

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4414399号  
(P4414399)

(45) 発行日 平成22年2月10日(2010.2.10)

(24) 登録日 平成21年11月27日(2009.11.27)

(51) Int.Cl. F I  
**G06F 3/06 (2006.01)** G06F 3/06 305F  
 G06F 3/06 301G

請求項の数 5 (全 20 頁)

<p>(21) 出願番号 特願2006-20090 (P2006-20090)                  (22) 出願日 平成18年1月30日(2006.1.30)                  (65) 公開番号 特開2007-200171 (P2007-200171A)                  (43) 公開日 平成19年8月9日(2007.8.9)                  審査請求日 平成18年12月25日(2006.12.25)</p>	<p>(73) 特許権者 000005223                  富士通株式会社                  神奈川県川崎市中原区上小田中4丁目1番                  1号                  (74) 代理人 100087848                  弁理士 小笠原 吉義                  (74) 代理人 100083297                  弁理士 山谷 皓榮                  (72) 発明者 佐藤 弘章                  神奈川県川崎市中原区上小田中4丁目1番                  1号 株式会社富士通コンピュータテクノ                  ロジーズ内                  審査官 菅原 浩二</p>
---	---

最終頁に続く

(54) 【発明の名称】 ディスク制御装置

(57) 【特許請求の範囲】

【請求項1】

ホストプロセッサとのインタフェースを持つ1または複数のチャンネルアダプタと、ディスク装置とのインタフェースを持つ1または複数のデバイスインタフェースと、1または複数の通信用の転送回路とをそれぞれ有し、ホストプロセッサとディスク装置との間のデータ転送を制御する複数のコントローラモジュールを備え、かつ前記複数のコントローラモジュールは、前記転送回路を用いた通信経路によって互いに通信する機能を持つディスク制御装置であって、

前記コントローラモジュールの少なくとも一つは、

前記通信経路上で異常が生じた場合に、その異常が生じた部位を閉塞し、その部位の閉塞によって前記複数のコントローラモジュール間で通信ができなくなった場合に、通信ができなくなったコントローラモジュールのうち、二重化されていないデータを保持するコントローラモジュールを二重化されているデータを保持するコントローラモジュールよりも優先的に切り離し対象として決定する経路閉塞判定手段と、

前記経路閉塞判定手段によって切り離し対象として決定されたコントローラモジュールを使用不可の状態に設定する切り離し制御手段とを備える

ことを特徴とするディスク制御装置。

【請求項2】

請求項1記載のディスク制御装置において、

前記経路閉塞判定手段は、二重化されているデータを保持しているか否かにより切り離

10

20

し対象のコントローラモジュールを決定できなかった場合に、前記通信経路上において閉塞した部位の位置とその数によって切り離し対象のコントローラモジュールを決定することを特徴とするディスク制御装置。

【請求項 3】

請求項 2 記載のディスク制御装置において、

前記経路閉塞判定手段は、前記通信経路上において閉塞した部位の位置とその数によって切り離し対象のコントローラモジュールを決定できなかった場合に、さらにあらかじめ定められたマスターとなる一つのコントローラモジュール以外のコントローラモジュールを優先的に切り離し対象として決定する

ことを特徴とするディスク制御装置。

10

【請求項 4】

請求項 1、請求項 2 または請求項 3 記載のディスク制御装置において、

前記経路閉塞判定手段は、自分のコントローラモジュールが切り離し対象の候補となった場合に、前記通信経路を用いた通信とは異なる通信手段によって、他の切り離し対象の候補となっているコントローラモジュールの生存を確認し、生存が確認できなかった場合には、その生存が確認できなかったコントローラモジュールを切り離し対象として決定する

ことを特徴とするディスク制御装置。

【請求項 5】

ホストプロセッサとのインタフェースを持つ 1 または複数のチャンネルアダプタと、ディスク装置とのインタフェースを持つ 1 または複数のデバイスインタフェースと、複数の通信用の転送回路とをそれぞれ有し、ホストプロセッサとディスク装置との間のデータ転送を制御する複数のコントローラモジュールを備え、かつ前記複数のコントローラモジュールは、前記転送回路と該転送回路を接続するスイッチング機能を持つ経路設定装置を用いた複数の通信経路によって互いに通信する機能を持つディスク制御装置であって、

20

前記通信経路上で異常が生じた場合に、その異常が生じた部位を閉塞し、その部位の閉塞によって前記複数のコントローラモジュール間で通信ができなくなった場合に、その部位の閉塞によってすべての通信経路が閉塞になったコントローラモジュールを切り離し対象として決定し、これにより決定できない場合に、前記通信できなくなったコントローラモジュールのうち、二重化されていないデータを保持するコントローラモジュールを二重化されているデータを保持するコントローラモジュールよりも優先的に切り離し対象として決定する経路閉塞判定手段と、

30

前記経路閉塞判定手段によって切り離し対象として決定されたコントローラモジュールを使用不可の状態に設定する切り離し制御手段とを備える

ことを特徴とするディスク制御装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ディスク制御装置に関し、特にディスク制御装置内の複数のコントローラモジュール（以下、CM という）間でのデータ転送経路に異常が発生した場合に、システムに影響が最小となるように閉塞する CM を決定し閉塞させるようにした経路閉塞判定機能を持つディスク制御装置に関するものである。

40

【背景技術】

【0002】

図 19 は、従来のディスク制御装置の例を示す。図 19 において、200 はディスク制御装置、201 はそれぞれホストプロセッサとディスク装置との間のデータ転送制御やエラー処理等を行うコントローラモジュール（CM）、202 はスイッチング機能を持つ経路設定装置、203 はホストプロセッサとのインタフェースを持つチャンネルアダプタ（CA）、204 はディスク装置とのインタフェースを持つデバイスアダプタ（DA）である。なお、経路設定装置 202 は、フロントエンドラウタ（FRT）と呼ぶこともある。

50

## 【 0 0 0 3 】

ホストプロセッサ（図示省略）からのディスク装置（図示省略）に対する入出力要求は、チャンネルアダプタ 2 0 3 から経路設定装置 2 0 2 を経て、複数の C M 2 0 1 の一つに送信され、C M 2 0 1 から経路設定装置 2 0 2 を経て、デバイスアダプタ 2 0 4 に送られ、ディスク装置への入出力が実行される。

## 【 0 0 0 4 】

図 1 9 に示すような従来のディスク制御装置 2 0 0 では、チャンネルアダプタ 2 0 3 とデバイスアダプタ 2 0 4 とが経路設定装置 2 0 2 上に載っており、各 C M 2 0 1 と、チャンネルアダプタ 2 0 3 と、デバイスアダプタ 2 0 4 とが、経路設定装置 2 0 2 を通して、随時任意に接続される構成になっているため、どこかのモジュールに故障等の障害が発生した場合、その部分だけを閉塞して切り離せば、障害の影響範囲を局所化することを比較的簡単に行うことができた。

10

## 【 0 0 0 5 】

なお、下記の特許文献 1 に記載されているディスク制御装置では、複数のチャンネルアダプタモジュール、複数のスイッチモジュールおよび一つのキャッシュモジュールで構成されるクラスタのそれぞれに異なる供給元から電源を供給し、各チャンネルアダプタモジュールまたはディスクアダプタモジュールは、特定のキャッシュモジュールとの間のデータ転送のための複数のパスを設定可能としている。この特許文献 1 に示されているディスク制御装置も、障害時におけるモジュールの閉塞、切り離しに関しては、基本的に図 1 9 に示したディスク制御装置と同様な方式を用いている。

20

【特許文献 1】特開 2 0 0 1 - 2 7 9 7 2 号公報

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【 0 0 0 6 】

図 2 0 は、本発明の課題を説明するための図である。本発明者等は、図 2 0 に示すような、従来のディスク制御装置とは異なる新しいディスク制御装置の開発に取り組んでいる。図 2 0 において、1 0、1 1 はそれぞれホストプロセッサとディスク装置との間のデータ転送制御やエラー処理等を行うコントローラモジュール（C M）、7 0、7 1 はスイッチング機能を持つ経路設定装置（F R T）、4 0、4 1 はホストプロセッサとのインタフェースを持つチャンネルアダプタ（C A）、5 0、5 1 はディスク装置とのインタフェースを持つデバイスインタフェース（D I）であり、図 1 9 のデバイスアダプタに相当するもの、6 0、6 1 は D M A チップ等によって構成されるデータ転送を行う転送回路を表す。

