

(19) 日本国特許庁(JP)

再公表特許(A1)

(11) 国際公開番号

W02012/063294

発行日 平成26年5月12日 (2014.5.12)

(43) 国際公開日 平成24年5月18日 (2012.5.18)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 11/20 (2006.01)	G06F 11/20 310A	5B034
G06F 9/46 (2006.01)	G06F 9/46 350	

審査請求 有 予備審査請求 未請求 (全 34 頁)

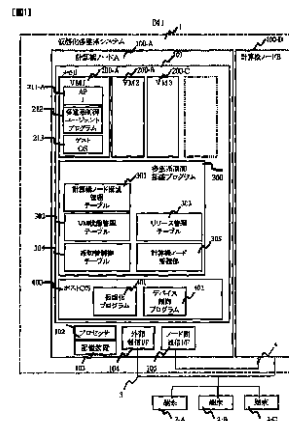
出願番号 特願2012-542705 (P2012-542705)	(71) 出願人 000005108
(21) 国際出願番号 PCT/JP2010/006654	株式会社日立製作所
(22) 国際出願日 平成22年11月12日 (2010.11.12)	東京都千代田区丸の内一丁目6番6号
(81) 指定国 AP (BW, GH, GM, KE, LR, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), EA (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PE, PG, PH, PL, PT, RO, RS, RU, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW	(74) 代理人 100100310 弁理士 井上 学
	(74) 代理人 100098660 弁理士 戸田 裕二
	(74) 代理人 100091720 弁理士 岩崎 重美
	(72) 発明者 金 成昊 神奈川県横浜市戸塚区吉田町292番地 株式会社日立製作所システム開発研究所内
	(72) 発明者 西島 英児 神奈川県横浜市戸塚区吉田町292番地 株式会社日立製作所システム開発研究所内
	Fターム(参考) 5B034 BB01 CC01 DD02

(54) 【発明の名称】 計算機システム

(57) 【要約】

物理計算機上で搭載されている複数の仮想マシンが同一の動作モードとして動作するよう、各仮想マシンの主系/従系を切り替える。

計算機システムは複数の計算機ノードを有し、各計算機ノードは複数の仮想計算機と仮想計算機を制御する制御基盤部とを有する。各仮想計算機は、自計算機ノードとは別の他の計算機ノード上で動作している他の仮想計算機と多重化グループを構成し、いずれか一方が主系として動作し他方が従系として動作する。制御基盤部は、各仮想計算機について、主系/従系のいずれとして動作しているのかを管理し、各仮想計算機各々の状態を監視している。制御基盤部は、自計算機ノード上で主系仮想計算機として動作する仮想計算機の障害を検知した場合に、予め定められた規則に従って、障害が生じた仮想計算機とともに自計算機ノード上で動作する他の仮想計算機も主系仮想計算機から従系仮想計算機へ切り替えるかどうか判定する。



- 1 Virtualization multiplexing system
- 2-A, 2-B, 2-C Terminal
- 100-A Computer node A
- 100-B Computer node B
- 401 Memory
- 102 Processor
- 103 Storage Device
- 104 External communication interface
- 105 Inter-node communication interface
- 212 Multiplexing control agent program
- 218 Host OS
- 300 Multiplexing control base program
- 301 Computer node structure management table
- 302 VM state management table
- 303 Resource management table
- 304 System switch control table
- 305 Computer node monitoring unit
- 400 Host OS
- 401 Virtualization program
- 402 Device control program

【特許請求の範囲】

【請求項 1】

複数の計算機ノードを有する計算機システムであって、

前記複数の計算機ノードは各々、複数の仮想計算機と、当該複数の仮想計算機を制御する制御基盤部とを有しており、

各仮想計算機は、当該仮想計算機が動作している自計算機ノードとは別の他の計算機ノード上で動作している他の仮想計算機と多重化グループを構成し、当該仮想計算機若しくは当該他の仮想計算機のいずれか一方が主系仮想計算機として動作し他方が従系仮想計算機として動作しており、

前記制御基盤部は、当該制御基盤部が備えられている自計算機ノードの複数の仮想計算機各々について、主系仮想計算機若しくは従系仮想計算機のいずれとして動作しているのかを管理し、当該複数の仮想計算機各々の状態を監視しており、

