

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4491482号
(P4491482)

(45) 発行日 平成22年6月30日(2010.6.30)

(24) 登録日 平成22年4月9日(2010.4.9)

(51) Int. Cl. F 1
G 0 6 F 11/20 (2006.01) G 0 6 F 11/20 3 1 0 F

請求項の数 11 (全 18 頁)

(21) 出願番号	特願2007-307106 (P2007-307106)	(73) 特許権者	000005108
(22) 出願日	平成19年11月28日 (2007.11.28)		株式会社日立製作所
(65) 公開番号	特開2009-129409 (P2009-129409A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成21年6月11日 (2009.6.11)	(74) 代理人	100075513
審査請求日	平成21年8月6日 (2009.8.6)		弁理士 後藤 政喜
		(74) 代理人	100114236
			弁理士 藤井 正弘
		(74) 代理人	100120260
			弁理士 飯田 雅昭
		(72) 発明者	松本 洋和
			神奈川県横浜市戸塚区戸塚町5030番地
			株式会社日立製作所 ソフトウェア事業部内

最終頁に続く

(54) 【発明の名称】 障害回復方法、計算機、クラスタシステム、管理計算機及び障害回復プログラム

(57) 【特許請求の範囲】

【請求項1】

業務処理を実行する第1の計算機と、前記第1の計算機によって処理されるデータの複製を保持する第2の計算機とを含むクラスタシステムにおいて、前記第1の計算機で発生した障害を回復する方法であって、

前記第1の計算機は、第1のプロセッサと、前記第1のプロセッサに接続される第1の記憶部と、前記第2の計算機に接続される第1のインタフェースとを備え、

前記第2の計算機は、第2のプロセッサと、前記第2のプロセッサに接続される第2の記憶部と、前記第1の計算機に接続される第2のインタフェースとを備え、

前記第1の記憶部は、前記業務処理で使用されるデータを記憶し、

前記クラスタシステムは、当該クラスタシステムの状態を含むシステム情報を保持し、

前記障害回復方法は、

前記第1の記憶部に記憶されたデータを、前記第2の計算機に送信し、

前記第1の計算機から送信されたデータを、前記第2の記憶部に記憶し、

前記第1の計算機に障害が発生した場合には、前記システム情報に基づいて、前記障害が発生した処理を前記第1の計算機で再開するか、又は、前記障害が発生した処理を前記第2の計算機が実行するか、を判定し、

前記障害が発生した処理を前記第1の計算機で再開する場合には、前記第2の記憶部に格納されたデータを前記第2の計算機から前記第1の計算機に送信し、前記第1の計算機に送信されたデータを前記第1の記憶部に記憶し、前記障害が発生した処理を再開し

10

20

、
前記障害が発生した処理を前記第2の計算機が実行する場合には、前記第2の計算機が、前記障害が発生した処理を実行することを特徴とする障害回復方法。

【請求項2】

前記システム情報は、前記第2の計算機の数を含み、

前記障害回復方法は、前記第1の計算機の障害発生時に、前記第2の計算機の数が所定の閾値よりも小さい場合には、前記障害が発生した処理を前記第1の計算機で再開することを特徴とする請求項1に記載の障害回復方法。

【請求項3】

前記システム情報は、前記第1の計算機で実行される処理を構成する各モジュールによって使用されるデータ量を含み、

前記障害回復方法は、

前記第1の計算機の障害発生時に、障害が発生したモジュールを特定し、

前記特定されたモジュールによって使用されるデータ量を前記システム情報から取得し

10

、
前記取得されたデータ量が所定の閾値よりも小さい場合には、前記障害が発生した処理を前記第1の計算機で再開することを特徴とする請求項1に記載の障害回復方法。

【請求項4】

前記システム情報は、前記第1の計算機及び前記第2の計算機の負荷情報を含み、

前記障害回復方法は、前記第1の計算機の障害発生時に、前記第1の計算機の負荷が所定の閾値よりも小さい場合には、前記障害が発生した処理を前記第1の計算機で再開することを特徴とする請求項1に記載の障害回復方法。

20

【請求項5】

前記障害回復方法は、

前記第1の計算機の障害発生時に、前記第1の計算機の負荷が所定の閾値以上の場合には、最も負荷の少ない計算機を選択し、

前記選択された計算機が、前記障害が発生した処理を実行することを指示することを特徴とする請求項4に記載の障害回復方法。

【請求項6】

業務処理を実行する第1の計算機と、前記第1の計算機によって処理されるデータの複製を保持する第2の計算機とを含むクラスタシステムに含まれる第1の計算機であって、プロセッサと、前記プロセッサに接続される記憶部と、前記第2の計算機に接続されるインタフェースとを備え、

30

前記記憶部は、

前記業務処理で使用されるデータを記憶し、

前記クラスタシステムの状態を含むシステム情報を記憶し、

前記プロセッサは、

前記記憶部に記憶されたデータを、前記第2の計算機に送信し、

前記第1の計算機に障害が発生した場合には、前記システム情報に基づいて、前記障害が発生した処理を再開するか、又は、前記障害が発生した処理を前記第2の計算機が実行するか、を判定し、

40

前記障害が発生した処理を再開する場合には、前記第2の計算機から前記第1の計算機によって処理されるデータの複製を取得し、前記取得されたデータを前記記憶部に記憶し、前記障害が発生した処理を再開し、

前記障害が発生した処理を前記第2の計算機が実行する場合には、前記第2の計算機に前記障害が発生した処理を実行するように指示することを特徴とする計算機。

【請求項7】

前記システム情報は、前記第2の計算機の数を含み、

前記プロセッサは、前記第1の計算機の障害発生時に、前記第2の計算機の数が所定の閾値よりも小さい場合には、前記障害が発生した処理を再開することを特徴とする請求

50

項 6 に記載の計算機。

【請求項 8】

前記システム情報は、前記第 1 の計算機で実行される処理を構成する各モジュールによって使用されるデータ量を含み、

前記プロセッサは、

前記第 1 の計算機の障害発生時に、障害が発生したモジュールを特定し、

前記特定されたモジュールによって使用されるデータ量を前記システム情報から取得し、

前記取得されたデータ量が所定の閾値よりも小さい場合には、前記障害が発生した処理を再開することを特徴とする請求項 6 に記載の計算機。

10

【請求項 9】

前記システム情報は、前記第 1 の計算機及び前記第 2 の計算機の負荷情報を含み、

前記プロセッサは、前記第 1 の計算機の障害発生時に、前記第 1 の計算機の負荷が所定の閾値よりも小さい場合には、前記障害が発生した処理を再開することを特徴とする請求項 6 に記載の計算機。

【請求項 10】

前記プロセッサは、

前記第 1 の計算機の障害発生時に、前記第 1 の計算機の負荷が所定の閾値以上の場合には、最も負荷の少ない計算機を選択し、

前記選択された計算機に、前記障害が発生した処理を実行するように指示することを特徴とする請求項 9 に記載の計算機。

20

【請求項 11】

業務処理を実行する第 1 の計算機と、前記第 1 の計算機によって処理されるデータの複製を保持する第 2 の計算機と、前記第 1 の計算機及び前記第 2 の計算機を管理する管理計算機とを含むクラスタシステムであって、

