格納装置は冗長性 (redundancy) を提供する2つの連結ポート (two connection ports) を含む。例えば、デュアルPCIeポートを含むデータ格納装置は第1サーバーノードと連結された第1ポート及び第2サーバーノードと連結された第2ポートを含む。このような方式により、サーバーノードの中の1つが故障の場合、データ格納装置は他のサーバーノードを通じて相変わらずアクセスされる。

【0040】

しかし、このような方式において、サーバーノードの故障は帯域幅制限 (bandwidth

50

th limitations)を悪化させる。上述した例を続いて参照すれば、2つのサーバーノードの全てが連結されたデータ格納装置に対する十分な帯域幅を提供できる反面、サーバーノードの中の1つが故障の場合、データ格納装置へのノードからのすべてのトラフィックが、生存サーバーノード(surviving server node)によって管理される。生存サーバーノードは追加的なトラフィックを管理するのに十分な帯域幅を有しないことがある。特に、データ格納装置の帯域幅要求が既にサーバーノードへの2つのリンクを飽和させた場合、サーバーノードの中の1つの故障は約50%のパフォーマンス減少を発生させる。

【0041】

本発明の一部の実施形態は複数のサーバーノードを通じてデータ格納装置へのアクセスを提供するシステム及び方法に係る。この時、システム及び方法はサーバーノードの故障状況においてパフォーマンス低下無しにデータ格納装置の可能な最大限のパフォーマンス(full performance potential)を維持する。さらに詳細には、本発明の実施形態は、複数のスイッチを通じてサーバーノードと(マルチポートデータ格納装置のような)コンピューティング資源を連結する適応的多重経路ファブリック(adaptive multipath fabric)と関連する。この時、多重経路ファブリックはシステムの故障状況においてデータトラフィックに対する代替経路(alternate paths)を提供する。本発明の実施形態は、要求された特定のシステムパフォーマンスプロファイルを達成するために、多様な複数のサーバーノード及びコンピューティング資源(例えば、データ格納装置)のセットに適用される。本発明の実施形態は、ノード故障状況においても、特定のパフォーマンスプロファイルを維持する一方で、初期構成に対してデータ格納装置又はサーバーノードをさらに追加することによって、コンピューティング資源(例えば、データ格納容量)及び一般コンピューティング能力の双方の増大(scaling)を可能にする。

【0042】

説明を簡易にするために、本発明の実施形態は、以下において、PCIeスイッチ及びPCIeファブリックを通じてホストサーバーノードと連結されたNVMe(non-volatile memory express)インターフェイスを含むソリッドステートドライブを参照して説明する。しかし、本発明の実施形態はこれに限定されず、ファブリックの基本構造(underlying architecture)はイーサネット(登録商標)(Ethernet(登録商標))、IB(Infiniband(登録商標))、ファイバチャネル(Fibre Channel)、SCSI(small computer system interface)、SAS(serially attached SCSI)等の他のインターフェイスに適用される。追加的に、本発明の実施形態はハードディスクドライブ、テープドライブ、DRAM(dynamic random access memory)のような揮発性メモリの他の形態の資源、及びベクトルプロセッサ、グラフィック処理ユニット(GPUs)、デジタル信号プロセッサ(DSPs; digital signal processors)、及びFPGA(field programmable gate array)のような演算ユニットに高可用性を提供するのに適用される。

【0043】

<多重経路ファブリック構造(Multipath fabric structure)>
本発明の実施形態は、コンピューティング資源がデュアルポートデータ格納装置セットであり、連結及びスイッチがPCIe連結及びスイッチであり、インターフェイスがネットワークインターフェイスであるデータ格納システムの特定事例を参照して以下に説明する。しかし、本発明の実施形態はこれに限定されず、他の形態のコンピューティング資源、連結プロトコル、及びインターフェイスに適用できる。

【0044】

図1は本発明の一実施形態に係る適応的多重経路ファブリックを含むデータ格納システム100を示すブロック図である。図1の実施形態はサーバーノード20をデュアルポート

10

20

30

40

50

データ格納装置セット 30 と連結するファブリック 10 を含む。ファブリック 10 は印刷回路基板上のパターン (trace)、複数の電気配線 (例えば、リボンケーブル、 mini - SAS HD ケーブル、 OCuLink ケーブル等)、及びそれらの組合せのような多様な方式により具現される。ファブリック 10 はスイッチ 40 の間のクロス接続 12 AC、12 BC、12 BD、12 AD を含む。ノード故障の場合に帯域幅バランシング (balancing) の提供のために追加的に、クロス接続 12 AC、12 BC、12 BD、12 AD はエンドポイント再割当及び帯域幅バランシングに使用されて入力 / 出力 (I/O) 負荷及び非平衡ノード CPU 使用をカウンティングするだけでなく、データ格納装置セット 30 の間のピアツーピア通信 (例えば、サーバーノード 20 からの干渉無しに、第 1 データ格納装置セット 30 A 及び第 2 データ格納装置セット 30 B の間の直接メモリアクセス伝送 (direct memory access transfers)) を提供する。

10

【0045】

デュアルポートデータ格納装置セット 30 の各々は 1 つ以上のデュアルポートデータ格納装置 32 を含む。各デュアルポートデータ格納装置 32 は資源リンク 16 を通じて 2 つの他のスイッチと連結される。本文に使用する ‘ ‘ スイッチ (switch) ’ ’ の用語は通信のためにスイッチと連結された装置に対して複数の通信経路を提供する電氣的な構成要素を示す。スイッチは装置の間のトラフィックの経路を設定し、スイッチと連結された通信装置の間の連結を設定する。

【0046】

図 1 に図示したように、各サーバーノード 20 は 1 つ以上のプロセッサ 24 (例えば、Intel (登録商標) Xeon (登録商標) プロセッサ) を含む。1 つ以上のプロセッサ 24 は PCIe 連結を通じてネットワークインターフェイスカード (NIC) 26 (例えば、イーサネット (登録商標) NIC) と連結され、ホストリンク 14 (例えば、他の PCIe 連結又は複数の PCIe 連結) を通じて対応するスイッチ 40 と連結される。各サーバーノード 20 は高可用性ピア (HA peer; high availability peer) と指称される他の 1 つのサーバーノード 20 とペアをなす。例えば、サーバーノード 20 A、20 B は HA ピアであり、ペア 22 AB を形成する。共通 HA ピアのペア 22 はファブリック 10 を通じて 2 つの予備経路 (redundant paths) により 1 つ以上のデータ格納装置セット 30 をアクセスする。例えば、HA ピア 22 AB はデータ格納装置セット 30 A、30 B へアクセスする。HA ピアノードの他のペア 22 は他のデータ格納装置セット 30 に他の予備経路によりアクセスする。例えば、HA ピア 22 CD はデータ格納装置セット (30 C、30 D) へアクセスする。