前記制御基盤部は、自計算機ノード上で主系仮想計算機として動作する仮想計算機の障害を検知した場合に、予め定められた規則に従って、障害が生じた仮想計算機とともに自計算機ノード上で動作する他の仮想計算機も主系仮想計算機から従系仮想計算機へ切り替えるかどうか判定し、判定結果に従って切替の制御を行うことを特徴とする計算機システム。

10

【請求項 2】

請求項 1 記載の計算機システムであって、

前記制御基盤部は、前記自計算機ノード上で主系仮想計算機として動作している仮想計算機の数と従系仮想計算機として動作している仮想計算機の数との比率に応じて、障害が生じた仮想計算機とともに自計算機ノード上で動作する他の仮想計算機も主系仮想計算機から従系仮想計算機へ切り替えるかどうか判定することを特徴とする計算機システム。

20

【請求項 3】

請求項 2 記載の計算機システムであって、

前記制御基盤部は、前記自計算機ノード上で主系仮想計算機として動作している仮想計算機の数、当該自計算機ノード上で動作している仮想計算機の総数の半数未満である場合に、障害が生じた仮想計算機とともに当該自計算機ノード上で主系仮想計算機として動作している他の仮想計算機も、主系仮想計算機から従系仮想計算機へ切り替えるよう制御を行うことを特徴とする計算機システム。

30

【請求項 4】

請求項 2 記載の計算機システムであって、

前記制御基盤部は、前記自計算機ノード上で主系仮想計算機として動作している仮想計算機の数、当該自計算機ノード上で動作している仮想計算機の総数の半数以上である場合には、障害が生じた仮想計算機を主系仮想計算機から従系仮想計算機へ切り替えるよう制御し、

前記制御基盤部は、前記障害が回復した場合に、前記仮想計算機を従系仮想計算機から主系仮想計算機へ戻すよう制御することを特徴とする計算機システム。

【請求項 5】

請求項 1 記載の計算機システムであって、

前記制御基盤部は、自計算機ノード上で動作する仮想計算機を主系仮想計算機から従系仮想計算機へ切り替える場合に、当該仮想計算機へ従系仮想計算機への切替要求を通知すると共に、当該仮想計算機と多重化グループを構成している他の仮想計算機を従系仮想計算機から主系仮想計算機へ切り替えるために、当該他の仮想計算機が動作する他の計算機ノードへ通知を送信することを特徴とする計算機システム。

40

【請求項 6】

請求項 1 記載の計算機システムであって、

前記制御基盤部は、前記自計算機ノード上で動作する各仮想計算機が使用するハードウェアリソースを管理しており、主系仮想計算機として動作する複数の仮想計算機に共有されるハードウェアリソースの障害を検知した場合には、当該複数の仮想計算機を共に主系

50

仮想計算機から従系仮想計算機へ切り替えるよう制御することを特徴とする計算機システム。

【請求項 7】

請求項 1 記載の計算機システムであって、

前記制御基盤部は、前記自計算機ノード上で動作する複数の仮想計算機が論理的なグループを構成している場合に、当該複数の仮想計算機の一に障害が生じた場合は、当該論理的なグループを構成している他の仮想計算機も共に主系仮想計算機から従系仮想計算機へ切り替えるよう制御することを特徴とする計算機システム。

【請求項 8】

複数の計算機ノードを有する計算機システムにおいて実行される方法であって、

前記複数の計算機ノードは各々、複数の仮想計算機と、当該複数の仮想計算機を制御する制御基盤部とを有しており、

各仮想計算機は、当該仮想計算機が動作している自計算機ノードとは別の他の計算機ノード上で動作している他の仮想計算機と多重化グループを構成し、当該仮想計算機若しくは当該他の仮想計算機のいずれか一方が主系仮想計算機として動作し他方が従系仮想計算機として動作しており、

前記制御基盤部は、当該制御基盤部が備えられている自計算機ノードの複数の仮想計算機各々について、主系仮想計算機若しくは従系仮想計算機のいずれとして動作しているのかを管理し、当該複数の仮想計算機各々の状態を監視しており、