前記第 1 の計算機は、第 1 のプロセッサと、前記第 1 のプロセッサに接続される第 1 の記憶部と、前記第 2 の計算機に接続される第 1 のインタフェースとを備え、

前記第 2 の計算機は、第 2 のプロセッサと、前記第 2 のプロセッサに接続される第 2 の記憶部と、前記第 1 の計算機に接続される第 2 のインタフェースとを備え、

前記管理計算機は、第 3 のプロセッサと、前記第 3 のプロセッサに接続される第 3 の記憶部と、前記第 1 の計算機及び前記第 2 の計算機に接続される第 3 のインタフェースとを備え、

30

前記第 1 の記憶部は、前記業務処理で使用されるデータを記憶し、

前記第 3 の記憶部は、前記クラスタシステムの状態を含むシステム情報を記憶し、

前記第 1 の計算機は、前記第 1 の記憶部に記憶されたデータを前記第 2 の計算機に送信し、

前記第 2 の計算機は、前記第 1 の計算機から送信されたデータを、前記第 2 の記憶部に記憶し、

前記管理計算機は、

前記第 1 の計算機に障害が発生した場合には、前記システム情報に基づいて、前記障害が発生した処理を前記第 1 の計算機で再開するか、又は、前記障害が発生した処理を前記第 2 の計算機が実行するか、を判定し、

40

前記障害が発生した処理を前記第 1 の計算機で再開する場合には、前記第 1 の計算機に前記障害が発生した処理の再開を指示し、

前記第 1 の計算機は、前記第 2 の記憶部に格納されたデータを前記第 2 の計算機から取得し、前記第 2 の計算機から取得したデータを前記第 1 の記憶部に記憶し、前記障害が発生した処理を再開し、

前記管理計算機は、

前記障害が発生した処理を前記第 2 の計算機が実行する場合には、前記システム情報に基づいて、前記障害が発生した処理を継続する第 2 の計算機を選択し、

50

前記選択された第2の計算機に、前記障害が発生した処理を実行することを指示し、
前記選択された第2の計算機は、前記障害が発生した処理を実行することを特徴とする
クラスタシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、現用系の計算機と待機系の計算機を含むクラスタシステムの障害を回復する
技術に関する。

【背景技術】

【0002】

従来、不揮発性の共有ディスクによって処理データを保持し、現用系の計算機と待機系
の計算機とを含むクラスタシステムでは、現用系のプロセスに障害が発生した場合、プロ
セスの再開又は待機系への系切り替えを実行することによって障害を回復させていた。

【0003】

処理性能を向上させるために不揮発性の共有ディスクの代わりに揮発性メモリを用いる
クラスタシステムでは、現用系にプロセス障害が発生すると、データが消失してしまうた
め、回復処理を実行することができない。そこで、現用系のプロセスに障害が発生した場
合の回復手段として、別の計算機に再開のために必要となるデータの複製を転送し、プロ
セスの再開時には、別の計算機に複製したデータを利用して再開を実行する技術が
開示されている（特許文献1参照）。特許文献1に開示された技術では、データを複製す
るために、転送元となる計算機と転送先となる計算機が循環して配置され、すべての計算
機でデータを二重化している。

【特許文献1】特開平9-168015号公報

【発明の開示】

【発明が解決しようとする課題】

【0004】

しかし、特許文献1に開示された技術では、データを二重に保護しているのみであるた
め、プロセス再開時完了までにデータ複製先に障害が発生すると、回復処理を実行す
ることができなくなってしまう。

【0005】

また、現用系のプロセス障害時には必ず同じ系によるプロセス再開を実行し、別系か
らのデータ転送を試みるため、待機系への系切り替えと比較して処理時間が増大してしま
う可能性があった。

【0006】

本発明の目的は、処理データ消失の可能性を最大限に抑えつつ、回復処理の高速化を図
るプロセス回復方法を提供する。

【課題を解決するための手段】

【0007】

本発明の代表的な一形態では、業務処理を実行する第1の計算機と、前記第1の計算機
によって処理されるデータの複製を保持する第2の計算機とを含むクラスタシステムにお
いて、前記第1の計算機で発生した障害を回復する方法であって、前記第1の計算機は、
第1のプロセッサと、前記第1のプロセッサに接続される第1の記憶部と、前記第2の計
算機に接続される第1のインタフェースとを備え、前記第2の計算機は、第2のプロセッ
サと、前記第2のプロセッサに接続される第2の記憶部と、前記第1の計算機に接続され
る第2のインタフェースとを備え、前記第1の記憶部は、前記業務処理で使用されるデー
タを記憶し、前記クラスタシステムは、当該クラスタシステムの状態を含むシステム情報
を保持し、前記障害回復方法は、前記第1の記憶部に記憶されたデータを、前記第2の計
算機に送信し、前記第1の計算機から送信されたデータを、前記第2の記憶部に記憶し、
前記第1の計算機に障害が発生した場合には、前記システム情報に基づいて、前記障害が
発生した処理を前記第1の計算機で再開するか、又は、前記障害が発生した処理を前記

10

20

30

40

50

第2の計算機が実行するか、を判定し、前記障害が発生した処理を前記第1の計算機で再開始する場合には、前記第2の記憶部に格納されたデータを前記第2の計算機から前記第1の計算機に送信し、前記第1の計算機に送信されたデータを前記第1の記憶部に記憶し、前記障害が発生した処理を再開始し、前記障害が発生した処理を前記第2の計算機が実行する場合には、前記障害が発生した処理を前記第2の計算機が実行する。

【発明の効果】

【0008】

本発明の一形態によれば、システムの状態に基づいて、プロセスの再開始と系切り替えのいずれを実行するかを判定することによって、障害復旧処理の高速化及び高信頼化を実現することができる。

10

【発明を実施するための最良の形態】

【0009】

以下、本発明の実施の形態を、図面を参照して説明する。

【0010】

(第1の実施の形態)

図1は、本発明の第1の実施の形態のクラスタシステムの一例を示すシステム構成図である。

【0011】

本発明の第1の実施の形態のクラスタシステムは、現用系の計算機1、及び、複数の待機系の計算機2～nを含む。

20

【0012】

現用系及び待機系の各計算機は、処理データ管理部101、負荷情報管理部201、及びクラスタ情報管理部301を有する。クラスタシステムに含まれる現用系及び待機系の各計算機は、同じ構成である。また、現用系の計算機が系切り替えによって待機系の計算機に処理が引き継がれると、処理を引き継いだ待機系の計算機は、以降、現用系の計算機として稼働する。また、現用系として稼働していた計算機は、可能であれば、待機系の計算機として稼働させてもよい。

【0013】

処理データ管理部101は、処理実行部102及び処理データ103を有する。処理実行部102は、要求された処理を実行する。処理データ103は、処理実行部102によって実行される処理に必要なデータである。また、処理データ103は、処理を高速化させるために揮発性のメモリに記憶されている。なお、処理データ103は、データベースに格納されていてもよい。