20

30

【0047】

ファブリック 10 は複数の HA ピアのペア 22 を連結して、1 つのサーバーノード故障状況において、データ格納装置 32 及びすべてのノードにわたった均衡ある帯域幅のための連結を提供する。N 個のサーバーノード 20 (又は N 個のスイッチ 40) を具備するシステムにおいて、1 つのスイッチ 40 から他の (N - 2) 個のスイッチ 40 の各々への追加的なクロス接続帯域幅の量は正常、非故障モードの動作において各サーバーノード 20 によって支援される帯域幅 (正常帯域幅 ‘ ‘ normal bandwidth ’ ’) の $1 / (N - 1)$ 倍である。結果的に、1 つのノードの故障による帯域幅の損失は残る (N - 1) 個のノードによって減少される。

【0048】

さらに詳細には、図 1 は 4 つのサーバーノード 20 (N = 4) を示す。説明を簡易にするために、図 1 は帯域幅ユニット B で表した帯域幅を示し、この時、B はデータ格納装置 32 の中の 1 つの帯域幅要求を示す。図 1 において、データ格納装置セット 30 の各々は 6 個のデータ格納装置 32 を含む。したがって、データ格納装置セット 30 の各々は 6 B の帯域幅 (各データ格納装置 32 当たり 1 B) を必要とする。図 1 に図示したように、4 つのサーバーノード 20 はシステムにおいて 24 個のデータ格納装置 32 に総 24 B の処理量 (ノード当たり 6 B) を提供する。

50

【 0 0 4 9 】

サーバーノード 2 0 及びそれに対応するスイッチ 4 0 の間の各連結の実際 (a c t u a l) の帯域幅の容量 (c a p c i t y) は 8 B であるので、データ格納装置セット 3 0 の各々において要求される 6 B の帯域幅の容量を 2 B ぐらい超過する。また、図 1 の実施形態において、ホストプロセッサ 2 4 及びネットワークインターフェイス 2 6 の間の連結はホストプロセッサ 2 4 及びスイッチ 4 0 の間の帯域幅、例えば 8 B と少なくとも同一の帯域幅を有する。

【 0 0 5 0 】

図 1 に図示した実施形態において、各ノードから非 H A ピアノードの各々へのクロス接続帯域幅は $6 B / 3 = 2 B$ である。例えば、第 1 サーバーノード 2 0 A と対応する第 1 スイッチ 4 0 A 及びその非 H A ピアノード 2 0 C、2 0 D のスイッチ 4 0 C、4 0 D の間のクロス接続 1 2 A C、1 2 A D の各々のクロス接続帯域幅は 2 B である。同様に、第 2 サーバーノード 2 0 B と対応する第 2 スイッチ 4 0 B 及びその非 H A ピアノード 2 0 C、2 0 D のスイッチ 4 0 C、4 0 D の間のクロス接続 1 2 B C、1 2 B D の各々のクロス接続帯域幅は 2 B である。図 1 に図示したように、クロス接続 1 2 はスイッチ 4 0 の間に形成される。しかし、本発明の実施形態がこれに限定されず、上述したように、クロス接続 1 2 の各々の最小帯域幅は正常、非故障モードの動作において各サーバーノード 2 0 によって支援される帯域幅をサーバーの個数 (N) から 1 を差し引いた値により割った帯域幅に設定される。その他の実施形態において、システムが複数のサーバーノード 2 0 の故障を許容するように設計された場合、クロス接続 2 0 の各々の最小帯域幅は正常、非故障モードの動作において各サーバーノード 2 0 によって支援される帯域幅をサーバーの個数 (N) から許容される故障の個数を差し引いた値により割った帯域幅に設定される。

【 0 0 5 1 】

図 2 は本発明の一実施形態に係る適応的多重経路ファブリックを含むデータ格納システムを示すブロック図であり、 $B = X 4$ (例えば、4 レーン P C I e リンクの帯域幅) である場合の適応的多重経路を含むデータ格納システム 1 0 0 を示すブロック図である。P C I e 3 . 0 の場合、 $X 4$ リンクは約 $4 G B / s$ の最大帯域幅又は処理量を提供する。類似の構成に図 1 と同様の参照番号が付与され、このような構成に対する説明は省略する。

【 0 0 5 2 】

図 2 に図示した実施形態において、ノード当たり 6 4 個の P C I e 3 . 0 レーンを具備する 4 つのサーバーノード 2 0 は各々セット当たり 6 個のドライブを含む 4 つのセットに配置された 2 4 個のデュアルポート N V M e S S D と連結される。2 4 個のデュアルポート N V M e S S D の各々は 2 つのエンドポイント (即ち、2 つの S S D エンドポイント) を含む。この時、各エンドポイントはサーバーノードの中のいずれか 1 つに割り当てられる。例えば、第 1 データ格納装置セット 3 0 A において S S D のエンドポイントの各々は第 1 サーバーノード 2 0 A に割り当てられる。P C I e 3 . 0 を使用する場合、図 2 に図示したシステムは $9 6 G B / s$ 使用者データ処理量 (P C I e 3 . 0 帯域幅の 9 6 レーン) のエンドツーエンドシステムパフォーマンスプロファイル及び高可用性 (1 つのサーバーノード故障の状況においてすべての S S D への最大限のパフォーマンスアクセスが維持されること) を保障する。

【 0 0 5 3 】

図 2 の例示的な実施形態において、全負荷により動作中である場合、N V M e S S D 3 2 の各々は、約 $4 G B / s$ によりデータを伝送し単一 $X 4$ リンクによって提供される。また、デュアルポート N V M e S S D の 2 つのポートの各々は $X 4$ リンクを提供する。結果的に、N V M e S S D の 2 つのポートの中の 1 つの故障はデータ格納装置が相変わらず最大限のパフォーマンスにて動作できるようにする。N V M e S S D の各々が 4 つの P C I e レーンの帯域幅を要求するので、データ格納装置の全体セット 3 0 の最大限のパフォーマンスを維持するために、6 個のデュアルポート N V M e S S D のセットの各々は $6 X 4 = X 2 4$ すなわち 2 4 レーンの帯域幅を必要とする。

【 0 0 5 4 】

図 2 に図示したように、データ格納装置セット 30 は P C I e スイッチ 40 と連結される。各スイッチ 40 は P C I e 連結を通じて対応するホストサーバーノード 20 と連結される。図 2 の実施形態において、P C I e スイッチとサーバーノード 20 のホストプロセッサ 24 との間に X 3 2 リンクが存在する。