前記制御基盤部は、自計算機ノード上で主系仮想計算機として動作する仮想計算機の障害を検知した場合に、予め定められた規則に従って、障害が生じた仮想計算機とともに自計算機ノード上で動作する他の仮想計算機も主系仮想計算機から従系仮想計算機へ切り替えるかどうか判定し、判定結果に従って切替の制御を行うことを特徴とする方法。

【請求項 9】

請求項 8 記載の方法であって、

前記制御基盤部は、前記自計算機ノード上で主系仮想計算機として動作している仮想計算機の数と従系仮想計算機として動作している仮想計算機の数との比率に応じて、障害が生じた仮想計算機とともに自計算機ノード上で動作する他の仮想計算機も主系仮想計算機から従系仮想計算機へ切り替えるかどうか判定することを特徴とする方法。

【請求項 10】

請求項 9 記載の方法であって、

前記制御基盤部は、前記自計算機ノード上で主系仮想計算機として動作している仮想計算機の数、当該自計算機ノード上で動作している仮想計算機の総数の半数未満である場合に、障害が生じた仮想計算機とともに当該自計算機ノード上で主系仮想計算機として動作している他の仮想計算機も、主系仮想計算機から従系仮想計算機へ切り替えるよう制御を行うことを特徴とする計算機システム。

【請求項 11】

請求項 9 記載の方法であって、

前記制御基盤部は、前記自計算機ノード上で主系仮想計算機として動作している仮想計算機の数、当該自計算機ノード上で動作している仮想計算機の総数の半数以上である場合には、障害が生じた仮想計算機を主系仮想計算機から従系仮想計算機へ切り替えるよう制御し、

前記制御基盤部は、前記障害が回復した場合に、前記仮想計算機を従系仮想計算機から主系仮想計算機へ戻すよう制御することを特徴とする方法。

【請求項 12】

請求項 8 記載の方法であって、

前記制御基盤部は、自計算機ノード上で動作する仮想計算機を主系仮想計算機から従系仮想計算機へ切り替える場合に、当該仮想計算機へ従系仮想計算機への切替要求を通知すると共に、当該仮想計算機と多重化グループを構成している他の仮想計算機を従系仮想計算機から主系仮想計算機へ切り替えるために、当該他の仮想計算機が動作する他の計算機

10

20

30

40

50

ノードへ通知を送信することを特徴とする方法。

【請求項 1 3】

請求項 8 記載の方法であって、

前記制御基盤部は、前記自計算機ノード上で動作する各仮想計算機が使用するハードウェアリソースを管理しており、主系仮想計算機として動作する複数の仮想計算機に共有されるハードウェアリソースの障害を検知した場合には、当該複数の仮想計算機を共に主系仮想計算機から従系仮想計算機へ切り替えるよう制御することを特徴とする方法。

【請求項 1 4】

請求項 8 記載の計算機システムであって、

前記制御基盤部は、前記自計算機ノード上で動作する複数の仮想計算機が論理的なグループを構成している場合に、当該複数の仮想計算機の一に障害が生じた場合は、当該論理的なグループを構成している他の仮想計算機も共に主系仮想計算機から従系仮想計算機へ切り替えるよう制御することを特徴とする方法。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、サーバ仮想化環境において可用性を上げるために複数の計算機ノードを多重化して構成された計算機システムに関し、特に計算機ノード上のホストマシンおよび各仮想マシンの障害監視と計算機ノードの系切り替えを行う方法に関する。

【背景技術】

20

【0002】

サーバ仮想化技術は、1 台の物理計算機上で複数の仮想マシンを実行し、物理計算機の台数を削減することで、運用コスト、計算機の設置スペース及び消費電力等を削減できるメリットがあり、近年、注目を浴びている。ただし、物理計算機で運用していたシステムと同様な信頼性の確保のため、仮想化環境を提供する物理計算機の多重化とともに、各物理計算機で制御・運用される仮想マシンの多重系構成も要求されている。