30

## 【 0 0 0 7 】

図 2 0 に示すディスク制御装置が、図 1 9 に示したような従来のディスク制御装置 2 0 0 と大きく異なる点は、ホストインタフェースであるチャンネルアダプタと、ディスク装置とのインタフェースであるデバイスインタフェースとが、C M 1 0、1 1 内に存在することである。これは、通常のアクセス時におけるホストプロセッサからディスク装置への C M 1 0 または C M 1 1 を介するデータ転送の経路をできるだけ単純化させ、制御の簡易化を図るとともに、チャンネルアダプタおよびデバイスインタフェースを含めた C M の低コスト化を実現するためである。

40

## 【 0 0 0 8 】

なお、C M 1 0、1 1 間のデータ転送は、必要に応じて転送回路 6 0、6 1 および経路設定装置 7 0 または 7 1 を介して行われるようになっている。

## 【 0 0 0 9 】

このように、図 2 0 に示すディスク制御装置では、チャンネルアダプタとデバイスインタフェースとが C M 1 0、1 1 内に存在するため、例えば C M 1 1 が障害により閉塞すると、そのチャンネルアダプタ 4 1 およびデバイスインタフェース 5 1 も閉塞してしまうことになり、システムへの影響が大きくなってしまいう問題がある。なお、閉塞とは、装置や回路等の特定の部位が使用不可の状態になることをいう。

## 【 0 0 1 0 】

50

本発明は上記問題点の解決を図り、CMが異常の場合でもCMが閉塞した場合の影響を考え、内部の異常部品を部分的に閉塞させることにより、部品またはCMの切り離しによる影響が小さくなるようにし、できるだけ運用の継続を可能にすることを目的とする。

【課題を解決するための手段】

【0011】

本発明は、ホストプロセッサとのインタフェースを持つ1または複数のチャンネルアダプタと、ディスク装置とのインタフェースを持つ1または複数のデバイスインタフェースと、1または複数の通信用の転送回路とをそれぞれ有し、ホストプロセッサとディスク装置との間のデータ転送を制御する複数のコントローラモジュール(CM)を備え、かつ前記複数のCMは、前記転送回路を用いた通信経路によって互いに通信する機能を持つディスク制御装置において、装置構成や部品の閉塞具合、データ保持状態かどうかなどの装置状態から判断して運用を継続するために閉塞させるCMを決定することをもっとも主要な特徴とする。

10

【0012】

すなわち、本発明に係るディスク制御装置内のCMの少なくとも一つは、CM間通信用の通信経路上で異常が生じた場合に、その異常が生じた部位を閉塞し、その部位の閉塞によって前記複数のCM間で通信ができなくなった場合に、通信ができなくなったCMのうち、二重化されていないデータを保持するCMを二重化されているデータを保持するCMよりも優先的に切り離し対象として決定する経路閉塞判定手段と、前記経路閉塞判定手段によって切り離し対象として決定されたCMを使用不可の状態に設定する切り離し制御手段とを備える。これによりCMの可用性が大きくなり、障害によるシステムへの影響が小さくなる。

20

【0013】

さらに、経路閉塞判定手段は、二重化されているデータを保持しているか否かにより切り離し対象のCMを決定できなかった場合に、前記通信経路上において閉塞した部位の位置とその数によって切り離し対象のCMを決定することもできる。これにより、使用できる可能性の小さいCMを優先的に切り離すことができる。

【0014】

またさらに、経路閉塞判定手段は、前記通信経路上において閉塞した部位の位置とその数によって切り離し対象のCMを決定できなかった場合に、あらかじめ定められたマスターとなる一つのCM以外のCMを優先的に切り離し対象として決定することもできる。これにより、マスターCMと比較して切り離しによる影響の小さいスレーブCMを優先的に切り離すことができる。

30

【0015】

また、経路閉塞判定手段は、自分のCMが切り離し対象の候補となった場合に、通信経路を用いた通信とは異なる通信手段によって、他の切り離し対象の候補となっているCMの生存を確認し、生存が確認できなかった場合には、その生存が確認できなかったCMを切り離し対象として決定することもできる。これにより、使用できる可能性の小さいCMを優先的に切り離すことができる。

【0016】

40

CM間の通信経路が各CMに付随する転送回路と該転送回路を接続するスイッチング機能を持つ経路設定装置とによって構成されるディスク制御装置では、前記経路閉塞判定手段は、異常が生じた通信経路上の部位を閉塞し、その部位の閉塞によってCM間で通信ができなくなった場合に、その部位の閉塞によってすべての通信経路が閉塞になったCMを切り離し対象として決定し、これにより決定できない場合に、通信できなくなったCMのうち、二重化されていないデータを保持するCMを二重化されているデータを保持するCMよりも優先的に切り離し対象として決定するように構成される。これにより、特に3個以上のCMを持つディスク制御装置において、効果的に切り離し対象のCMを選ぶことができるようになる。

【発明の効果】

50

## 【 0 0 1 7 】

本発明によれば、CMにおいて異常が検出された場合に、直ちにそのCMを閉塞するのではなく、内部の異常部品を部分的に閉塞させ、CM間の通信経路が存在しCM間通信ができる限りはCMを切り離さないで運用を継続する。CM間の通信経路がなくなり、CM間通信ができなくなったときには、保持しているデータが二重化されておりデータロストにならないCMを優先的に切り離し対象として選択して運用を継続する。したがって、障害によるシステムへの影響を小さくすることができるようになる。

## 【 発明を実施するための最良の形態 】

## 【 0 0 1 8 】

図1は、本発明の実施の形態に係るディスク制御装置の構成例を示す図である。図1に示すディスク制御装置1は、複数のコントローラモジュール(CM)10, 11, ...と、各CMを接続するためのスイッチング機能を持つ経路設定装置(FRT)70, 71を備える。CM10, 11は、ホストプロセッサとディスク装置との間のデータ転送制御やエラー処理等を行うものであり、それぞれホストプロセッサとのインタフェースを持つチャネルアダプタ(CA)40, 41, ディスク装置とのインタフェースを持つデバイスインタフェース(DI)50, 51およびCM間のデータ転送やメッセージ通信を行うためのDMAチップ等によって構成される転送回路60, 61を備える。

10

## 【 0 0 1 9 】

ディスク制御装置1におけるCMの数が2の場合、すなわち2CM構成の場合には、経路設定装置70, 71を設けず、一方のCM10の転送回路60と他方のCM11の転送回路61とを直接接続した構成にすることもできる。

20

## 【 0 0 2 0 】

本実施の形態では、CM10は、データ転送制御やエラー処理等を行うためのCPUを備え、CPUとプログラムとによって実現される経路閉塞判定手段20と、切り離し制御手段30とを備える。CM11も同様である。

## 【 0 0 2 1 】

経路閉塞判定手段20および切り離し制御手段30が行う処理の概要を、図2～図4に従って説明する。

## 【 0 0 2 2 】

図2は、通信経路閉塞の処理を説明する図である。なお、各CM10～13は、それぞれ2つの転送回路60～63を備え、経路設定装置70または71によって設定された通信経路を通してCM-CM間通信ができるようになっている。本実施の形態において、通信経路閉塞は、次のように行われる。

30

## 【 0 0 2 3 】

CM-CM間通信で異常を検出した場合、CM-CM間のデータ転送やメッセージ通信を行うDMAチップによって構成される転送回路を切り離すことにより、異常を検出した通信経路を使わずに、反対側の通信経路を使ってCM-CM間通信を行うことで、CMを閉塞させることなく運用を継続する。すなわち、図2(A)に示すように、CM11の転送回路61に障害が発生し、経路設定装置70を通してCM10またはCM13からのCM11への通信ができなくなった場合、CM11を閉塞するのではなく、CM11における経路設定装置70側の転送回路61を閉塞し、切り離す。これによって、CM10またはCM13からCM11への通信は、図2(B)に示すように、他の経路設定装置71を通して行われ、運用を継続させることができる。