【 0 0 5 5 】

図 2 に図示した構成において、サーバーノード 20 A、20 B は 1 2 個の S S D 0 0 乃至 1 1 (第 1 データ格納装置セット 30 A、第 2 データ格納装置セット 30 B に該当する 2 個の S S D セット) にデュアルポートアクセスを提供する H A ピアノードである。同様に、サーバーノード 20 C、20 D は他の 1 2 個の S S D 1 2 乃至 2 3 (2 個のデータ格納装置セット 30 C、30 D に該当する 2 個の S S D セット) にデュアルポートアクセスを提供する H A ピアノードである。

10

【 0 0 5 6 】

図 2 に図示した P C I e ファブリック 10 はホストリンク 14 をさらに含む。ホストリンク 14 はスイッチ 40 及び資源 (データ格納装置 32) の間の資源連結 16 に追加的に、4 つのサーバーノード 20 からスイッチ 40 への総 1 2 8 個の P C I e 3 . 0 レーン (各ノード当たり 3 2 レーン) を連結する。ホストリンク 14 の各々は複数の連結を含む。例えば、デュアル - プロセッササーバーノードの場合、各サーバーノードからの 3 2 個のレーンはサーバーノードの第 1 C P U ソケットと連結された 1 6 個のレーン及びサーバーノードの第 2 C P U ソケットと連結された 1 6 個のレーンを含む。また、サーバーノード 20 の各々はサーバーノードを複数のスイッチ 40 と連結する複数のホストリンク 14 を含む。図 2 の実施形態に図示したように、資源連結 16 は 2 4 個のデュアルポート N V M e S S D にわたった総 1 9 2 個のレーンを含む。この時、各 S S D は X 4 レーンの帯域幅 (最大 4 G B / s) の処理容量を有する。各 S S D は処理容量 (S S D 当たり X 4 レーン) に比べて 2 倍多いポート連結 (S S D 当たり 2 X 4 レーン) を含む。2 4 個のドライブにわたった総 9 6 個のレーンに対する S S D 当たり X 4 レーンの全処理量 (f u l l t h r o u g h p u t) を達成するために、各サーバーノード 20 は X 3 2 レーンの処理容量の中の X 2 4 レーンの帯域幅を提供する。

20

【 0 0 5 7 】

一部の状況において、デュアルポート N V M e S S D の各ポートは S S D の全処理容量 (f u l l t h r o u g h p u t c a p a b i l i t y) より低い帯域幅を有する。例えば、一部デュアルポート N V M e S S D は 2 つの X 2 ポートのみを含み、これは各ポートは X 2 レーンのみを支援することを意味する。結果的に、S S D の最大処理量は X 4 レーンであると仮定すれば、ポートの中の 1 つが故障であるか、又はポートの中の 1 つと連結されたスイッチが故障の場合、S S D は X 2 レーン (即ち、S S D の処理容量の半分) のみの連結が可能である。

30

【 0 0 5 8 】

1 つのサーバーノードが故障の状況において、9 6 G B / s の持続可能な帯域幅を支援するために、ファブリック 10 のクロス接続 12 はノード 20 A、20 C、ノード 20 A、20 D、ノード 20 B、20 C、及びノード 20 B、20 D の各々の間に 2 4 G B / s / (N - 1) = 8 G B / s のクロス接続帯域幅 (c r o s s - c o n n e c t i o n b a n d w i d t h) を提供する。

40

【 0 0 5 9 】

P C I e スイッチ 40 によって提供されるレーンの最小個数は連結された構成要素、即ち 2 つのデータ格納装置セット 30 (例えば、第 1 P C I e スイッチ 40 A は第 1 データ格納装置セット 30 A のデータ格納装置及び第 2 データ格納装置セット 30 B のデータ格納装置と連結される)、ホストプロセッサ 24、及びファブリック 10 のクロス接続 12 の必要条件に依存する。図 2 に図示した実施形態において、全体 2 4 + 2 4 + 3 2 + 8 + 8 = 9 6 レーンに対して、データ格納装置セット 30 の各々は 2 4 レーンを必要とし、ホストプロセッサ 24 は 3 2 レーンを必要とし、ファブリック 10 への 2 つのクロス接続 12 の各々は 8 レーンを必要とする。図 2 に図示した実施形態において、P C I e スイッチ 4

50

0の各々はX96スイッチすなわち96レーンスイッチである。しかし、本発明の実施形態がこれに限定されることはなく、スイッチは連結された構成が必要とする数より多いレーンを含む。このような特定実施形態において、PCIeスイッチは96レーン以上を含む。

【0060】

<ノード故障における帯域幅リバランシング (Rebalancing bandwidth under node failure)>

図3は本発明の一実施形態に係るスイッチの故障及びデータ伝送トラフィックのリバランシング (rebalancing) を示すブロック図である。サーバーノード及びスイッチの間の予備連結 (redundant connections) がないので、スイッチ40の故障状況において、スイッチ40と連結されたサーバーノード20はシステムの残りの部分との連結が切断される。したがって、本発明の一部の実施形態において、スイッチ40の故障はそれと連結されたサーバーノードの損失を実質的に発生させる。しかし、本発明の実施形態がこれに限定されることはなく、一部の実施形態において、サーバーノードは複数のスイッチと連結される。

10

【0061】

図3の実施形態において、故障のスイッチ40Bと連結されたサーバーノード20Bはデータ格納システム100から連結が切断され、それによって、スイッチ40B及びサーバーノード20Bによって管理される作業負荷又はデータトラフィックが残りの(N-1)個のスイッチ40A、40C、40D及び(N-1)個のサーバーノード20A、20C、20Dに分散される。故障のスイッチ又はノード20Bと連結されたデータ格納装置セット30に/からのデータ伝送は重複して連結されたスイッチ40Aを経由する。このような作業負荷のリバランシングはシステムの24個SSD全体の処理量(SSD当たりX4リンク)を維持する。

20

【0062】

システムに内装された超過容量により、故障のスイッチと連結されたドライブセットに/からのデータ伝送のパフォーマンスが維持される。特に、帯域幅の一部はHAペアの生存メンバーに直接連結されたサーバーノードからもたらされ、帯域幅の残る部分はファブリック10を通じて連結された他のサーバーノード(例えば、20A、20C、20D)によって提供される。また、残りの(N-1)個のサーバーノード(例えば、20A、20C、20D)は追加的な負荷を収容するように、各々のスイッチ(例えば、40A、40C、40D)との連結により十分な帯域幅を有する。