【0003】

特許文献 1 には、物理計算機上の各仮想マシン単位で、障害を検知し系を切替える方法が開示されている。特許文献 1 では、図 20 に示すように、物理計算機 2001 と 2031 が、LAN を介して多重化されており、物理計算機 A 2001 はホストマシン 2002 および仮想マシン 2005 と 2015 を備えたサーバ仮想化環境を備えており、物理計算機 B 2031 も同様な構成を備えている。また、障害検知方法および系切替え方法として、物理計算機 A 2001 上の仮想マシン A 2005 上にゲストクラスタプログラム 2007、仮想マシン B 2015 上にゲストクラスタプログラム 2017、および、ホストマシン 2002 上にホストクラスタプログラム 2003 を設けている。物理計算機 B 2031 も同様な構成である。

30

【0004】

ゲストクラスタプログラム 2007 と 2017 または 2037 と 2047 は、アプリケーションプログラム 2006、2016、2036、2046 の稼働状況を監視しつつ、ホストクラスタプログラム 2003 または 2033 に監視結果を通知する。また、ゲストクラスタプログラム 2007、2017 は、ホストクラスタプログラム 2003 から系切替え命令を受けた場合には、各々対応する仮想マシン 2005、2015 に対する主系 / 従系モードの系切替えを実施する。仮想マシンが主系モードの時には、アプリケーションが入力されたデータに対して処理を行い処理結果を出力する。一方、仮想マシンが従系モードの時には、アプリケーションが入力されたデータに対して処理を行うが処理結果の出力は行わない。また、ゲストクラスタプログラム 2037、2047 も同様に、ホストクラスタプログラム 2033 から系切替え命令を受けた場合には、各々対応する仮想マシン 2035、2045 に対する主系 / 従系モードの系切替えを実施する。一方、ホストクラスタプログラム 2003 または 2033 は、同じ物理計算機上で稼働している各仮想マシン（仮想マシン 2005 と 2015 または仮想マシン 2035 と 2045）の稼働状

40

50

況を監視し、例えば、主系モードの仮想マシンA2005が停止していることを検知した場合には、仮想マシンA2035を従系モードから主系モードへ切り替える。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2008-269332号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

特許文献1では、物理計算機上の仮想マシン単位での障害を検知し障害発生時には主系と従系との間の系切り替えを実施する。例えば、図20において、物理計算機A2001上の各仮想マシン2005と2015が全て主系モードで、物理計算機B2031上の各仮想マシン2035と2045が全て従系モードの場合に、物理計算機A2001上の仮想マシンA2005が障害とする。この場合には、物理計算機B上の仮想マシンA2035が従系モードから主系モードに切り替わって稼働する。この場合物理計算機Aには主系モードの仮想マシンと従系モードの仮想マシンが混在した状態となる。物理計算機Bも同様に主系モードおよび従系モードの仮想マシンが混在した状態となる。

10

【0007】

しかしながら、物理計算機に対してハードウェア保守を行いたい場合には、物理計算機上の全ての仮想マシンが従系モード（または停止モード）でなければ、物理計算機を一旦停止させることができない。なぜなら、主系モードの仮想マシンは実際の処理を行った結果を出力しているから、処理中の仮想マシンを停止することができない。つまり、物理計算機上に主系モードの仮想マシンおよび従系モードの仮想マシンが混在している状態では、ハードウェア保守を行うことができない問題がある。

20

【0008】

そこで、仮想化計算機の多重系構成制御において、物理計算機上に搭載されている複数の仮想マシンがなるべく同じの動作モードとなるように、各仮想マシンの系切り替えを実施する技術を提供する。

【課題を解決するための手段】

【0009】

計算機システムは複数の計算機ノードを有し、各計算機ノードは複数の仮想計算機と仮想計算機を制御する制御基盤部とを有する。各仮想計算機は、自計算機ノードとは別の他の計算機ノード上で動作している他の仮想計算機と多重化グループを構成し、いずれか一方が主系仮想計算機として動作し他方が従系仮想計算機として動作している。制御基盤部は、各仮想計算機について、主系/従系仮想計算機のいずれとして動作しているのかを管理し、各仮想計算機各々の状態を監視している。制御基盤部は、自計算機ノード上で主系仮想計算機として動作する仮想計算機の障害を検知した場合に、予め定められた規則に従って、障害が生じた仮想計算機とともに自計算機ノード上で動作する他の仮想計算機も主系仮想計算機から従系仮想計算機へ切り替えるかどうか判定する。