40

## 【 0 0 2 4 】

しかし、転送回路の切り離しにより複数の転送回路が閉塞していくと、通信経路がなくなり、通信できないCMが出てきてしまう。この問題を解決するため、通信できないCMを切り離すことにより運用を継続する。図3(A)の例では、CM11における経路設定装置70側の転送回路61と、CM12における経路設定装置71側の転送回路62とが閉塞している。この場合、CM11とCM12との間は通信経路がないため、通信ができない。したがって、CM11とCM12のうち重要でないほうのCM(例えばCM12)

50

を、図3(B)に示すように閉塞する。これにより、CM11が通信できないCMが存在しなくなるため、運用の継続が可能となる。

【0025】

図4は、ディスク制御装置1が2CM構成時の場合のCM閉塞の例を示している。ディスク制御装置1が経路設定装置70、71なしで、転送回路60と転送回路61とが直結された装置構成の場合、図4(A)に示すように、CM10とCM11の転送回路60、61が片方ずつ閉塞してしまうと、CM10とCM11の通信経路がなくなってしまう。そこで、このような場合、図4(B)に示すように、CM11を切り離す。こうすることにより、ディスク制御装置1は、一つのCM10で運用を継続することが可能となる。

【0026】

なお、この場合にCM10を切り離すかCM11を切り離すかによって、システムへの影響が大きく異なることがある。そこで本発明では、後に詳述するように、システムの影響が小さいCMを優先的に切り離し対象として選択することを行う。

【0027】

図5は、コントローラモジュールのブロック構成図である。図1に示すCM10(CM11も同様)は、1または複数のCPU100と、メモリ制御ハブ101と、メモリ110と、ホストプロセッサとのインタフェースを持つ複数のチャンネルアダプタ(CA)40と、ディスク装置とのインタフェースを持つ複数のデバイスインタフェース(DI)50と、DMAチップ等によってデータ転送を行う複数の転送回路60と、モジュール管理制御部(MMC: Module Management Controller)102を備える。

【0028】

メモリ制御ハブ101は、CPU100と、メモリ110と、CA40と、DI50と、転送回路60と、モジュール管理制御部102とを接続する機能を持つ。メモリ110には、ディスクキャッシュとして用いられるキャッシュ領域111、転送データが一時的に格納されるバッファ領域112、各種制御用のテーブルデータが格納されるテーブル領域113、CPU100が実行するための制御用のプログラムが格納されるプログラム領域114が設けられる。

【0029】

モジュール管理制御部102は、電源監視、温度監視などを行い、また保守管理用の通信機能を持つサービスコントローラ(SVC: Service Controller)80との通信によって、ディスク制御装置1の各種の保守機能を実現するものである。

【0030】

転送回路60は、ディスク制御装置1が2台のCMで構成される場合には、相手側のCMの転送回路に直結され、またディスク制御装置1が3台以上のCMで構成される場合には、図1に示す経路設定装置70、71に接続される。

【0031】

図1に示す経路閉塞判定手段20および切り離し制御手段30は、CPU100と、メモリ110のプログラム領域114に格納されたプログラムとによって実現されるが、以下、経路閉塞判定手段20および切り離し制御手段30が行う処理について具体的に説明する。

【0032】

[2CM構成での経路閉塞発生時の処理]

ディスク制御装置1が2台のCMで構成される場合、障害により通信経路がなくなると、どちらかのCMを切り離して1CM構成として運用を行う。この場合、以下の(a)~(d)の条件に従い、切り離すCMを決定する。条件(a)~(d)は、優先順位を(a)>(b)>(c)>(d)とし、優先順位の高いものから条件を探していく。

【0033】

(a) 二重化データを持っていないCM

活性保守でのCMの組み込みや、データの書き戻しが完了していないデータを保持するCMが存在することにより、データが二重化できていない状態で、両通信経路が閉塞した

10

20

30

40

50

場合、データが二重化されていないCMを切り離すとデータロストとなってしまうため、データが二重化されているCMを切り離す。

【0034】

データが二重化されていないCMは、経路閉塞を検出した場合に、相手CMを切り離す。データが二重化されているCMは、基本的には二重化されていないCMから切り離されることを期待し、相手CMからの切り離しを待つが、相手CMがハングアップしたことにより、実際には経路障害ではなくCM障害であるにもかかわらず、経路閉塞に見えている可能性もあるので、SVC間通信を使ったCM生存確認を行い、相手CMの生存が確認できれば、これ以上無駄な動きをしないため自己Panicする。

【0035】

Panicとは、ソフトウェア的にリセットすることであり、通常は、ソフトウェアの矛盾を検出したときなどに処理継続不能となるために行われるものである。自己Panicは、相手CMから切り離されることを期待しているが、それ以上不要な動作をしてシステムに悪影響を出さないように、自らソフトウェア的にリセットし、現在の処理を打ち切ることである。

【0036】

相手CMの生存が確認できない場合、自分が切り離されることが期待できないため、相手CMを切り離す。データが二重化されていないCMを切り離すことによりマシンドウンとする。

【0037】

図6に、ダーティ(Dirty)のないCMの切り離しの例を示す。ここで、ダーティのないCMとは、データが二重化できていない状態のデータを持っていないCMをいう。図6の例では、CM10とCM11とが両経路の閉塞により通信できない状態になっているが、CM10が二重化されていない一重化データを持っているので、CM11を切り離し、CM10により運用を継続する。

【0038】

(b) 転送回路ノースポート(NP)閉塞CM

転送回路は、自分のCMの内部側へのポートと、相手側のCM(または経路設定装置)側へのポートを持つが、前者をノースポート(NP)、後者をサウスポート(SP)と呼ぶ。転送回路NP異常は、図5に示す転送回路60とメモリ制御ハブ101間のCM10内部の異常であり、そのCMを保守しない限り異常を取り除くことはできない。そこで、転送回路NPが閉塞しているCMがある場合には、そのCMを切り離す。

【0039】

ただし、両CMとも転送回路NPに異常がある場合には、次の判定に従う。

- ・転送回路NP異常がないCMは、経路閉塞を検出したならば相手CMを切り離す。
- ・転送回路NP異常があるCMは、SVC間通信を使ったCM生存確認を行い、相手CMの生存が確認できれば自己Panicする。生存が確認できない場合、相手CMを切り離す。

【0040】

図7に、転送回路NP閉塞時のCMの切り離しの例を示す。図7の例では、CM10とCM11とが両経路の閉塞により通信できない状態になっているが、CM11の転送回路61にNP異常が発生している。そこで、CM11を切り離し、CM10により運用を継続する。

【0041】

(c) その他に閉塞しているものがあるCM

転送回路NP閉塞がない場合や、両CMとも同数の転送回路NP閉塞がある場合には、その他の部品が閉塞しているCMを切り離す。以下の優先順で判断する。

- ・優先順1：両転送回路が閉塞になるCM。
- ・優先順2：オンラインのCAが存在しないCM。
- ・優先順3：DIポート閉塞が多いCM。

10

20

30

40

50

## 【 0 0 4 2 】

その他の部品が閉塞していないCMは、経路閉塞を検出したときに相手CMを切り離す。その他の部品が閉塞しているCMは、SVC間通信を使ったCM生存確認を行い、相手CMの生存が確認できれば自己Panicする。生存が確認できない場合、相手CMを切り離す。

## 【 0 0 4 3 】

図8(A)は、転送回路SP閉塞判定によるCMの切り離しの例を示している。図8(A)の例では、CM11の二つの転送回路61がSP(サウスポート)閉塞になっている。したがって、上記優先順1の判断に従ってCM11を切り離す。

## 【 0 0 4 4 】

図8(B)は、CA閉塞判定によるCMの切り離しの例を示している。図8(B)の例では、CM11のすべてのCA41が閉塞になっている。したがって、上記優先順2の判断に従ってCM11を切り離す。

## 【 0 0 4 5 】

図8(C)は、DIポート(Port)閉塞判定によるCMの切り離しの例を示している。図8(C)の例では、CM11のDIポート(ディスク装置への接続部)の閉塞数が、CM10のものより多い。そこで、上記優先順3の判断に従ってCM11を切り離す。

## 【 0 0 4 6 】

## (d) スレーブCM

以上の条件に当てはまらない場合には、マスターCMを優先して残す。マスターCMは、ディスク制御装置1に存在するすべてのCMの状態を管理するCMであり、ディスク制御装置1内の特定の一つのCMが担当する。マスターCM以外のCMは、スレーブCMと呼ばれる。マスターCMが障害により切り離される場合には、スレーブCMの中の一つのCMが新たにマスターCMとして選ばれる。