30

【0014】

処理データ管理部101は、自系が現用系か待機系かを示すクラスタ状態304を回復判断テーブル303に通知する。処理実行部102は、処理管理部100の各モジュールの回復に必要なデータ量105を計算する。さらに、クラスタ情報管理部301にデータ量105を通知し、回復判断テーブル303に記録する。処理実行部102は、さらに、各モジュールの稼働状態を監視する機能と、障害が発生した場合にクラスタ情報管理部301にプロセス障害を通知する機能を備える。プロセス障害の通知には、障害が発生したモジュールの情報を含む。モジュールについては、図6にて詳細に説明する。

40

【0015】

処理データ管理部101は、さらに、他系とデータを送受信するデータ転送部104を有する。データ転送部104は、処理実行部102によって処理された処理データ103を、他の計算機に転送し、又は、他の計算機から転送された処理データを受信する。なお、本発明の第1の実施の形態では、現用系の計算機のメモリに格納されている処理データ103は、すべての待機系の計算機に転送され、当該計算機のメモリに格納される。

【0016】

データ転送部104による処理データの転送方法は、各計算機に個別にデータを送信するユニキャストであってもよいし、システム内のすべての計算機に対して同時にデータを

50

送信するマルチキャストであってもよい。マルチキャストによって、転送データ量の削減を図ることができる。

【 0 0 1 7 】

また、データ転送部 1 0 4 は、データ転送量に応じて、事前又は転送時にデータを圧縮することなどによって転送量を抑制したり、転送経路を複数使用することによって転送経路を他の処理よりも優先的に利用したりしてもよい。

【 0 0 1 8 】

さらに、本発明の第 1 の実施の形態では、データ転送部 1 0 4 によって他系に処理データを同期転送する。処理データを非同期で転送する場合には、障害発生時に一部の処理データが失われる可能性がある。したがって、データの再生が可能な場合など一部の処理データの欠損が許容されるシステムである場合、又は、さらに上位のシステムなどからデータの再送が可能であれば適用可能である。非同期転送の場合には、他系に処理データを転送後、処理データの格納の完了を待たずに自系の処理を継続できるため、処理性能を向上させることができる。

10

【 0 0 1 9 】

負荷情報管理部 2 0 1 は、負荷情報判断部 2 0 2 及び負荷情報転送部 2 0 3 を有する。負荷情報判断部 2 0 2 は、計算機の負荷情報を判断する。負荷情報転送部 2 0 3 は、負荷情報を他系に転送したり、他系から転送された負荷情報を受信したりする。さらに、負荷情報転送部 2 0 3 は、自系又は他系の負荷情報である負荷量 2 0 4 を回復判断テーブル 3 0 3 に一定間隔で通知する。なお、負荷量を一定間隔で他系に通知するのではなく、障害発生時に他系に負荷量を通知してもよい。この場合には、引き継ぎ先の系を他系が判断する構成としてもよい。

20

【 0 0 2 0 】

クラスタ情報管理部 3 0 1 は、クラスタ情報転送部 3 0 2 及び回復判断テーブル 3 0 3 を有する。クラスタ情報転送部 3 0 2 は、クラスタ情報を他系に転送したり、他系から転送されたクラスタ情報を受信したりする。回復判断テーブル 3 0 3 は、処理実行部 1 0 2 によって処理されたデータ量 1 0 5、クラスタ状態 3 0 4、及び、自系及び他系の負荷量 2 0 4 を格納する。

【 0 0 2 1 】

クラスタ情報管理部 3 0 1 は、処理データ管理部 1 0 1 を監視することによって、自系のプロセス障害を検出する。処理データ管理部 1 0 1 の監視は、データ量 1 0 5 の通信をハートビートとして利用する方法であってもよいし、負荷量 2 0 4 の通信によって負荷量を測定できたか否かを検出する方法であってもよい。また、他の通信によって直接的又は間接的に監視する方法であってもよい。

30

【 0 0 2 2 】

クラスタ情報管理部 3 0 1 は、自系のプロセス障害を検出すると、後述する判断基準に基づいて、プロセス再開又は系切り替えのいずれかを実行するかを判断する。プロセス再開を実行する場合には、処理データ管理部 1 0 1 にプロセスの再開を指示する。処理データ管理部 1 0 1 は、プロセスを再開する指示を受け付けると、データ転送部 1 0 4 を介して、他系に複製されているデータの転送を要求することによって、プロセス再開に必要なデータを取得する。データ取得後、障害が発生した処理データ管理部 1 0 1 の全部又は一部のプロセスを再開し、回復を完了する。

40

【 0 0 2 3 】

一方、クラスタ情報管理部 3 0 1 は、系切り替えを実行する場合には、クラスタ情報転送部 3 0 2 を介して、系切り替え先となる他系に引継ぎを指示する。引継ぎを指示された他系は、データ転送部 1 0 4 によって複製されているデータを取得し、プロセスを実行することによって系切り替えによる回復を完了する。

【 0 0 2 4 】

さらに、クラスタ情報管理部 3 0 1 は、クラスタ情報転送部 3 0 2 によって他系からのクラスタ情報を一定時間受信できなかった場合には、クラスタ情報を受信できなかった他

50

系に障害が発生したものと認識する。他系に障害が発生した場合には、複製されている処理データを利用して、プロセスを起動することによって、系切り替えを実行する機能を有する。ここで、他系の障害検出によって実行される系切り替え処理が、障害が発生した系のプロセス再開又は系切り替えを指示する処理と重複して実行されないように制御する必要がある。例えば、障害が発生した系のプロセス再開又は系切り替えが完了するために必要な時間だけ待機してもよいし、障害が発生した系で回復処理が実行されていないことを確認してから系切り替えを実行するようにしてもよい。さらに、同時に複数の計算機でプロセスが引き継がれないように、共有ディスク及びIPアドレスなどの共有されるリソースが排他制御される仕組みであってもよい。

【0025】

図2は、本発明の第1の実施の形態のハードウェアの構成を示す図である。

【0026】

現用系及び待機系の各計算機は、前述したように同じ構成である。各計算機は、CPU21、ディスプレイ装置22、キーボード23、マウス24、ネットワークインタフェースカード(NIC)25、ハードディスク装置26及びメモリ27を備える。CPU21、ディスプレイ装置22、キーボード23、マウス24、ネットワークインタフェースカード25、ハードディスク装置26及びメモリ27は、バス28によって接続される。

【0027】

現用系及び待機系の各計算機は、NIC25を介してネットワークに接続し、他の計算機と相互に通信する。

【0028】

CPU21は、メモリ27に記憶されたプログラムを実行する。メモリ27は、CPU21によって実行されるプログラム及び当該プログラムの実行に必要なデータを記憶する。メモリ27は、処理管理部100、オペレーティングシステム30、処理データ管理部101、負荷情報管理部201、クラスタ情報管理部301、処理データ103、及び、回復判断テーブル303を記憶する。メモリ27は、前述のように、揮発性のメモリである。