30

【発明の効果】

40

【0010】

仮想化計算機の多重系構成制御において、物理計算機上に搭載されている複数の仮想マシンがなるべく同じの動作モードとなるように、各仮想マシンの系切り替えを実施することができる。

【図面の簡単な説明】

【0011】

【図1】仮想化多重システムの構成例を示す図である。

【図2】計算機ノードの構成例を示す図である。

【図3】計算機ノード構成管理テーブルの構成例を示す図である。

【図4】VM状態管理テーブルの構成例を示す図である。

50

- 【図 5】リソース管理テーブルの構成例を示す図である。
- 【図 6】系切替制御テーブルの構成例を示す図である。
- 【図 7】VM 1 を対象とした系切替制御テーブルの例を示す図である。
- 【図 8】計算機ノード間通信フォーマットの一例を示す図である。
- 【図 9】多重系制御エージェントプログラムと多重系制御基盤プログラム間の通信フォーマットの一例を示す図である。
- 【図 10】仮想化多重系システムで実行される処理の一例を示すフローチャートである。
- 【図 11】リソース監視処理の一例を示す図である。
- 【図 12】リソース監視処理の一例を示すフローチャートである。
- 【図 13】系切替通知処理の一例を示すフローチャートである。
- 【図 14】メッセージ送信処理の一例を示すフローチャートである。
- 【図 15】ノード監視処理の一例を示すフローチャートである。
- 【図 16】生死監視処理の一例を示すフローチャートである。
- 【図 17】全 VM 切替処理の一例を示すフローチャートである。
- 【図 18】主系処理の一例を示すフローチャートである。
- 【図 19】従系処理の一例を示すフローチャートである。
- 【図 20】計算機システムの一例を示す図である。

10

【発明を実施するための形態】

【0012】

以下、実施の形態を図面を用いて説明する。なお、本発明が実施の形態に制限されることは無く、本発明の思想に合致するあらゆる応用例が本発明に該当する。

20

【実施例 1】

【0013】

< 図 1 : 全体構成例 >

図 1 は、仮想化多重システムの全体構成例である。

【0014】

仮想化多重系システム 1 は一つまたは複数の端末 2 と外部ネットワーク 3 を介して接続される。仮想化多重系システム 1 は二台の計算機ノード 100 (図 1 では計算機ノード A 100 - A と計算機ノード B 100 - B) と内部ネットワーク 4 より構成されている。

30

【0015】

計算機ノード 100 はメモリ 101、プロセッサ 102、記憶装置 103、外部通信インタフェース 104 (図中、インタフェースを I / F と略記)、ノード間通信 I / F 105 を備える。計算機ノード A 100 - A と計算機ノード B 100 - B は共に同様の構成を有する。

【0016】

メモリ 101 内にはホストオペレーティングシステム 400 (図中、ホスト OS と略記)、多重系制御基盤プログラム 300 と一つまたは複数の仮想マシン環境 200 (図中、VM と略記) が保持される。ホスト OS 400 内には仮想化プログラム 401 とデバイス制御プログラム 402 が保持される。多重系制御基盤プログラム 300 は計算機ノード構成管理テーブル 301、VM 状態管理テーブル 302、リソース管理テーブル 303、系切替制御テーブル 304 と計算機ノード監視部 305 から構成される。VM 200 はアプリケーション 211 (図中、AP と略記)、多重系制御エージェントプログラム 212、ゲストオペレーティングシステム 213 (図中、ゲスト OS と略記) から構成される。

40

【0017】

尚、図 1 では端末 2 として三台の端末 2 - A、2 - B、2 - C が仮想化多重系システム 1 に接続されており、また計算機ノード A 100 - A 内には VM 200 として VM 1 200 - A、VM 2 200 - B、VM 3 200 - C が備えられている例を示したが、特に台数は限定されるものではない。

【0018】

50

仮想化多重系システム1内の計算機ノードA 100-A、計算機ノードB 100-Bは内部ネットワーク4を介して接続される。計算機ノード100では、外部ネットワーク3は外部通信I/F104に接続し、内部ネットワーク4はノード間通信I/F105に接続する。