## 【 0 0 4 7 】

マスターCMは、経路閉塞を検出したならば、スレーブCMを切り離す。スレーブCMは、マスターCMが異常である可能性もあるので、自己PanicするのではなくSVC間通信を使ってマスターCMの生存確認を行う。生存確認で異常がなければ、切り離されるので自己Panicする。この生存確認で異常を検出した場合、マスターCM異常であるのでマスターCMを切り離す。

## 【 0 0 4 8 】

図9は、スレーブCMの切り離しの例を示している。図9の例では、マスターCMであるCM10と、スレーブCMであるCM11との両経路が閉塞しているが、マスターCMを優先して残すため、スレーブCMであるCM11を切り離し、CM10だけで運用を継続する。

## 【 0 0 4 9 】

## [ 3CM以上の構成での経路閉塞発生時の処理 ]

2CM構成の場合とは違い、互い違いに転送回路が閉塞することにより、あるCMとは通信できるが、あるCMとは通信できないというような状態になる場合がある。転送回路閉塞によりCM間通信ができないCMが出てきた場合、通常はすでに転送回路が閉塞しているCMを切り離す。

## 【 0 0 5 0 】

図10は、4CM構成の場合のCMの切り離しの例を示している。図10(A)の例では、CM11の経路設定装置70側の転送回路(#0)61が、すでに切り離し済みになっていたとする。ここで、CM12の経路設定装置71側の転送回路(#1)62に異常が発生すると、CM11とCM12間に通信経路がなくなるため、CM11とCM12は通信できなくなる。この場合、図10(B)に示すように、すでに転送回路が切り離されているCM11を切り離す。

## 【 0 0 5 1 】

以上のように、通常、すでに転送回路が閉塞したCMを切り離すが、そのCMを切り離

10

20

30

40

50

した場合にデータロストとなるようなケースでは、後から転送回路閉塞を検出したCMを切り離す。後から転送回路閉塞を検出したCMを切り離しても運用不能状態やデータロストになる場合には、マシンダウンとしてシステム全体を停止する。

【0052】

図11は、4CM構成の場合のCMの切り離しの他の例を示している。図11(A)の例では、CM11の経路設定装置70側の転送回路(#0)61が、すでに切り離し済みになっているが、CM11はデータが二重化できていないダーティな状態であったとする。ここで、CM12の経路設定装置71側の転送回路(#1)62に異常が発生すると、CM11とCM12間に通信経路がなくなるため、CM11とCM12は通信できなくなる。この場合、図11(B)に示すように、すでに転送回路が切り離されているCM11ではなく、CM12を切り離す。このように、すでに転送回路が閉塞されているCMのデータが二重化されていない場合、切り離すとデータロストとなるため、後から転送回路異常を検出したCMを切り離す。

10

【0053】

同一経路設定装置(FRT)側の転送回路が閉塞しているCMが複数ある状態で、反対側のパスの転送回路が閉塞した場合、または、すでに転送回路が閉塞しているCMがミラーを組んでいない場合、すなわち切り離してもデータロストにはならない場合には、すでに転送回路が閉塞していたCMを切り離す。すでに転送回路が閉塞していたCMを切り離すことによりデータロストとなる場合には、後から転送回路の閉塞を検出したCMを切り離す。

20

【0054】

図12は、4CM構成の場合のCMの切り離しの他の例を示している。図12(A)に示す状態においては、CM10とCM11、CM11とCM12、CM12とCM13、CM13とCM10とが、それぞれミラーを組んでいたとする。ミラーを組んでいるとは、同一データを保持してデータを二重化している状態をいう。また、CM11とCM13の経路設定装置70側の転送回路61、63が、すでに切り離し済みになっていたとする。

【0055】

ここで、CM12の経路設定装置71側の転送回路62に異常が発生すると、CM11とCM12間、CM12とCM13間は通信できなくなる。この場合、CM11とCM13とは、ミラーを組んでいないので、これらを切り離してもデータロストとはならない。そこで、図12(B)に示すように、切り離してもデータロストとはならないCM11とCM13を切り離す。

30

【0056】

同様に、図13(A)に示す状態においては、CM10とCM11、CM11とCM12、CM12とCM13、CM13とCM10とが、それぞれミラーを組んでいたとする。また、CM11とCM12の経路設定装置70側の転送回路61、62が、すでに切り離し済みになっていたとする。この状態で、CM13の経路設定装置71側の転送回路63に異常が発生すると、CM11とCM13間、CM12とCM13間は通信できなくなる。この場合、CM11とCM12とは、ミラーを組んでいるため、これらを切り離すとデータロストとなる。そこで、図13(B)に示すように、切り離してもデータロストとはならないCM13を切り離す。

40

【0057】

次に、図14～図16に従って、本発明の実施の形態に係る経路閉塞判定の処理フローについて説明する。

【0058】

転送回路または経路設定装置の異常が検出され、切り離し事象が検出されると(S10)、閉塞判定が必要かどうかの判定1を行う(S11)。今回切り離し対象となった経路の反対側で、すでに閉塞している転送回路等がない場合、または、切り離し中の転送回路もしくは経路設定装置が存在しない場合には、ステップS13へ進み、通常の切り離しを

50

実施する。また、ディスク制御装置 1 が 1 C M 構成 ( 判定 2 ) の場合にも ( S 1 2 ) , ステップ S 1 3 へ進み、通常の切り離しを実施する。ここで、通常の切り離しとは、異常になった部位だけを閉塞させることをいう。

【 0 0 5 9 】

それ以外の場合にはステップ S 1 4 へ進み、ディスク制御装置 1 が 2 C M 構成かどうかの判定 3 を行う ( S 1 4 ) 。ディスク制御装置 1 が 3 C M 以上で構成される場合、図 1 6 のステップ S 3 0 へ進む。

【 0 0 6 0 】

ディスク制御装置 1 が 2 C M 構成の場合、相手 C M のデータが二重化されているかどうかの判定 4 を行う ( S 1 5 ) 。相手 C M のデータが二重化されている場合、図 1 5 のステップ S 2 0 へ進む。

10

【 0 0 6 1 】

判定 4 の結果、相手 C M のデータが二重化されていない場合、S V C 間通信で相手 C M の生存を確認する ( S 1 6 ) 。すなわち、2 C M 構成の場合、自分が切り離し対象となると判断できたとしても、実際に相手 C M が正常に動作するかどうか判断できない。通信経路は両方閉塞しているため、サービスコントローラ ( S V C ) 8 0 を経由した通信を用いて、相手 C M の生存を確認する。なお、この S V C 8 0 を経由した通信は、基本的には運用中は使用しない。

【 0 0 6 2 】

この S V C 通信の結果、相手 C M から応答があるかどうかの判定 5 を行い ( S 1 7 ) , 相手 C M から応答があれば、相手 C M は生存していると判断し、ステップ S 1 8 へ進む。一方、相手 C M から応答がなければ、相手 C M は生存していないと判断し、ステップ S 1 9 へ進む。

20

【 0 0 6 3 】

相手 C M から応答があった場合には、自 C M が相手 C M から切り離されることになるが、それ以上不要な動作をしないために自己 P a n i c する ( S 1 8 ) 。なお、自己 P a n i c とは、前述のように、相手 C M から切り離されることを期待しているが、それ以上不要な動作をしてシステムに悪影響を出さないように、自らソフトウェア的にリセットすることである。

【 0 0 6 4 】

また、相手 C M から応答がない場合、本来は自 C M を切るべきであるが自 C M 以外に生存している C M が存在しないため、相手 C M を切り離す ( S 1 9 ) 。

30

【 0 0 6 5 】

判定 4 において相手 C M のデータが二重化されていることがわかった場合、図 1 5 のステップ S 2 0 へ進み、自分のデータも二重化されているかどうかの判定 6 を行う。自分のデータが二重化されていない場合、ステップ S 2 9 へ進み、相手 C M を切り離す。

【 0 0 6 6 】

自分のデータも二重化されている場合、相手 C M と自分との転送回路 N P 閉塞数を比較する判定 7 , 判定 8 を行う ( S 2 1 , S 2 2 ) 。相手 C M のほうが転送回路 N P 閉塞が多い場合、ステップ S 2 9 へ進み、相手 C M を切り離す。また、自分のほうが転送回路 N P 閉塞が多い場合、ステップ S 2 6 へ進む。