30

【0063】

本発明の一部の実施形態において、適応的多重経路ファブリック10のスイッチ40はこのようなノード故障を自動的に検出し、続いてSSDエンドポイントをサーバーノード20に自動的に再割りし、生存サーバーノードを通じて帯域幅をリバランシングするようにプログラムされる。言い換えれば、スイッチ40はスイッチの現在の構成に基づいて各SSDにどのように連結されるかに対する情報を維持し、エラー条件に対して物理及びリンク階層により他のスイッチ40、データ格納装置セット30、及びサーバーノード20の間の連結をモニターリングし、このようなエラーを管理システム(例えば、サーバーノード20の中の1つ又は他の専用管理プロセッサ)に報告する。管理システムは報告されたエラーに基づいてリンク又はサーバーノードが故障であるか否かを判別し、SSDのSSDエンドポイントをサーバーノード20の中の到達できるノードに再割りするようにスイッチ40を再構成する。図3に図示した1つのノード故障による帯域幅の低下はないが、ノードの間のクロス接続12はPCIeスイッチの1つの追加的なレベルを経由し、これによって追加的な遅延が発生する。しかし、PCIeスイッチを通じた遅延はSSDに/からのデータアクセスの全体遅延と比較して一般的に小さく、無視できる。

40

【0064】

図4は本発明の一実施形態に係るサーバーノードの故障及びデータ伝送のリバランシング (rebalancing) を示すブロック図である。図4を参照すれば、サーバーノード

50

ド 20B が故障であるが、対応するスイッチ 40B は生存した状況であって、データは機能スイッチ 40B を通じて相変わらず、経由できるが、3つの生存サーバーノード 20A、20C、20Dのみを通じてアクセスされる。このような状況において、適応的多重経路ファブリックの管理システム（例えば、サーバーノード 20 の中の 1 つ又は他の専用管理プロセッサ）はサーバーノード 20B の故障を自動的に検出するようにプログラムされ、図 3 に図示した実施形態のように、SSD エンドポイントを生存エンドポイント（surviving endpoints）に自動的に再割当する。本発明の他の実施形態において、エラーの検出及びスイッチ 40 の自動的再構成はスイッチ自体により（例えば、スイッチ 40 に集積された処理ユニットによって）具現される。

【0065】

本発明の一部の実施形態において、個別のスイッチは 2 個のデータ格納装置セット 30 が最大限のパフォーマンス（full performance）にて動作するのに十分な帯域幅を提供しなくともよい。例えば、スイッチ 40B が故障の場合、データ格納装置セット 30A、30B への唯一の経路は生存スイッチ 40A を通じる経路である。生存スイッチ 40A が X96 レーンより小さいレーンを含む場合、データ格納装置セット 30A、30B は最大限のパフォーマンスにて動作するのに十分な帯域幅を有さない。しかし、スイッチ 40B が故障ではなく、関連されたサーバーノード 20B のみが故障の場合、スイッチ 40B が正常サーバーノード 20C、20D にデータを再ルーティングすることに参加できる。このような一部の実施形態において、データ格納装置セット 30 の全部が最大限のパフォーマンスにて続いて動作できる十分な帯域幅が提供される。

【0066】

< ノード間通信（Inter-node communication） >

本発明の一部の実施形態において、補助ファブリック 50（secondary fabric）がサーバーノード 20 の間の通信のために含まれる。図 5 は本発明の一実施形態に係るサーバーノード 20 を連結する補助ファブリック 50 をさらに含むシステムのブロック図である。補助ファブリック 50 は補助スイッチ 54（例えば、54A、54B、54C、及び 54D）を他の 1 つに連結する相互連結 52（inter-connections）（例えば、52AB、52AC、52AD、52BC、及び 52BD）を含む。この時、補助スイッチ 54 の各々は対応する 1 つのサーバーノード 20 と直接連結される。例えば、補助スイッチ 54A はサーバーノード 20A と直接連結される。ファブリック 10 と同様に、補助ファブリック 50 は印刷回路基板上のパターン（trace）、複数の電気配線（例えば、リボンケーブル、mini-SASHD ケーブル、OCuLink ケーブル等）、及びそれらの組合せのような多様な方式により具現される。本発明の一実施形態によれば、補助ファブリック 50 はサーバーノード CPU NTB（non-transparent bridge）ポートと連結される。補助ファブリック 50 はサーバーノード 20 の間のメタデータを同期化するのに使用され、サーバーノード 20 の間の低遅延内部通信（low-latency internal communication）を提供する。

【0067】

< ファブリックの漸進的な拡張（Incrementally expanding the fabric） >

図 1、図 2、図 3、図 4、及び図 5 は 4 つのデータ格納装置セットへの高可用性及び高性能アクセスを提供するサーバーノード 20 の 2 つのペア 22 を具備するシステムを示しているが、本発明の実施形態がこれに限定されることではない。

【0068】

本発明の実施形態は特定アプリケーションの作業負荷の必要条件によってデータ格納装置、スイッチ、及びサーバーノードの構成を含む。

【0069】

図 6 は本発明の一実施形態に係る 1 つのデータ格納装置セット 30A 及び 2 つのサーバーノード 20A、20B を含む小さい高可用性構成を示すブロック図である。図 6 の構成は

10

20

30

40

50

1つのデータ格納装置セット30Aのデータ格納容量が現在作業負荷に対して十分な場合、及び予備サーバーノードを通じた高可用性が適切な場合に有用である。

【0070】

図6を参照すれば、1つのデータ格納装置セット30Aは第1スイッチ40A及び第2スイッチ40Bの両方と連結される。図1の実施形態と同様に、第1スイッチ40Aは第1サーバーノード20Aと連結され、第2スイッチ40Bは第2サーバーノード20Bと連結される。第1サーバーノード20A又は第2サーバーノード20Bの中のいずれか1つが故障の場合、データ格納装置セット30Aは生存ノードを通じてアクセスを維持する。

【0071】

データ格納需要が増加する場合、追加的なデータ格納装置セットが図6のシステムに追加できる。例えば、第1スイッチ40A及び第2スイッチ40Bの全てに追加データ格納装置を連結することによって、追加的な1つのデータ格納装置セットが追加されて、他のスイッチへの相互連結無しに、第1サーバーノード20A、第2サーバーノード20B、第1スイッチ40A、第2スイッチ40B、第1データ格納装置セット30A、及び第2データ格納装置セット30Bと実質的に類似に構成される。前述のように、サーバーノード20A、20Bの中の1つ又はスイッチ40A、40Bの中の1つが故障の場合、第1及び第2データ格納装置セット30A、30Bの全てはアクセス可能であるように維持される。