【0029】

処理管理部100は、オペレーティングシステム30上で実行されるプログラムである。処理データ管理部101、負荷情報管理部201、及びクラスタ情報管理部301は、処理管理部100によって呼び出されるプログラムである。処理データ管理部101、負荷情報管理部201及びクラスタ情報管理部301については、図1にて説明した処理を実行する。

【0030】

処理データ103は、業務処理に必要なデータである。処理データ103は、前述したように、データベース管理システムによって管理されていてもよい。この場合、データベース管理システムは、メモリ27に記憶される。回復判断テーブル303は、図1にて説明したように、現用系の計算機で発生した障害を回復させるために必要なクラスタ情報などの情報を格納する。

【0031】

ディスプレイ装置22は、業務処理の実行結果など各種情報を表示する。キーボード23及びマウス24は、利用者からの入力を受け付ける。NIC25は、ネットワークに接続する。ハードディスク装置26は、メモリ27に格納される処理データ、及び、メモリ27にロードされるプログラムなどを格納する。

【0032】

図3は、本発明の第1の実施の形態の回復判断テーブル303の構成を示す図である。

【0033】

回復判断テーブル303は、クラスタ状態判断テーブル331、データ量判断テーブル311及び負荷状態判断テーブル321を含む。

【0034】

10

20

30

40

50

クラスタ状態判断テーブル331は、各計算機のクラスタ状態304、及び、利用者又はシステムによって設定された残台数の閾値情報を含む。本発明の第1の実施の形態では、クラスタ状態には、「現用系」、「待機系」及びプロセスダウンを含む「ダウン」の三状態が定義されているが、さらに詳細なクラスタ状態を定義してもよい。例えば、待機系として起動中である状態を含んでもよい。この場合には、起動後は待機系としての役割を果たすことから待機系として扱ってもよいし、現在は待機系の役割を果たしていないことから待機系として扱わなくてもよい。

【0035】

データ量判断テーブル311は、処理実行部102を構成するモジュールごとのデータ量、利用者又はシステムによって設定されたデータ量の閾値情報、及び、モジュール間の依存関係を表す情報を含む。依存関係は、例えば、図7に示すように、識別子の命名ルールによって表現してもよい。識別子の命名ルールによって依存関係を表す場合には、まず、メインモジュールから直接呼び出される下位モジュールである1段目の各モジュールは、英文字(A-Z)で示される識別子が付与される。さらに、1段目の各モジュールによって呼び出される2段目の各モジュールには、1段目のモジュールの識別子に数字(1-9)を付加した識別子を付与する。なお、モジュール間の依存関係を表す情報は、木構造などの他の手段によって表されてもよい。さらに、モジュール間の依存関係を表す情報は、データ量判断テーブル311とは別のテーブルに保持されてもよい。

10

【0036】

負荷状態判断テーブル321は、各計算機の負荷量204を保持する。負荷状態判断テーブル321は、利用者又はシステムによって設定された負荷量の閾値情報、及び、各計算機の負荷量を含む。負荷量は、例えば、処理対象のデータ量又は処理の所要時間であってもよいし、データ量などの情報を変数とする計算式によって算出される値であってもよい。

20

【0037】

図4は、本発明の第1の実施の形態のクラスタ情報管理部301によるクラスタ状態判断テーブル331に基づいて障害を回復させる処理の手順を示す図である。

【0038】

図4に示す障害回復処理では、クラスタ情報管理部301によって、すべての計算機が障害などによって停止することによる処理データ103の消失を防ぐため、待機系の残り台数が閾値よりも少なくならないように制御する。

30

【0039】

CPU21は、自系(現用系)で障害発生を検知した場合には(ステップ401)、クラスタ状態判断テーブル331を参照し、待機台数の合計を算出する(ステップ402)。さらに、クラスタ状態判断テーブル331から残台数閾値情報を取得する(ステップ403)。

【0040】

CPU21は、待機台数が0か否かを判定する(ステップ404)。待機台数が0の場合には(ステップ404の結果が「Y」)、回復に必要なデータが存在しないため、システム回復が不可能と判断し(ステップ405)、本処理を終了する。なお、ステップ405の処理では、終了以外に、処理データ103を不揮発性ディスクに複写するなどのデータ保護処理を実行してもよい。

40

【0041】

CPU21は、待機台数が0よりも大きい場合には(ステップ404の結果が「N」)、待機台数の合計が残台数閾値情報以下であるか否かを判定する(ステップ406)。待機台数の合計が残台数閾値情報以下の場合には(ステップ406の結果が「Y」)、プロセスの再開を試みる(ステップ407)。さらに、プロセスの再開が成功したか否かを判定する(ステップ408)。

【0042】

一方、CPU21は、待機台数の合計が残台数閾値情報よりも大きい場合(ステップ4

50

06)、又は、プロセスの再開に失敗した場合には(ステップ408の結果が「N」)、待機系に系を切り替える(ステップ409)。待機系への系切り替えが完了、又は、プロセスの再開が成功すると(ステップ408の結果が「Y」)、システムを回復させることができる(ステップ410)。

【0043】

図5A及び図5Bは、本発明の第1の実施の形態のクラスタ状態判断テーブル331に基づいて障害を回復させる処理の一例を示す図である。

【0044】

図5Aに示すケース1では、計算機1は現用系、計算機2~4は待機系となっている。計算機1のクラスタ状態判断テーブル331には、各計算機のクラスタ状態が格納されている。

10

【0045】

ここで、現用系の計算機1に障害が発生した場合には、残台数閾値情報と稼働中の待機系の台数とを比較する。ケース1では、待機系の残り台数は3台であり、閾値(2台)よりも大きいため、待機系への系切り替えを実行する。

【0046】

図5Bに示すケース2では、計算機1は現用系、計算機2、3は待機系、計算機4は障害によるダウン中となっている。現用系の計算機1に障害が発生すると、待機系の残り台数2台であり、残台数閾値情報の値(2台)以下であるため、計算機1は待機系からデータを取得し、プロセスの再開を試みる。

20

【0047】

図6は、本発明の第1の実施の形態のクラスタ情報管理部301によるデータ量判断テーブル311に基づいて障害を回復させる処理の手順を示す図である。

【0048】

本発明の第1の実施の形態では、処理実行部102は、機能ごとにモジュール単位で分割された構成となっている。処理実行時に最初に行われるモジュールをメインモジュールとする。また、各モジュールは、機能ごとに階層構造となっており、上位のモジュールが下位のモジュールを作成し、さらに、下位のモジュールに障害が発生したか否かを監視する。処理実行部102は、障害が発生した場合には、クラスタ情報管理部301に障害の発生したモジュールを通知する。

30

【0049】

次に、クラスタ情報管理部301は、データ量判断テーブル311を参照し、回復が必要となるモジュールを特定する。最下位のモジュールに障害が発生した場合は、当該モジュールを再作成することによって回復させる必要がある。また、下位モジュールを有するモジュールに障害が発生した場合には、すべての下位モジュールもあわせて回復させる必要がある。