40

【 0 0 6 7 】

相手および自分の両 C M とも閉塞していない場合、または転送回路 N P 閉塞の数が同数の場合には、ステップ S 2 3 へ進み、相手 C M と自分との閉塞部品数を比較する判定 9 , 判定 1 0 を行う ( S 2 3 , S 2 4 ) 。相手 C M のほうが閉塞部品を多く持っている場合、ステップ S 2 9 へ進み、相手 C M を切り離す。また、自分のほうが閉塞部品数を多く持っている場合には、ステップ S 2 6 へ進む。

【 0 0 6 8 】

相手および自分の両 C M とも閉塞部品がない場合、または閉塞部品数が同数の場合には、ステップ S 2 5 へ進み、自分がマスター C M かどうかの判定 1 1 を行う。自分がマスタ

50

ーCMの場合、ステップS29へ進み、相手CMを切り離す。自分がスレーブCMの場合、ステップS26へ進む。

【0069】

ステップS26では、SVC80を経由した通信を用いて、相手CMの生存を確認する。このSVC通信の結果、相手CMから応答があるかどうかの判定12を行い(S27)、相手CMから応答があれば、相手CMは生存していると判断し、自己Panicする(S28)。一方、相手CMから応答がなければ、相手CMは生存していないと判断し、ステップS29へ進み、相手CMを切り離す。

【0070】

図14の判定3において、ディスク制御装置1が3CM以上の構成であることがわかった場合、図16のステップS30へ進み、今回の閉塞で両経路閉塞になるかどうかの判定13を行う。今回閉塞対象となった転送回路または経路設定装置とは反対側の通信経路(以下、反対経路という)における転送回路切り離しを検出した場合、ステップS31へ進み、両経路閉塞となるCMを切り離す。

【0071】

両経路閉塞にならない場合、反対経路が閉塞しているCMのデータが二重化されているかどうかの判定14を行う(S32)。反対経路が閉塞しているCMのデータが二重化されていなければ、ステップS36へ進む。

【0072】

反対経路が閉塞しているCMのデータが二重化されていれば、その反対経路が閉塞しているCM数が1かどうかの判定15を行う(S33)。反対経路が閉塞しているCM数が1の場合、ステップS35へ進み、反対経路閉塞CMを切り離す。そのCM数が2以上の場合、反対経路が閉塞しているCMがミラーを組んでいるかどうかの判定16を行う(S34)。ミラーを組んでいない場合、すなわち切り離してもデータロスとならない場合、ステップS35へ進み、反対経路閉塞CMを切り離す。ミラーを組んでいて、そのCMを切り離すとデータロスとなる場合、ステップS36へ進む。

【0073】

ステップS36では、今回の検出が転送回路の切り離しかどうかの判定17を行う。今回の切り離しが転送回路ではなく、経路設定装置の切り離しの場合には、どこを切り離してもデータロスとなり、これ以上の運用継続は不可能になるため、ステップS39へ進みサブシステムダウンとする。なお、サブシステムダウンとは、データが二重化されていないCMの切り離しを検出した場合など、システムとしての運用継続が不可能な場合に影響を最小限に抑えるため、ディスク制御装置1のホストインタフェースを閉じて運用を停止させた状態にすることをいう。

【0074】

今回の検出が転送回路の切り離しの場合、今回、転送回路を切り離すCMのデータが二重化されているかどうかの判定18を行う(S37)。二重化されている場合、今回検出のCMを切り離す(S38)。二重化されていない場合、ステップS39へ進み、サブシステムダウンとする。

【0075】

なお、2CM構成のディスク制御装置1においても、CMの切り離しでは、切り離し対象のCMのデータが二重化されていない場合があり、このようなCMの切り離しを検出した場合にはサブシステムダウンとする。また、ステップS31、S35、S38において、切り離しを検出したCMの切り離し対象が自分であった場合、実際に通信できるかわからないが、マスターCMに対して切り離しの検出を通知する(マスターCMが切り離しを検出した場合、次にマスターCMになるCMに対して通知する)。通知の通信結果にかかわらず、その後の処理は何もしないでマスターCMから切り離されるのを待つ。

【0076】

図17は、本発明の実施の形態に係るCM切り離しの処理フローを示す。図1に示す切り離し制御手段30によるCMの切り離しは、次のように行われる。以下、図17に示す

10

20

30

40

50

(A) ~ (H)に従って説明する。この例では、CM10がマスターCMであり、他のCM11 ~ 13がスレーブCMであったとする。

【0077】

(A)スレーブCMであるCM11に異常が発生したとする。

【0078】

(B)CM11は、異常を検出すると、マスターCMであるCM10へ切り離しを通知する。なお、自分以外の他のCMの異常を検出した場合にも、マスターCMへの通知を行う。このとき、マスターCM以外のCMがマスターCMの異常を検出した場合には、次にマスターCMになるCMに異常を通知する。マスターCMが自分の異常を検出した場合には、そのまま切り離し処理を開始する。

【0079】

(C)マスターCMであるCM10は、他のスレーブCMであるCM12, 13に、切り離し開始を通知する。

【0080】

(D)CM12, 13は、CM10から切り離し開始の通知を受けたら、以降の切り離し対象CM11への通信を中止する。なお、マスターCMの切り離しの場合には、新規にマスターになるCMが、他のCMに対してマスターCMの変更情報を配信する。

【0081】

(E)マスターCMからの指示で、全CMへのホストI/Oをサスペンド状態にする。すなわち、ホストプロセッサの入出力を一時的に停止状態にする。

【0082】

(F)マスターCMであるCM10は、切り離し対象CM11のステータスを“デグレード”に変更した構成情報を他のCM12, 13に配信する。

【0083】

(G)マスターCMからの指示で、全CMへのホストI/Oをレジュームする。すなわち、一時的に停止状態にしていたホストプロセッサの入出力を再開する。

【0084】

(H)最後にマスターCMであるCM10は、切り離し対象CM11をハード的にリセットする。これにより、CM11の切り離し処理が完了する。

【0085】

図18は、CM10における経路閉塞判定手段20および切り離し制御手段30が用いる管理用のデータが格納されたテーブルの構成例を示す。これらのテーブルは、メモリ110におけるテーブル領域113に記憶される。これらのテーブルは、各CMが持つが、マスターCM主導で同期される。すなわち、テーブルの内容に変更があると、その変更情報がマスターCMから他のCMへ配信され、同一内容が保持されるようになっている。

【0086】

図18(A)に示すテーブルは、CMのCPU情報を保持するテーブルである。このテーブルにおいて、モジュールIDは、CPUのモジュールを識別する識別情報である。ロール(Role)は、該当CMの役割を示す情報であり、このCMがマスターCMであるかスレーブCMであるかを示す。

【0087】

ダーティフラグ(Dirty Flag)は、ダーティリカバリの状態を示すフラグであり、CMを切り離したり、組み込んだりした場合に、データが一重化状態で二重化されるまでの間にセットされる。このダーティフラグに値がセットされているCMが切り離されるとデータロスとなる。DIポートステータスは、デバイスインタフェース(DI)におけるディスク装置へのポートごとに、ポートが閉塞しているかどうかを示す情報である。

【0088】

CMのCPU情報テーブルは、他にも各種の状態情報や定義されている部品のタイプなどの多くの情報を持つが、本発明に関係しないものについては詳しい説明を省略する。

【0089】

10

20

30

40

50

図18(B)に示すテーブルは、CM単位に必要な情報を格納するCM情報テーブルである。このテーブルにおいて、マスターCPUのモジュールIDは、CMが複数のCPUを搭載する場合に、マスターとなるCPUを示すモジュールIDである。

【0090】

経路ステータスは、CMと経路設定装置(FRT)間の経路についてパス単位で切り離しを行うために用いる状態情報である。転送回路NPステータスは、転送回路ごとにNP閉塞中かどうかを示す状態情報であり、転送回路SPステータスは、転送回路ごとにSP閉塞中かどうかを示す状態情報である。その他にも、例えばモジュール管理制御部102などの各種の状態情報を持つ。

【0091】

図18(C)に示すテーブルは、チャンネルアダプタ(CA)単位に必要な情報を格納するCA情報テーブルである。ステータスとして、CAが正常であるか異常であるかなどを示す情報を保持する。