【0072】

図7は本発明の一実施形態に係る追加的なデータ格納装置セットを含む図6に図示した構成からの拡張を示す。図7に図示したように、第2データ格納装置セット30Bは第1スイッチ40A及び第2スイッチ40Bと連結される。図7の構成は追加的なサーバーノードの代わりに第3スイッチ40C及び第4スイッチ40Dの追加的なクロス接続をさらに含む。第3データ格納装置セット30Cは第3及び第4スイッチ40C、40Dと連結され、第4データ格納装置セット30Dは第3及び第4スイッチ40C、40Dと連結される。クロス接続12AC、12BCは第3スイッチ40Cを第1及び第2スイッチ40A、40Bと連結し、クロス接続12AD、12BDは第4スイッチ40Dを第1及び第2スイッチ40A、40Bと連結する。結果的に、第1及び第2サーバーノード20A、20Bは第3及び第4スイッチ40C、40Dを通じて第3及び第4データ格納装置セット30C、30Dをアクセスする。

【0073】

図7に図示した構成において、サーバーノード20A、20Bの中の1つが故障の場合、又は4つのスイッチ40A、40B、40C、40Dの中のいずれかが故障の場合において、データ格納装置の全部がアクセス可能であるように維持される。しかし、第3及び第4データ格納装置セット30C、30Dのデータ処理パフォーマンス(data throughput performance)はクロス接続12の帯域幅によって制限され、サーバーノード20及びその対応するスイッチ40の間のホストリンク14によって制限される。特に、図7の配列において、2つのサーバーノード20は4つのデータ格納装置セット30の全部に総16Bの帯域幅(各サーバーノード20毎に8B)を提供し、これはクロス接続12AC、12AD、12BD、12BDのみを通じて連結されるデータ格納装置セット30C、30Dへの可用である最大帯域幅と対応する、データ格納装置セット当たり4Bを意味する。言い換えれば、クロス接続当たり2Bであり、4つのクロス接続は2つのドライブセットと共有される総8B(ドライブセット当たり4B)を提供する。また、サーバーノードの中の一部又はスイッチの中の一部の故障はシステムのデータ処理量に追加的な影響を及ぼす。このような意味において図7の構成は、例えば十分な帯域幅が総格納容量より重要でない場合にさらに適合する。

【0074】

追加的な帯域幅が要求される場合、図7の構成はサーバーノードを第3及び第4スイッチ40C、40Dに連結するように拡張され、これにしたがって図2に図示したのと実質的に同一な構造になる。このような意味において、本発明の実施形態は使用者の増加する要

10

20

30

40

50

求条件を対応するために必要に応じてシステムの漸進的な拡張を許容する。

【0075】

本発明のその他の実施形態において、使用者によって要求される作業負荷はデータ集中 (data-intensive) より演算集中 (compute-intensive) による。図8は本発明の一実施形態に係る追加的なサーバーノードを含む図6に図示した構成からの拡張を示す。

【0076】

図6に図示したように6個のデータ格納装置の1つのセット30A及び2つのサーバーノードの基本構成に、第3サーバーノード20Cがクロス接続12AC、12BCを通じて第1及び第2スイッチ40A、40Bと連結された第3スイッチ40Cと共に追加される。又は、第3サーバーノード20Cは追加的なスイッチの代わりにパッシブ相互連結ボード (passive interconnect board) (例えば、サーバーノード20C及びスイッチ40A、40Bの間の効率的な直接連結) を通じて第1及び第2スイッチ40A、40Bと連結される。これは同一のデータ格納装置セットへのアクセスを維持しながら、システムの演算能力を向上させる。演算要求がさらに増加する場合、第4ノード20Dが第4スイッチ40D又はパッシブ連結ボードの中の1つを通じて第1及び第2スイッチ40A、40Bと連結される。

10

【0077】

本願の使用事例は、2つのサーバーノード20A、20B上において駆動するソフトウェアスタックがデータ格納装置32への最高帯域幅を達成する能力に影響を及ぼす隘路現象である場合であり、この場合、さらに多くのサーバーノードの追加がさらに多くのサーバー演算能力を提供する。追加的なストレージが要求される場合、追加的なデータ格納装置セットが図7に図示したのと類似な方式により漸進的に追加されて図2に図示したように24個のデータ格納装置及び4つのサーバーノードの構造が達成される。

20

【0078】

説明を簡易にするために、本発明の実施形態が4個以下のスイッチを含む構造により説明した。しかし、本発明の範囲がこれに限定されることではない。例えば、本発明の一部の実施形態は4個以上のスイッチを含む。

【0079】

類似の方式を使用する場合、本発明の実施形態に係る適応的多重経路ファブリックを含むデータ格納システム100の変形は、例えば、144GB/sの全体処理量に対して、6個のスイッチを通じて36個のデータ装置と連結された6個のノードを含む。他の例として、図9は本発明の一実施形態に係る192GB/sの全体処理量 (各々約4GB/sを有するデータ格納装置に対する上述した仮定に基づく) のために8個のスイッチ40A、40B、40C、40D、40E、40F、40G、40Hを通じて48個のデータ格納装置と連結された8個のサーバーノード20A、20B、20C、20D、20E、20F、20G、20Hを含むシステムを示すブロック図である。

30

【0080】

このように、本発明の実施形態に係る適応的多重経路ファブリック構造はサーバーノードのグループ及び高性能マルチポートNVMe SSDのグループを使用して均衡を成し、構成可能なエンドツーエンドシステムパフォーマンスプロファイル (balanced and configurable end-to-end system performance profile) を提供する。パフォーマンスはサーバーノード及びエンドポイントを通じて均衡を成し、パフォーマンスは1つのノードが故障しても維持され、均衡をなす。

40

【0081】

データ格納システムは、例えばイーサネット (登録商標) 連結を通じてSANを提供し、多重経路ファブリックはサーバーノードへの基本ドライブの間のパフォーマンスに相応しいイーサネット (登録商標) を通じてネットワークストレージパフォーマンスを提供する。

50

【 0 0 8 2 】

例えば、システムパフォーマンスプロファイルは100GB/s 使用者データ処理量と、サーバーノードの間のローレイテンシ内部通信と、エンドポイント上のデータへの高可用性アクセスとの中の1つ以上の組合せである。

【 0 0 8 3 】

本発明の実施形態は適応的多重経路ファブリック構造を提供する。適応的多重経路ファブリック構造はサーバーノードのグループ(P C I e ルート - コンプレックス)をマルチポートSSDのグループと連結し、

サーバーノードの個数、SSDの個数、エンドツーエンドパフォーマンス規定(e n d - t o - e n d p e r f o r m a n c e s p e c i f i c a t i o n)に関して柔軟性(f l e x i b i l i t y)を提供し、

サーバーノード及びSSDを通じてエンドツーエンド負荷均衡を支援し、SSDのマルチポートを通じて一对のサーバーノード(H A - p e e r s)からSSDのセットへの予備アクセスを提供し、