【0050】

各モジュールは、処理の実行時に処理データ103を必要とする。プロセスを再開する場合には、障害が発生したモジュールごとに必要なデータを待機系から取得する必要がある。各モジュールが必要とするデータ量が大きい場合には、データ転送における処理時間が増大し、系切り替えと比較して回復処理に必要な時間が大きくなる場合がある。したがって、データ転送量が多い場合には、系切り替えを実行したほうが高速な回復が可能となる。本処理では、データ転送量に基づいてプロセスを再開するか系切り替えを実行するかを判断し、システムを回復させる。

40

【0051】

CPU21は、自系(現用系)で障害発生を検知した場合には(ステップ401)、データ量判断テーブル311を参照し、障害モジュール及び障害モジュールに依存関係を有する下位モジュールを特定し、全モジュールのデータ量の合計を算出する(ステップ421)。さらに、データ量判断テーブル311からデータ量閾値情報を取得する(ステップ422)。

50

【 0 0 5 2 】

C P U 2 1 は、データ量の合計がデータ量閾値情報の値よりも小さいか否かを判定する（ステップ 4 2 3）。データ量の合計がデータ量閾値情報の値よりも小さい場合には（ステップ 4 2 3 の結果が「 Y 」）、待機系の計算機から転送されるデータ量が小さいため、プロセスの再開を試みる（ステップ 4 0 7）。さらに、プロセスの再開が成功したか否かを判定する（ステップ 4 0 8）。

【 0 0 5 3 】

一方、C P U 2 1 は、データ量の合計がデータ量閾値情報の値以上の場合（ステップ 4 0 6 の結果が「 N 」）、又は、プロセスの再開に失敗した場合には（ステップ 4 0 8 の結果が「 N 」）、待機系に系を切り替える（ステップ 4 0 9）。待機系への系切り替えが完了、又は、プロセスの再開が成功すると（ステップ 4 0 8 の結果が「 Y 」）、システムを回復させることができる（ステップ 4 1 0）。

10

【 0 0 5 4 】

図 6 では、障害が発生したモジュールに対してのみプロセスを再開するか否かを判断する例を示したが、依存関係を有するより上位のモジュールを対象として、再帰的にモジュールを再開させてもよい。例えば、系切り替えを実行するステップ 4 0 9 の処理の前に、上位モジュールの再開を再帰的に実行するようにすればよい。また、このように再帰的にモジュールを再開する場合には、データ量の合計を閾値と比較せずに、無条件にプロセスを再開するようにしてもよい。

20

【 0 0 5 5 】

また、検知された障害が自プロセス内のメモリ資源枯渇による障害であった場合には、プロセスの再開によるメモリ状態の初期化によって回復可能な場合がある。したがって、最初に、メインモジュール配下の全モジュールのデータ量を算出し、算出された値に基づいて、プロセスの再開又は系切り替えのいずれを実行するかを判断する処理を追加してもよい。

【 0 0 5 6 】

図 7 は、本発明の第 1 の実施の形態のデータ量判断テーブル 3 1 1 に基づいて障害を回復させる処理の一例を示す図である。

【 0 0 5 7 】

図 7 では、障害（ 1 ）及び障害（ 2 ）が発生した場合について説明する。また、図 7 の説明において、プロセスを再開するか否かを判定するために基準となるデータ量判断テーブル 3 1 1 は、図 3 に示した回復判断テーブル 3 0 3 のデータ量判断テーブル 3 1 1 を利用する。

30

【 0 0 5 8 】

障害（ 1 ）は、モジュール B に障害が発生した場合を示している。この場合、まず、障害が発生したモジュール B に下位モジュールが存在するか否かを、データ量判断テーブル 3 1 1 に含まれるモジュール間の依存関係に基づいて判断する。

【 0 0 5 9 】

データ量判断テーブル 3 1 1 を参照すると、モジュール B には、下位モジュールとしてモジュール B 1 及びモジュール B 2 が存在し、当該モジュールの処理データを待機系から転送する必要があることがわかる。そして、モジュール B、モジュール B 1 及びモジュール B 2 の処理データ量の合計を算出すると、150（= 30 + 70 + 50）となる。さらに、処理データ量の合計とデータ量判断テーブル 3 1 1 に格納された閾値と比較し、プロセスの再開が必要であるか否かを判断する。障害（ 1 ）では、各モジュールのデータ量の合計（150）が閾値（100）よりも大きいため、プロセスを再開せずに系を切り替える。

40

【 0 0 6 0 】

一方、障害（ 2 ）は、モジュール C に障害が発生した場合を示している。同様に、データ量判断テーブル 3 1 1 からモジュール C 及び下位モジュールであるモジュール C 1 の処理データの合計値を算出し、閾値と比較する。障害（ 2 ）では、各モジュールのデータ量

50

の合計(30)が閾値(100)よりも小さいため、プロセスの再開を実行する。

【0061】

図8は、本発明の第1の実施の形態の負荷情報判断部202によって障害を回復させる処理の手順を示す図である。

【0062】

プロセスの再開又は系切り替えによって障害を回復させる場合に、処理を実行する計算機の負荷が高いと、回復処理に要する時間が増大する可能性が高く、さらに、正常に回復処理を実行できない可能性がある。そこで、図8に示す障害回復処理では、できるだけ負荷の低い計算機で処理が継続されるように回復処理を実行する。

【0063】

各計算機の負荷量は、所定の基準に基づいて定められた方法によって決定される値である。例えば、負荷量は、一又は複数の情報に重み付けすることによって算出される。負荷量の基準としては、例えば、CPU使用率、ネットワーク使用率、処理が完了していないデータ量などが挙げられる。また、重み付けの方法としては、前述した負荷量の基準と過去の実行時間に基づいて算出された値を利用して、事前に定義された算出式を用いる方法などがある。

【0064】

負荷情報管理部201は、負荷量を一定間隔で算出し、負荷情報転送部203によって他系に転送する。他系からの負荷量が一定間隔に受信できなかった場合は、当該系の負荷量は高いと判断し、負荷量に最大値を設定する。また、負荷情報管理部201は、自系のクラスタ情報管理部301に算出された自系の負荷量及び受信した他系の負荷量を通知する。クラスタ情報管理部301は、通知された負荷量を回復判断テーブル303の負荷情報判断テーブル321に格納する。

【0065】

CPU21は、自系(現用系)で障害発生を検知した場合には(ステップ401)、負荷状態判断テーブル321を参照し、各計算機の負荷量を取得する(ステップ441)。さらに、負荷状態判断テーブル321から負荷量閾値情報を取得する(ステップ442)。

【0066】

CPU21は、自系の負荷量が負荷量閾値情報の値よりも小さいか否か、又は、自系の負荷量が最も低い場合かを判定する(ステップ443)。自系の負荷量が負荷量閾値情報の値よりも小さい場合、又は、自系の負荷量が最も低い場合には(ステップ443の結果が「Y」)、プロセスの再開を試みる(ステップ407)。さらに、プロセスの再開が成功したか否かを判定する(ステップ408)。

【0067】

一方、CPU21は、自系の負荷量が負荷量閾値情報の値以上の場合、かつ、自系の負荷量が最も低くない場合(ステップ443の結果が「N」)、又は、プロセスの再開に失敗した場合には(ステップ408の結果が「N」)、最も負荷の低い待機系に系切り替えを実行する(ステップ444)。待機系への系切り替えが完了、又は、プロセスの再開が成功すると(ステップ408の結果が「Y」)、システムを回復させることができる(ステップ410)。