【0092】

図18(D)に示すテーブルは、経路設定装置単位に必要な情報を格納するFRT情報テーブルである。ステータスとして、経路設定装置が正常であるか異常であるかなどを示す情報を保持する。

【0093】

図18(B)のCM情報テーブルにおける経路ステータスと、図18(D)のFRT情報テーブルにおけるステータスが、経路閉塞判定手段20において、経路閉塞状態の判断に用いられる。切り離し検出から、これらのステータスが変更されるまでには、多少なりともタイムラグが起こる。そのため、マスターCMは、これらの情報テーブル以外に、ローカルに転送回路と経路設定装置について切り離し中か否かの情報を覚えておき、ステータスが変更される前でも、今後ステータスが変化するかどうかを判断できるようにしている。

【0094】

図18(B)のCM情報テーブルにおける転送回路NPステータスが、経路閉塞判定手段20において、転送回路NP閉塞の判断に用いられ、また、これと転送回路SPステータス、および図18(A)のCMのCPU情報テーブルにおけるDIポートステータス、図18(C)のCA情報テーブルにおけるステータスが、経路閉塞判定手段20において、閉塞部品数の判断に用いられる。

【図面の簡単な説明】

【0095】

【図1】本発明の実施の形態に係るディスク制御装置の構成例を示す図である。

【図2】通信経路閉塞の処理を説明する図である。

【図3】3以上のCM構成時のCMの切り離し処理を説明する図である。

【図4】2CM構成時のCMの切り離し処理を説明する図である。

【図5】CMのブロック構成図である。

【図6】ダーティのないCMの切り離しの例を示す図である。

【図7】転送回路NP閉塞時のCMの切り離しの例を示す図である。

【図8】転送回路SP閉塞判定、CA閉塞判定、DIポート閉塞判定によるCMの切り離しの例を示す図である。

【図9】スレーブCMの切り離しの例を示す図である。

【図10】4CM構成の場合のCMの切り離しの例を示す図である。

【図11】4CM構成の場合のCMの切り離しの例を示す図である。

【図12】4CM構成の場合のCMの切り離しの例を示す図である。

【図13】4CM構成の場合のCMの切り離しの例を示す図である。

【図14】本発明の実施の形態に係る経路閉塞判定の処理フローを示す図である。

【図15】本発明の実施の形態に係る経路閉塞判定の処理フローを示す図である。

【図16】本発明の実施の形態に係る経路閉塞判定の処理フローを示す図である。

10

20

30

40

50

【図17】本発明の実施の形態に係るCM切り離しの処理フローを示す図である。

【図18】CMが用いるテーブルの構成例を示す図である。

【図19】従来のディスク制御装置の例を示す図である。

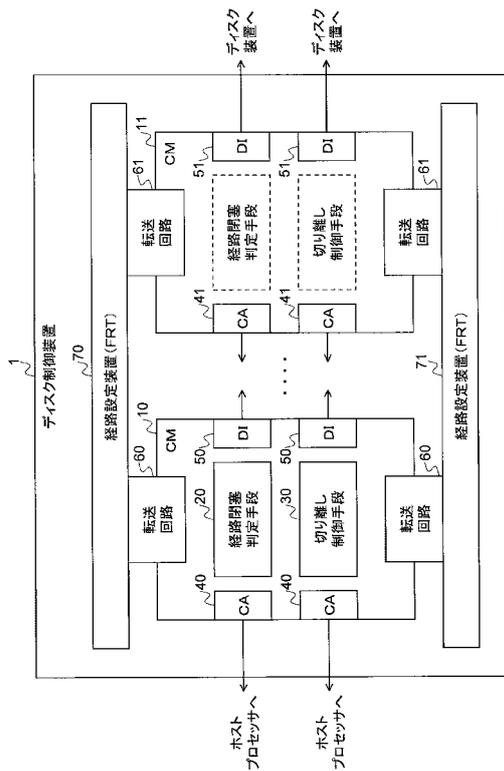
【図20】本発明の課題を説明するための図である。

【符号の説明】

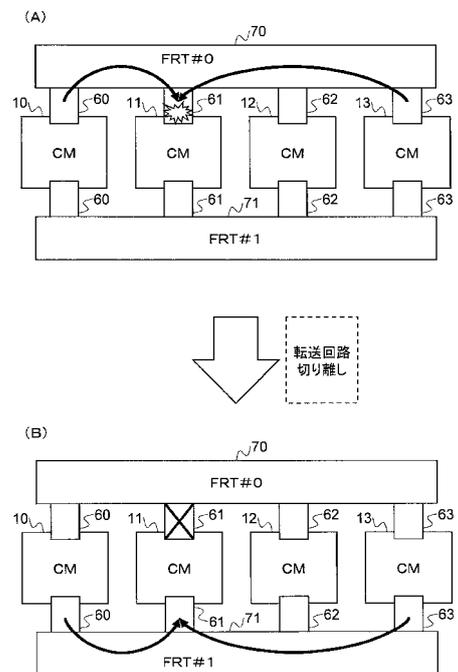
【0096】

- 1 ディスク制御装置
- 10, 11, 12, 13 コントローラモジュール(CM)
- 20 経路閉塞判定手段
- 30 切り離し制御手段
- 40, 41 チャンnelアダプタ(CA)
- 50, 51 デバイスインタフェース(DI)
- 60, 61, 62, 63 転送回路
- 70, 71 経路設定装置(FRT)

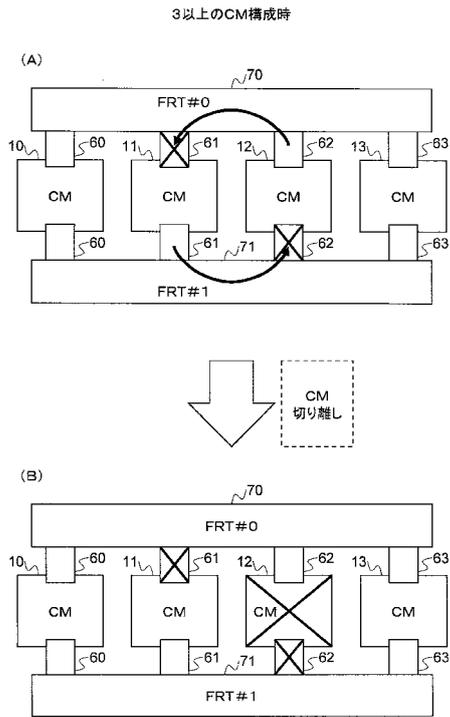
【図1】



【図2】

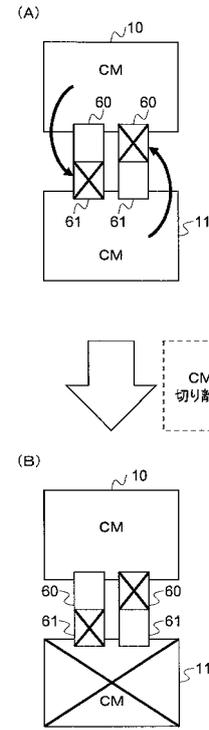


【 図 3 】

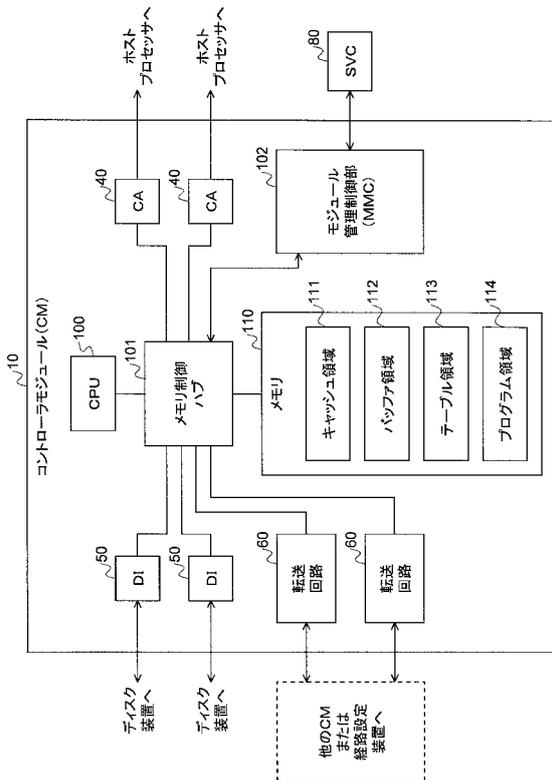


【 図 4 】

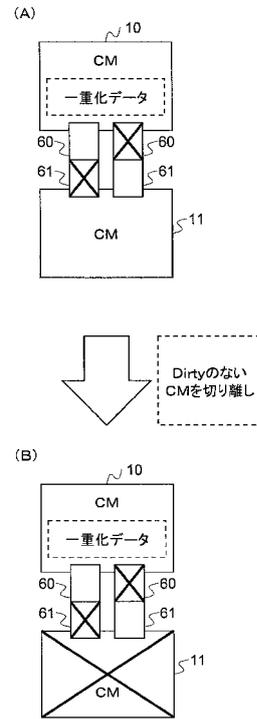
2CM構成時(FRTなしの場合)