様々なペアのHAピアノードの間のクロス接続帯域幅を提供して故障復旧及び負荷均等化シナリオによりすべてのサーバーノードに帯域幅をリバランシングし、すべての生存ノードによって帯域幅をリバランシングすることによってエンドツーエンドシステムパフォーマンスの低下無しに1つのノード故障に耐え、自動故障検出及びその後のサーバーノードへのSSDエンドポイントの再割当及び帯域幅のリバランシングに転じる能力を提供する。

【 0 0 8 4 】

本発明の実施形態は、柔軟な拡張可能な方式によりサーバーノードのようなP C I e ルート - コンプレックスをデュアルポートN V M e S S DのようなマルチポートP C I e エンドポイントと連結する機能と、

1つのルート - コンプレックス故障(H A)の状況においてパフォーマンス低下を制限しながら、すべてのP C I e エンドポイントを継続的にアクセスする機能と、ルート - コンプレックス及びエンドポイントの間の帯域幅の割当を調整する機能と、

故障復旧又は負荷バランシングシナリオにより、P C I e エンドポイントをルートコンプレックスに動的に再割当する機能と、

ホストルート - コンプレックス上のオーバーヘッド無しにエンドポイントの間のピアツーピアデータ伝送を遂行する機能と、HAピアサーバーノードの間のローレイテンシの高い処理量通信を遂行する機能を可能とする。

【 0 0 8 5 】

本発明を特定例示的な実施形態と関連して説明したが、本発明が記載した実施形態に限定されず、特許請求範囲の範囲及びその思想内に含まれる同等の配列及び多様な変形及びそれらの均等物を含むと意図する。

【 0 0 8 6 】

例えば、本発明の実施形態はP C I e スイッチを通じてサーバーノードに連結されるP C I e ポートを含むソリッドステートドライブとしてコンピューティング資源に関して説明したが、本発明の実施形態はこれに限定されない。例えば、本発明の実施形態において、ソリッドステートドライブは2以上のポートを含むが、代わりに多様な複数のエンドポイントポート(例えば、2以上のポート)を含むSSDの使用も含む。また、本発明の一部の実施形態において、サーバーノードは複数のポートを通じてファブリックに連結される。例えば、サーバーノードは多数のホストリンクを通じて1つ以上のスイッチと連結され、それによって、ホストリンク又はスイッチが故障の状況において、サーバーノード及びコンピューティング資源の間の予備連結を提供する。ホストリンク及びスイッチの帯域幅に応じて、このような予備連結は、故障状態においても、サーバーノードが最高帯域幅にて動作するようにする。

【 0 0 8 7 】

一部の実施形態によれば、ソリッドステートドライブはイーサネット(登録商標)、I B (I n f i n i b a n d (登録商標))、F C (F i b r e C h a n n e l)、S A

10

20

30

40

50

S (s e r i a l l y a t t a c h e d S C S I) 等の他のインターフェイスを使用する。例えば、イーサネット（登録商標）インターフェイスの場合に、P C I e スイッチはネットワーク（イーサネット（登録商標））スイッチに交替される。

【 0 0 8 8 】

一部の実施形態によれば、適応的多重経路ファブリックを使用して連結された資源はハードディスクドライブ、テープドライブ、D R A M のような揮発性メモリの他の形態のコンピューティング資源、及びベクトルプロセッサ、G P U、D S P、F P G A のようなコンピューティングユニットである。

【 0 0 8 9 】

本発明のその他の実施形態において、個別的なコンピューティング資源の各々は複数のポートを含むことを必要としない。例えば、各コンピューティング資源セットは個別的なコンピューティング資源の各々へのリンク及び2つの連結されたスイッチ40にリンクを提供するアダプター又はスイッチを含む。さらに詳細な例として、データ格納装置32の各々は各データ格納装置32がデータ格納ストレージ装置のセットと関連されたアダプターと連結される単一ポートデータ格納装置である。アダプターは2つのスイッチ40と連結される。このような方式により、個別的なデータ格納装置の各々が単一ポート装置であっても、データ格納装置セット30は複数のスイッチ40に相変わらず連結される。このような技法は上述したように他の形態のコンピューティング資源に適用できる。

【 0 0 9 0 】

説明を簡易にするために、インターフェイスの相対的な帯域幅は図示した実施形態において同一であるが（例えば、データ格納装置の各セットに対して6B、各サーバーノード及び対応するスイッチの間の連結に対して8B、及びスイッチの間の各クロス接続に対して2B）、本発明の実施形態がこれに限定されることはなく、本発明の実施形態は他の帯域幅（例えば、インターフェイスと関連されたコンピューティング資源の処理量と対応する帯域幅）を提供するインターフェイスを含む実施形態をさらに含む。1つのサーバーノードの故障状況において、適応的多重経路ファブリックと連結された資源の最大限のパフォーマンスを維持するように設計された本発明の一部の実施形態において、クロス接続の全体帯域幅（又はクロス接続帯域幅）は1つのノードによって一般的に提供される帯域幅と少なくとも同一であって十分である。本発明の一部の実施形態が本文において、クロス接続の全部が同一の帯域幅を有し、H A ペアが外部のすべてのスイッチに連結されると説明したが、本発明の実施形態がこれに限定されることではない。例えば、一部の実施形態において、クロス接続は他の帯域幅を有し、一部の実施形態において、クロス接続は他のスイッチの全部より少ない数にて形成される。

【 0 0 9 1 】

説明を簡易にするために、8個以下のホストを含む実施形態を説明したが、本発明の実施形態がこれに限定されることはなく、類似な概念が多様な数のホストにより具現される。

【 0 0 9 2 】

同様に、本発明の実施形態が正確に6個のデータ格納装置のセットに限定されることはなく、各セットにおいて多数のS S D（例えば、各セットにおいて同一の数又は各セットにおいて多様な数）を含む実施形態をさらに含む。