【0068】

図9A及び図9Bは、本発明の第1の実施の形態の負荷情報判断部202によって障害を回復させる処理の一例を示す図である。

【0069】

負荷量は、基準となる負荷量を100とした場合の相対的な値とし、値が大きいほど負荷が高いものとする。

【0070】

図9Aに示すケース1では、負荷量の高い計算機1に障害が発生した場合の例を示している。計算機1の負荷量(70)は閾値(40)よりも大きく、他の待機系の計算機のほ

10

20

30

40

50

うが計算機 1 よりも負荷量が小さいため、待機系に切り替える。この場合、系切り替え先は最も負荷量の低い計算機 3 となる。

【 0 0 7 1 】

図 9 B に示すケース 2 では、負荷量の低い計算機 1 に障害が起こった場合の例を示している。計算機 1 の負荷量 (2 0) は閾値 (4 0) よりも小さいため、プロセスの再開を実行する。なお、計算機 1 の負荷量が閾値以上の場合であっても、計算機 1 が最も負荷量の低い計算機であるため、プロセスの再開を実行する。

【 0 0 7 2 】

図 1 0 は、本発明の第 1 の実施の形態の現用系障害時の一連の回復処理の手順を示す図である。

【 0 0 7 3 】

図 1 0 に示した回復処理は、図 4、図 6 及び図 8 に示した手順を組み合わせたものである。各ステップの説明については、前述したとおりである。

【 0 0 7 4 】

本処理は、クラスタ情報管理部 3 0 1 によって、自系 (現用系) の計算機の障害が検知された場合に実行される (ステップ 4 0 1) 。

【 0 0 7 5 】

C P U 2 1 は、まず、クラスタ状態判断テーブル 3 3 1 を参照し、待機系の計算機の台数及び残台数閾値情報の値を比較する (ステップ 4 0 2 ~ 4 0 6) 。処理データ 1 0 3 の消失を防ぐことを最優先とするため、現用系のデータを保持する待機系の計算機が一定台数以上稼働するように制御する。

【 0 0 7 6 】

続いて、C P U 2 1 は、データ量判断テーブル 3 1 1 を参照し、障害を回復するために待機系から取得するデータ量とデータ量閾値情報の値とを比較する (ステップ 4 2 1 ~ 4 2 3) 。そして、転送されるデータ量がデータ量閾値情報の値よりも少ない場合には、プロセスの再開を試みる (ステップ 4 0 7) 。転送されるデータ量が少ないほど、プロセスを再開するために必要な時間が短くなるからである。

【 0 0 7 7 】

最後に、C P U 2 1 は、負荷状態判断テーブル 3 2 1 を参照し、各計算機の負荷量と負荷量閾値情報の値とを比較する (ステップ 4 4 1 ~ 4 4 3) 。自系の計算機の負荷量が負荷量閾値情報の値よりも小さい場合、又は、自系の計算機の負荷がシステム内で最も低い場合には、プロセスの再開を試みる。自系の負荷量が負荷量閾値情報の値以上の場合、かつ、自系の負荷量がシステム内で最も低くない場合、又は、プロセスの再開を失敗した場合には (ステップ 4 4 3 の結果が「N」) 、最も負荷の低い待機系を取得して系切り替えを実行する (ステップ 4 4 4) 。

【 0 0 7 8 】

本発明の第 1 の実施の形態によれば、現用系が回復するために必要なデータをすべての待機系が保持することによって、プロセス回復完了までに連続的に障害が発生した場合であってもデータ消失を防ぐことができる。

【 0 0 7 9 】

また、本発明の第 1 の実施の形態によれば、プロセスの障害回復手段として、システムの状態に基づいて、プロセスの再開と系切り替えのいずれかを実施することによって、障害復旧処理の高速化及び高信頼化を実現することができる。

【 0 0 8 0 】

(第 2 の実施の形態)

本発明の第 1 の実施の形態では、回復判断テーブル 3 0 3 を各計算機が保持していたが、本発明の第 2 の実施の形態では、管理計算機が回復判断テーブル 3 0 3 を保持する。さらに、管理計算機によってプロセスの障害回復方法が決定され、各計算機に指示される。

【 0 0 8 1 】

図 1 1 は、本発明の第 2 の実施の形態のクラスタシステムの一例を示すシステム構成図

10

20

30

40

50

である。

【0082】

本発明の第2の実施の形態のクラスタシステムは、現用系及び待機系の計算機(1~n)以外に管理計算機11を含む。現用系及び待機系の計算機(1~n)と管理計算機11とは、ネットワークを介して接続される。

【0083】

管理計算機11は、クラスタ状態判断テーブル331及び負荷状態判断テーブル321を保持し、現用系に障害が発生した場合に、プロセスを再開するか待機系に系切り替えを実行するかを判断する。また、待機系に切り替える場合には、処理を引き継ぐ計算機を選択する。

10

【0084】

管理計算機11のハードウェア構成は、図2に示した計算機のハードウェア構成と同様であって、CPU、メモリ、NIC及び入出力装置などを備える。なお、管理計算機11は、仮想計算機上で実行されるプログラムによって実現されてもよい。

【0085】

管理計算機11は、回復判断テーブル303、データ量取得部108、クラスタ情報転送部302、負荷情報転送部203及び障害回復部110を含む。

【0086】

回復判断テーブル303は、本発明の第1の実施の形態と同様に、データ量判断テーブル311、クラスタ状態判断テーブル331及び負荷状態判断テーブル321を含む。

20

【0087】

なお、データ量判断テーブル311は、他の情報と比較して更新頻度が多く、管理計算機に格納されたデータ量判断テーブル311を随時更新すると、ネットワークトラフィックが増大し、処理効率が悪化するおそれがあるため、各計算機に格納されている。本発明の第2の実施の形態では、現用系の計算機に格納されたデータ量判断テーブル311の情報を定期的に管理計算機11が取得することによって、ネットワークトラフィックの増大を抑制する。

【0088】

データ量取得部108は、現用系の計算機に格納されたデータ量判断テーブル311から情報を取得し、管理計算機11のデータ量判断テーブル311に情報を格納する。

30

【0089】

クラスタ情報転送部302は、現用系及び待機系の計算機から送信されたクラスタ情報を受信し、管理計算機11のクラスタ状態判断テーブル331に受信したクラスタ情報を格納する。

【0090】

負荷情報転送部203は、現用系及び待機系の計算機から送信された負荷情報を受信し、管理計算機11の負荷状態判断テーブル321に受信した負荷情報を格納する。

【0091】

障害回復部110は、現用系の計算機に障害が発生すると、回復判断テーブル303に格納された情報に基づいて、システムを回復させる。なお、管理計算機11で実行される回復処理は、図10に示した本発明の第1の実施の形態の回復処理と同様である。