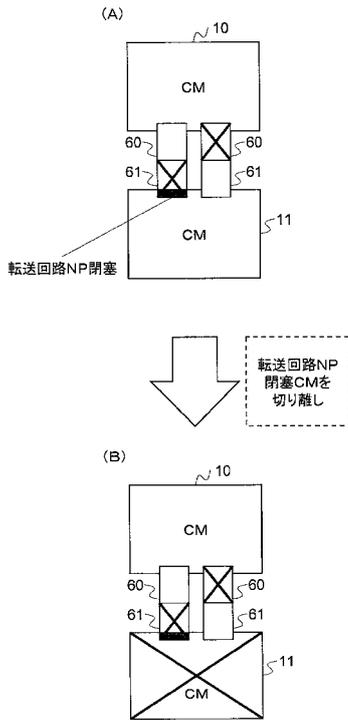
【 図 5 】



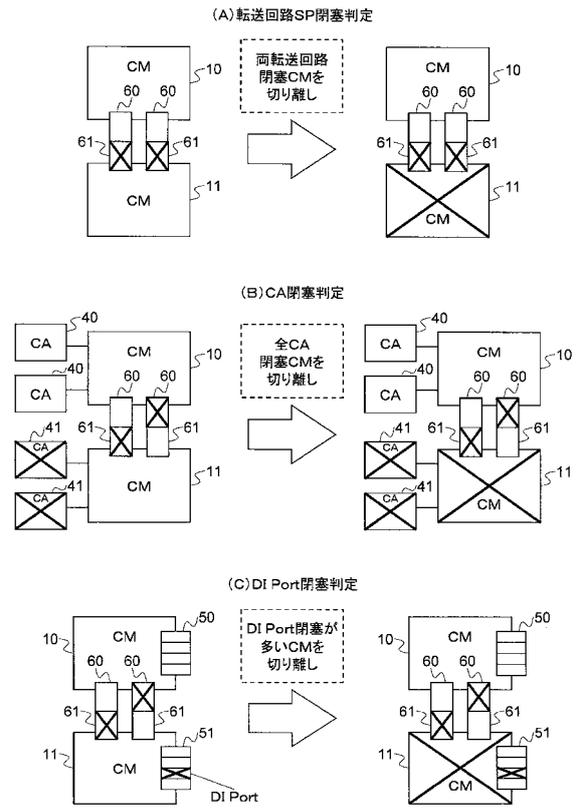
【 図 6 】



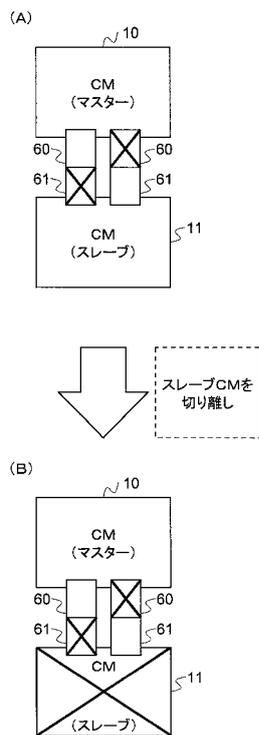
【図7】



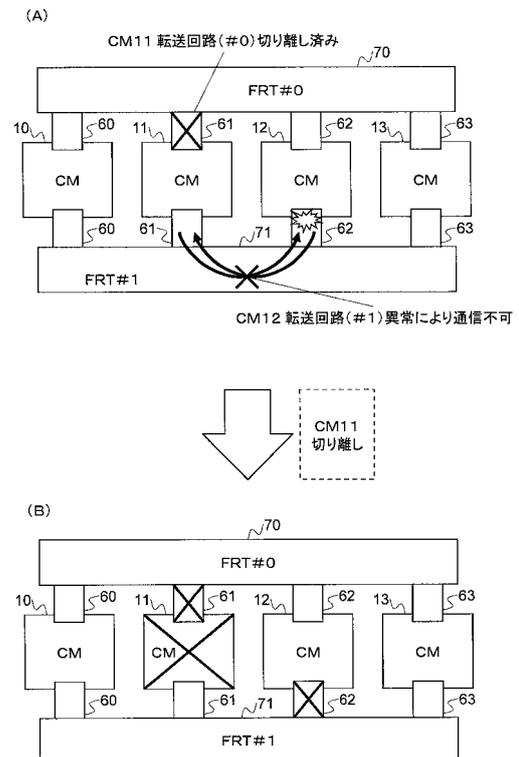
【図8】



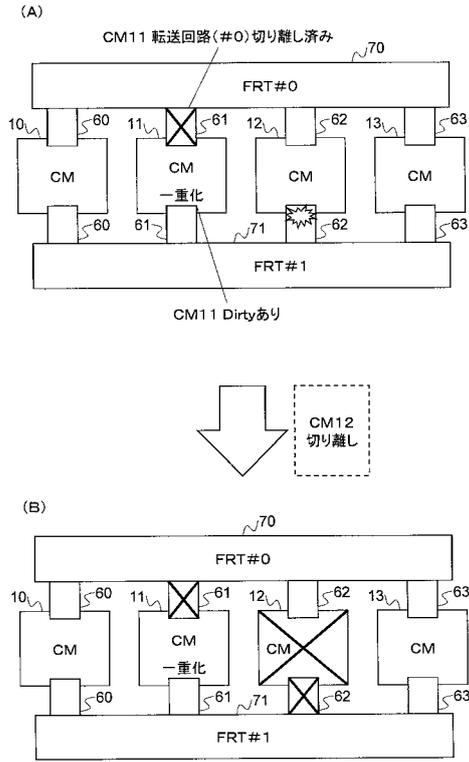
【図9】



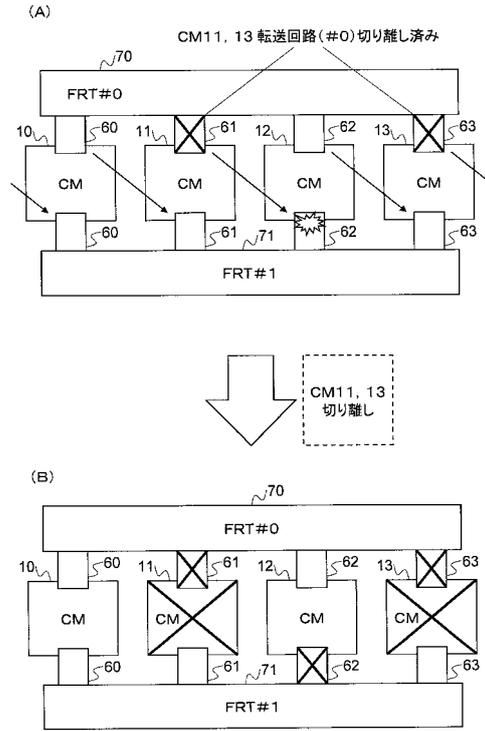
【図10】



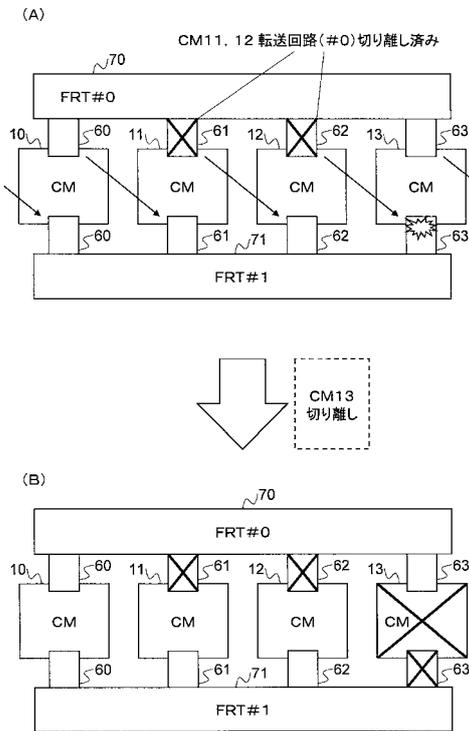
【図11】



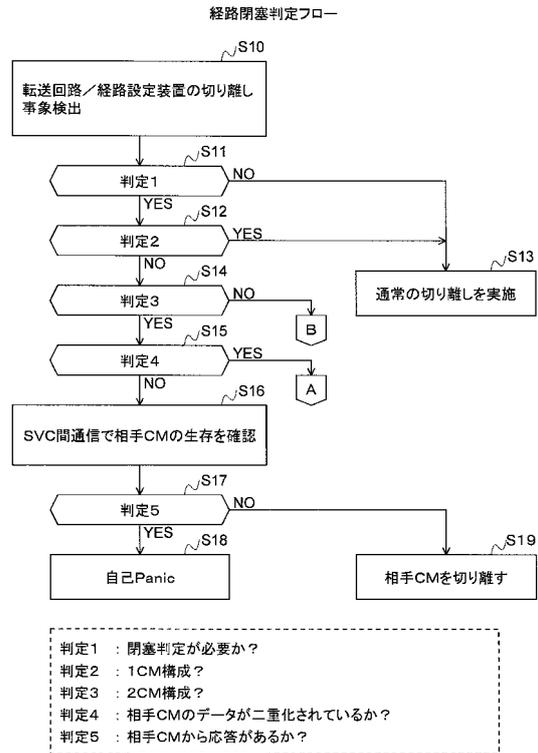
【図12】



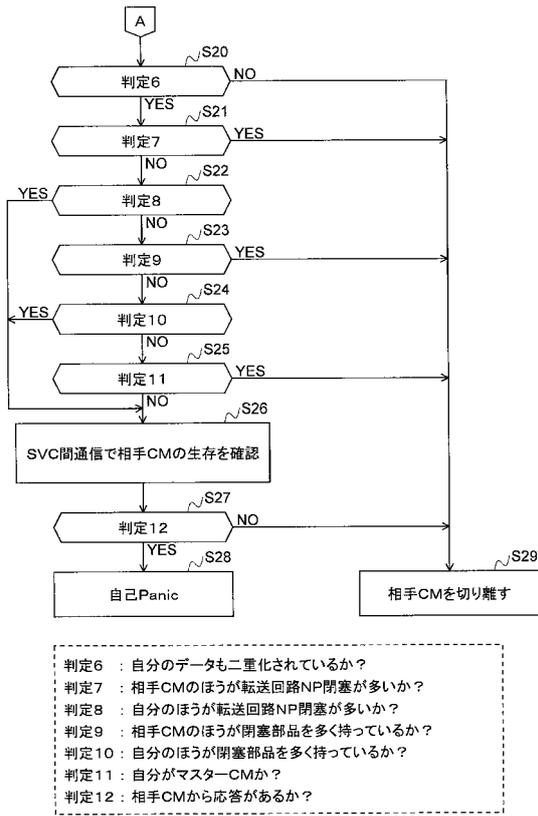
【図13】



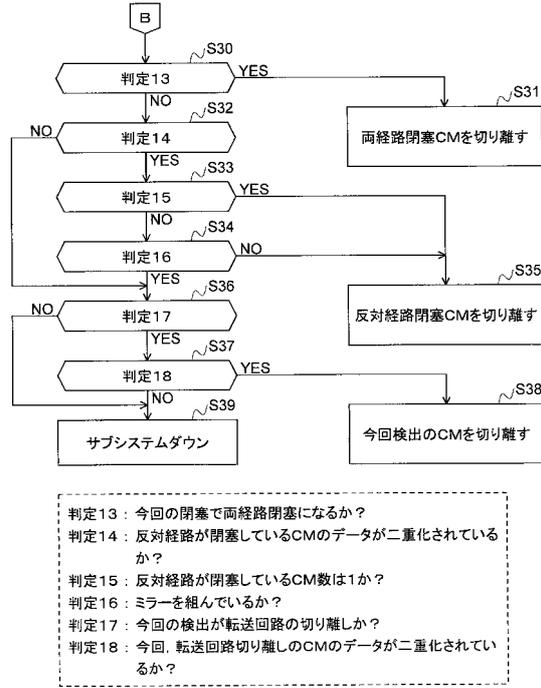
【図14】



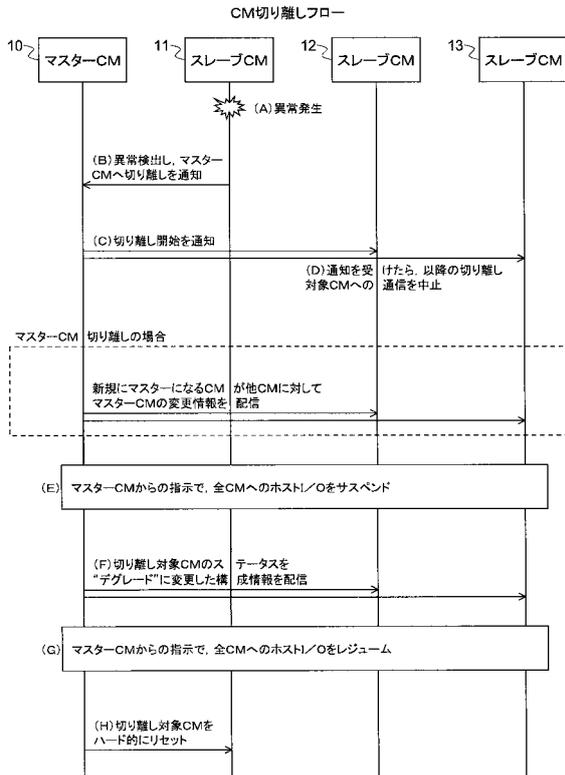
【図15】