【 0 0 9 3 】

本発明の実施形態は、ファブリック故障検出及び再構成機能を使用してストレージ管理ツールと通信してストレージ基盤施設管理を向上させる。

【 符号の説明 】

【 0 0 9 4 】

- 10 ファブリック
- 12 クロス接続
- 14 ホストリンク
- 16 資源リンク
- 20 サーバーノード

10

20

30

40

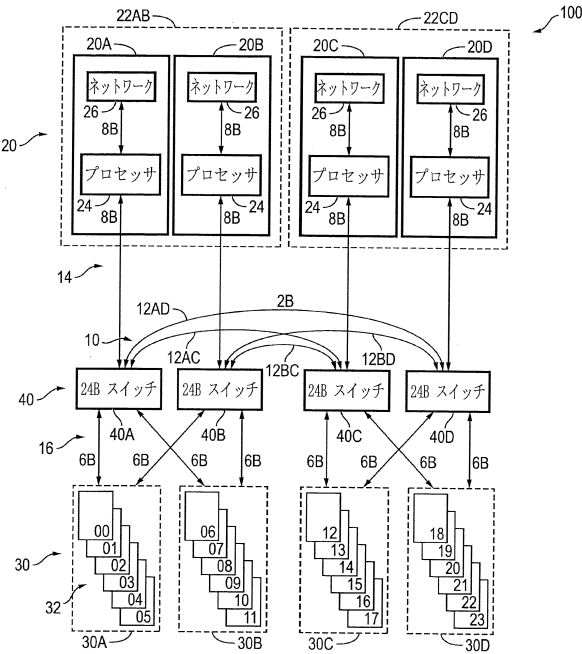
50

- 2 2 ペア
- 2 4 ホストプロセッサ
- 2 6 ネットワークインターフェイス
- 3 0 データ格納装置セット
- 3 2 データ格納装置
- 4 0 スイッチ
- 5 0 補助ファブリック
- 5 2 相互連結
- 5 4 補助スイッチ
- 1 0 0 データ格納システム