40

【0092】

本発明の第2の実施の形態によれば、本発明の第1の実地の形態と同様に、現用系が回復するために必要なデータをすべての待機系が保持することによって、プロセス回復完了までに連続的に障害が発生した場合であってもデータ消失を防ぐことができる。

【0093】

また、本発明の第2の実施の形態によれば、各計算機の情報を一元管理されるため、回復に必要な情報をすべての計算機で共有する必要がない。したがって、回復に必要な情報を転送するために必要なネットワークのトラフィックを軽減することができる。

【0094】

50

さらに、本発明の第2の実施の形態によれば、各計算機がシステム内の他の計算機を監視する必要がなくなるため、各計算機の負荷を軽減することができる。

【図面の簡単な説明】

【0095】

【図1】本発明の第1の実施の形態のクラスタシステムの一例を示すシステム構成図である。

【図2】本発明の第1の実施の形態のハードウェアの構成を示す図である。

【図3】本発明の第1の実施の形態の回復判断テーブルの構成を示す図である。

【図4】本発明の第1の実施の形態のクラスタ状態判断テーブルに基づいて障害を回復させる処理の手順を示す図である。

【図5A】本発明の第1の実施の形態のクラスタ状態判断テーブルに基づいて障害を回復させる処理の一例を示す図である（プロセス再開始）。

【図5B】本発明の第1の実施の形態のクラスタ状態判断テーブルに基づいて障害を回復させる処理の一例を示す図である（系切り替え）。

【図6】本発明の第1の実施の形態のデータ量判断テーブルに基づいて障害を回復させる処理の手順を示す図である。

【図7】本発明の第1の実施の形態のデータ量判断テーブルに基づいて障害を回復させる処理の一例を示す図である。

【図8】本発明の第1の実施の形態の負荷情報判断部によって障害を回復させる処理の手順を示す図である。

【図9A】本発明の第1の実施の形態の負荷情報判断部によって障害を回復させる処理の一例を示す図である（プロセス再開始）。

【図9B】本発明の第1の実施の形態の負荷情報判断部によって障害を回復させる処理の一例を示す図である（系切り替え）。

【図10】本発明の第1の実施の形態の現用系障害時の一連の回復処理の手順を示す図である。

【図11】本発明の第2の実施の形態のクラスタシステムの一例を示すシステム構成図である。

【符号の説明】

【0096】

- 1 ~ n 計算機
- 1 1 管理計算機
- 2 1 CPU
- 2 2 ディスプレイ装置
- 2 3 キーボード
- 2 4 マウス
- 2 5 ネットワークインタフェースカード
- 2 6 ハードディスク装置
- 2 7 メモリ
- 1 0 0 処理管理部
- 1 0 1 処理データ管理部
- 1 0 2 処理実行部
- 1 0 3 処理データ
- 1 0 4 データ転送部
- 1 0 8 データ量取得部
- 1 1 0 障害回復部
- 2 0 1 負荷情報管理部
- 2 0 2 負荷情報判断部
- 2 0 3 負荷情報転送部
- 3 0 1 クラスタ情報管理部

10

20

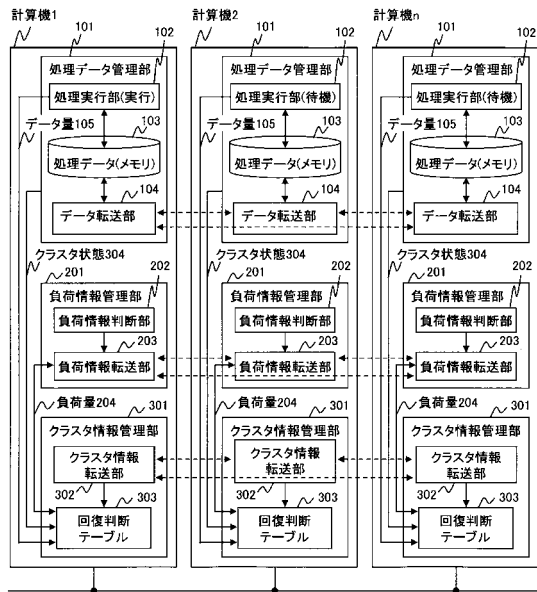
30

40

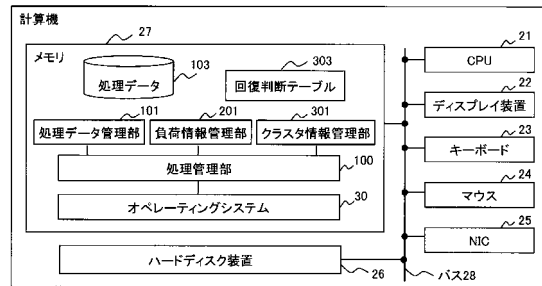
50

- 3 0 2 クラスタ情報転送部
- 3 0 3 回復判断テーブル
- 3 1 1 クラスタ状態判断テーブル
- 3 2 1 負荷状態判断テーブル
- 3 3 1 クラスタ状態判断テーブル