【図16】



【図17】



【図18】

(A) CMのCPU情報テーブル

モジュールID	ロール	ダーティフラグ	DIポートステータス	その他
⋮	⋮	⋮	⋮	⋮

(B) CM情報テーブル

マスターCPUのモジュールID	経路ステータス	転送回路NPステータス	転送回路SPステータス	その他
⋮	⋮	⋮	⋮	⋮

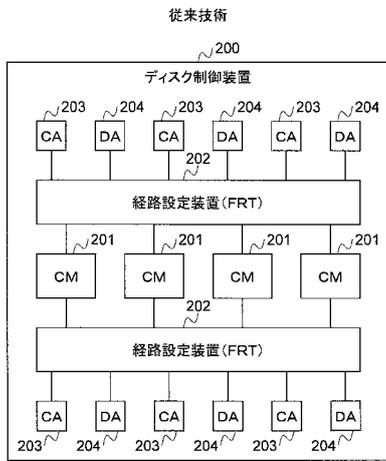
(C) CA情報テーブル

ステータス	その他
⋮	⋮

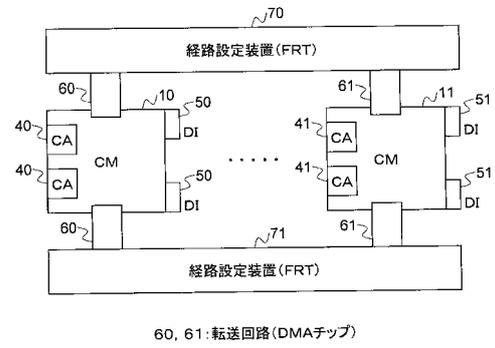
(D) FRT情報テーブル

ステータス	その他
⋮	⋮

【図19】



【図20】



---

フロントページの続き

- (56)参考文献 特開2006-11581(JP,A)  
特開2001-256003(JP,A)  
特開平7-152491(JP,A)  
特開平5-46502(JP,A)  
特開平9-146842(JP,A)  
特開2005-4791(JP,A)  
特開2006-285808(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G06F 3/06