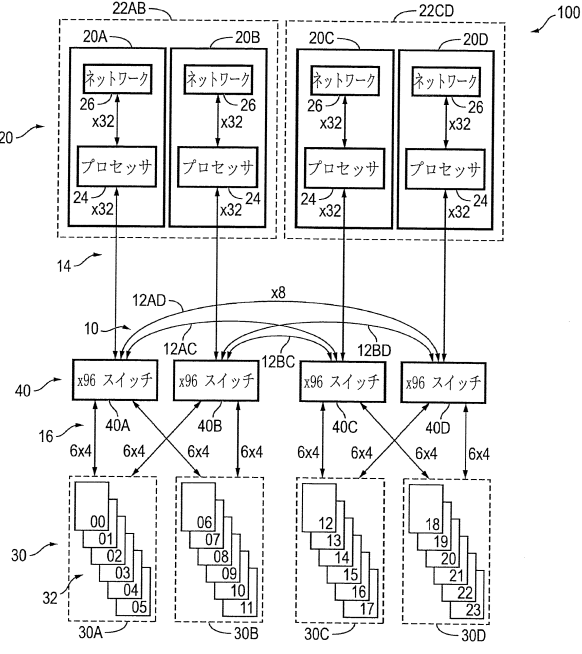
10

【図面】

【図 1】



【図 2】



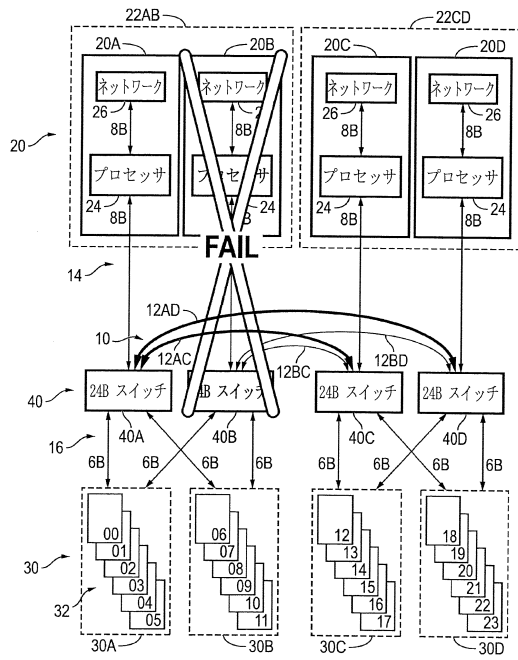
20

30

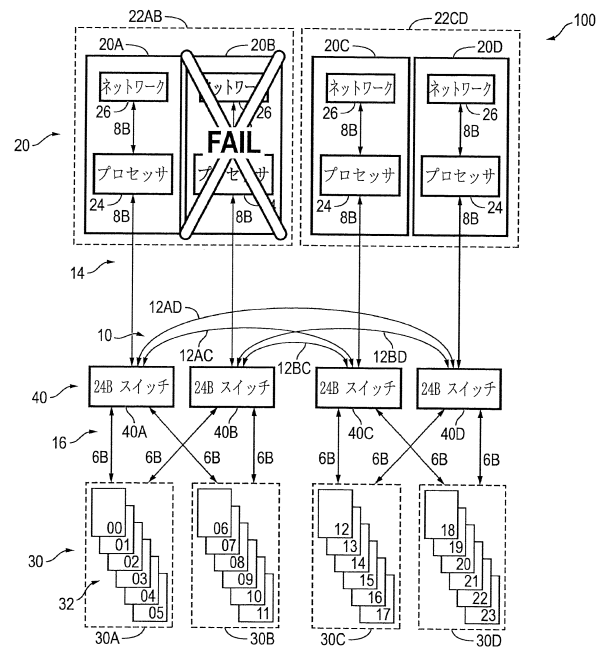
40

50

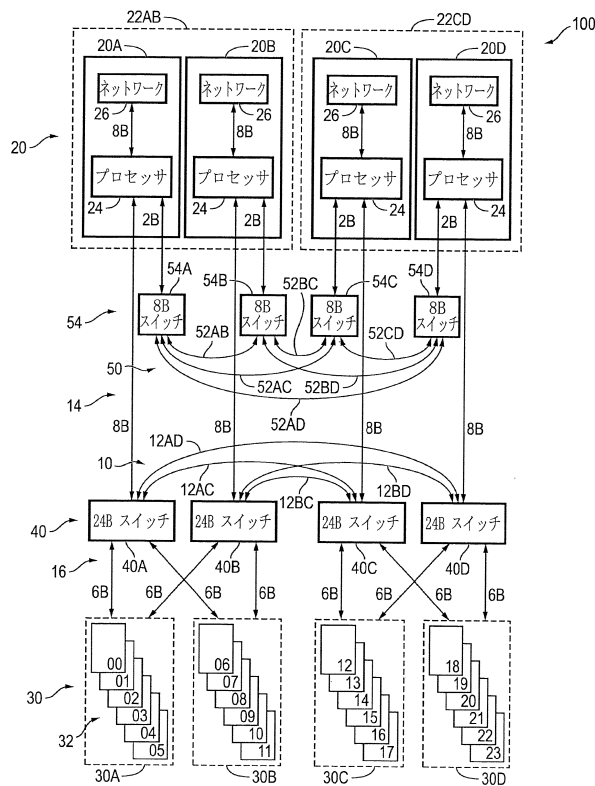
【図 3】



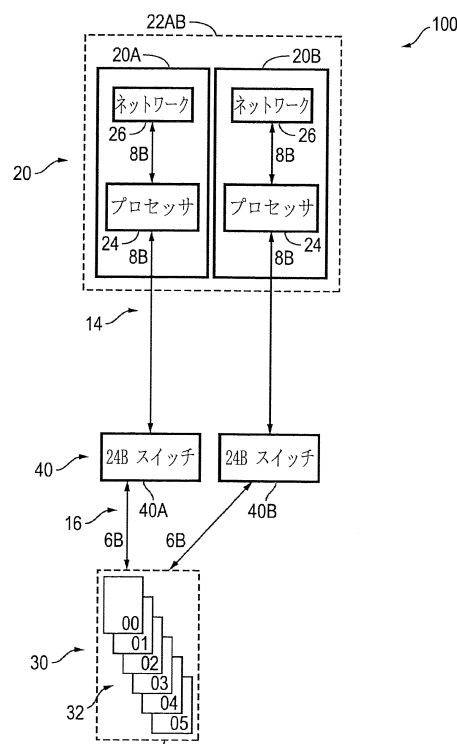
【図 4】



【図 5】



【図 6】



10

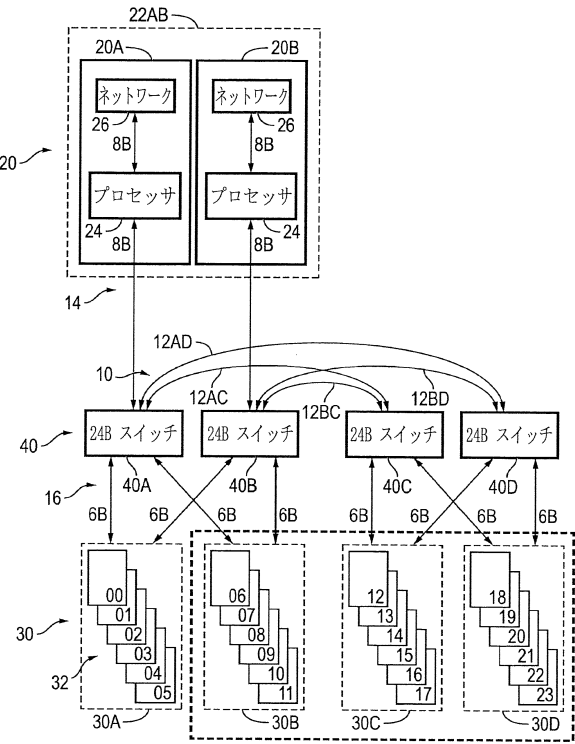
20

30

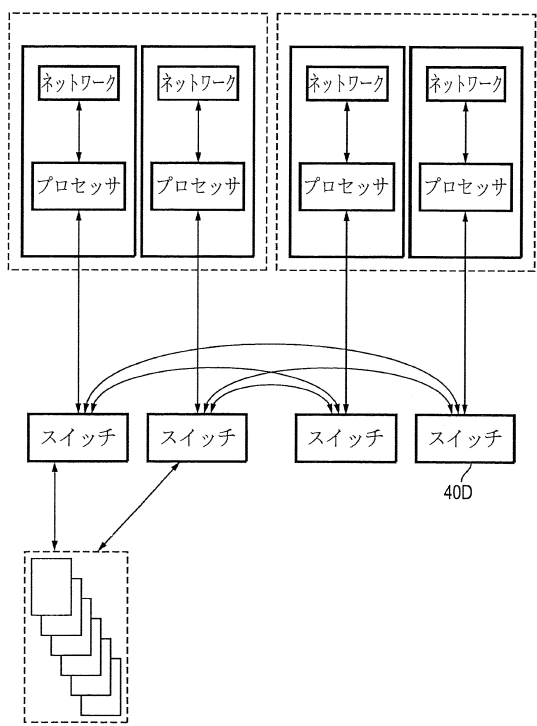
40

50

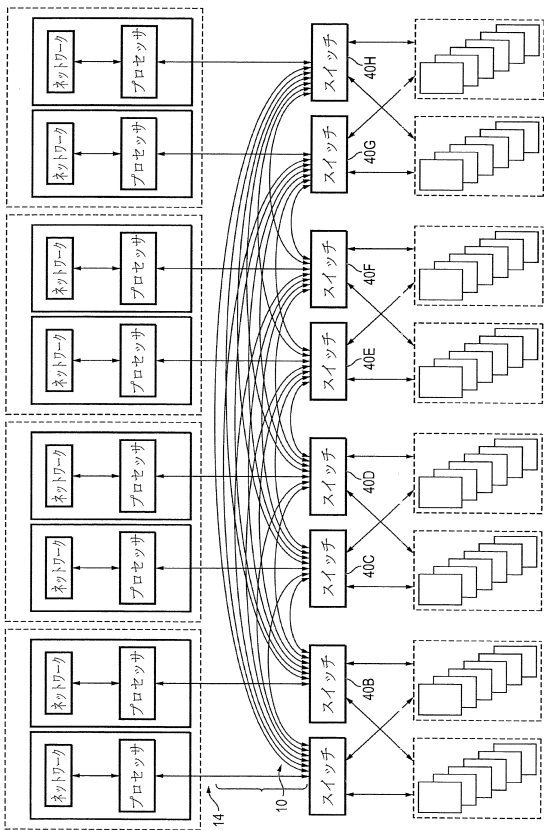
【図 7】



【図 8】



【図 9】



10

20

30

40

50

フロントページの続き

3 1 4

(72)発明者 ステファン ジー . フィッシャー

アメリカ合衆国 9 4 0 4 1 カリフォルニア州 マウンテン ビュー 2 6 3 ビラ エステー アパ
ート 1 6 0 0

(72)発明者 ジャン ピン

アメリカ合衆国 9 5 0 3 5 カリフォルニア州 ミルピタス ガルシア コート 3 1 6

(72)発明者 インディラ ジョシ

アメリカ合衆国 9 5 0 7 0 カリフォルニア州 サラトガ ボニー リッジ ウェイ 1 9 9 5 0

(72)発明者 ハリー ロジャース

アメリカ合衆国 9 5 1 2 5 カリフォルニア州 クレストフィールド ドライブ 1 3 3 5

審査官 大石 博見

(56)参考文献 米国特許出願公開第 2 0 0 9 / 0 2 0 4 7 4 3 (U S , A 1)

(58)調査した分野 (Int.Cl. , D B 名)

H 0 4 L 4 5 / 2 8

H 0 4 L 4 7 / 1 2 5

G 0 6 F 1 3 / 0 0