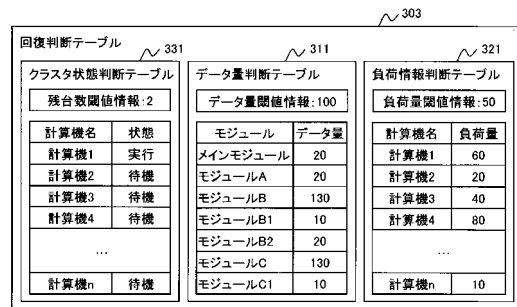
【図1】



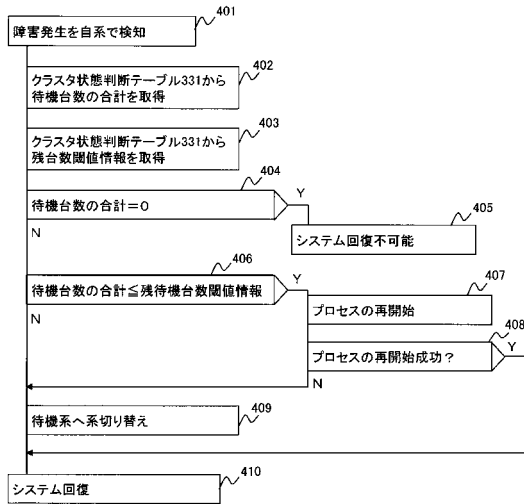
【図2】



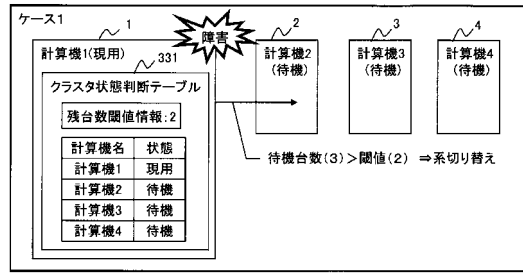
【図3】



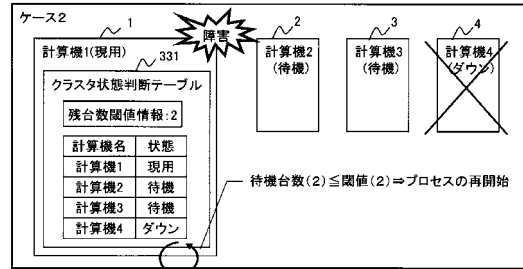
【図4】



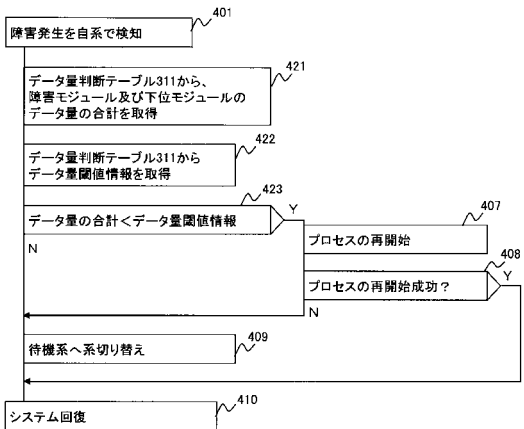
【図5A】



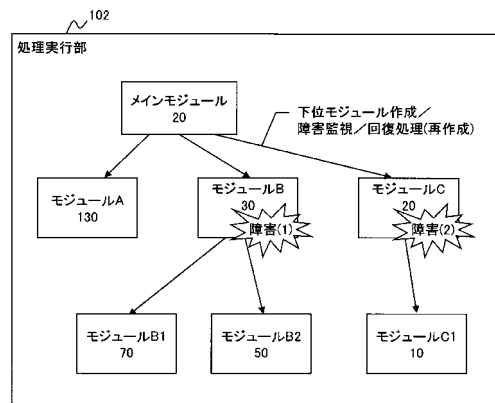
【図5B】



【図6】

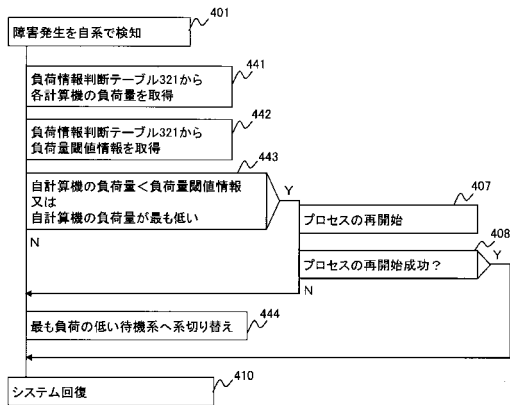


【図7】

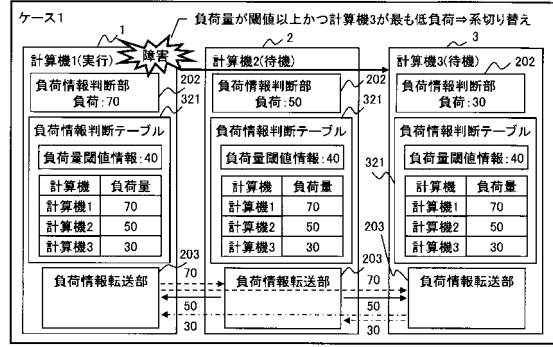


障害(1): 30+70+50 ≥ 100(閾値) ⇒ 系切り替え
 障害(2): 20+10 < 100(閾値) ⇒ プロセスの再開始

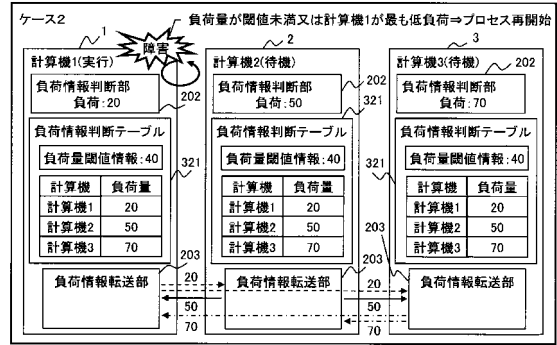
【図8】



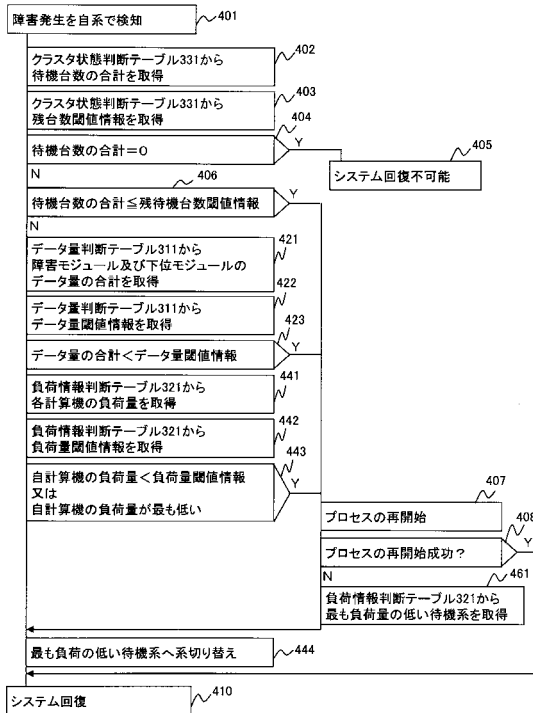
【図9A】



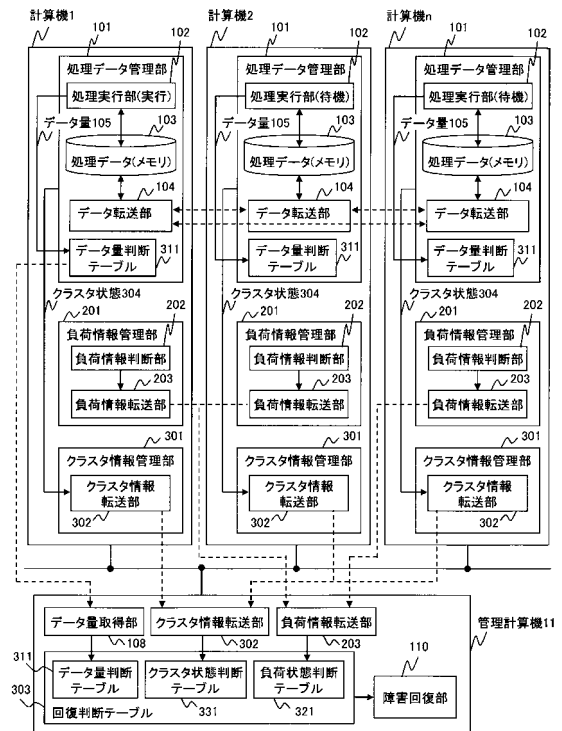
【図9B】



【図10】



【図11】



フロントページの続き

- (72)発明者 馬場 恒彦
東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所 中央研究所内
- (72)発明者 浜田 真二
神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所 ソフトウェア事業部内
- (72)発明者 市村 高志
神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所 ソフトウェア事業部内
- (72)発明者 高橋 則明
神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所 ソフトウェア事業部内

審査官 漆原 孝治

(56)参考文献 特開平05-250197(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F 11/20