【0019】

端末2は外部ネットワーク3を介して仮想化多重系システム1に要求を送信し、仮想化多重系システム1は要求を処理した結果を端末2に返信する。仮想化多重系システム1内の計算機ノードA 100-A、計算機ノードB 100-Bは、一方が「主系ノード」、他方が「従系ノード」として動作する。端末2からの要求は計算機ノードA 100-A、計算機ノードB 100-B両方で受取られ、処理されるが、正常な状態で動作する場合には、処理結果は「主系ノード」である計算機ノード100だけが端末2に返信する。外部通信I/F104などの故障で主系ノードの処理継続に支障が生じた時に、計算機ノード100間の主系ノード/従系ノードの関係は切替る。

10

【0020】

計算機ノード100の記憶装置103はホストOS400、多重系制御基盤プログラム300(プログラムと各種テーブル情報)、VM200を構成するAP211、多重系制御エージェントプログラム212、ゲストOS213を格納する。プロセッサ102は記憶装置103からホストOS400、多重系制御基盤プログラム300、VM200構成物をメモリ101に展開、実行し、また外部通信I/F104やノード間通信I/F105からの割込みを処理する。外部通信I/F104は外部ネットワーク3を通じて端末2とのデータの送受信を行う。ノード間通信I/F105は内部ネットワーク4を通じて他の計算機ノード100との間でデータの送受信を行う。

20

【0021】

仮想化プログラム401はVM200を作成、管理するためのプログラムであり、デバイス制御プログラム402は外部通信I/F104、ノード間通信I/F105や記憶装置103にアクセスするためのプログラムである。ホストOS400は仮想化プログラム401やデバイス制御プログラム402を用いて、多重系制御基盤プログラム300やVM200の実行を制御し、外部通信I/F104、ノード間通信I/F105や記憶装置103へのアクセス制御を行う。

【0022】

VM200は固有のプログラムやオペレーティングシステムが動作することが可能な、仮想化プログラム401により作成された仮想的な実行環境である。複数のVM200が存在する場合に、異なるVM200内で動作するプログラムやオペレーティングシステムは同じでも良いし、異なっても良い。また、異なるVM200内で動作するプログラムやオペレーティングシステムが互いに直接影響を与え合うことはない。VM200は仮想化多重系システム1内で一意の識別子を持ち、例えば、図1中のVM200は「VM1」、「VM2」、「VM3」の識別子を保有する。また異なる計算機ノード100内で動作する任意のVM200の対(ペア)は、「主系VM」/「従系VM」の関係を持つ。通常、主系ノードの計算機ノード100内のVM200は「主系VM」となり、従系ノードの計算機ノード100内のVM200は「従系VM」となる。「主系VM」/「従系VM」の関係を持つVM200の対は、二つの単体VM200間で構成される場合もあるし、一つの計算機ノード100内の複数のVMグループが協同してアプリケーションを実行してサービスを提供する場合には複数のVM200から構成されるVMグループ間で構成される場合もある。主系VMの動作に異常が生じた場合には、「主系VM」/「従系VM」の関係を入れ替える、系切替処理が行われる。

30

40

【0023】

AP211は端末2からの要求を処理する。AP211は自身が動作するVM200が主系VM、従系VMのいずれであっても端末2からの要求を読み取り、演算処理を行うが、端末2への応答動作は異なる。主系VMで動作するAP211は演算結果を端末2へ返答するが、従系VMで動作するAP211は演算結果を返答しない。AP211は「系切替

50

」を行うことで端末 2 への返答動作を切替える。

【 0 0 2 4 】

多重系制御エージェントプログラム 2 1 2 は多重系制御基盤プログラム 3 0 0 との間で情報通信を行うプログラムであり、A P 2 1 1 の動作状態を監視して多重系制御基盤プログラム 3 0 0 に通知したり、多重系制御基盤プログラム 3 0 0 からの指示に従って A P 2 1 1 の系切替を行ったりする。A P 2 1 1 の動作状態を監視する方法の詳細は図 1 1、図 1 2 と共に後述する。

【 0 0 2 5 】

ゲスト OS 2 1 3 は V M 2 0 0 の環境内で A P 2 1 1 や多重系制御エージェントプログラム 2 1 2 の動作を制御する。

10

【 0 0 2 6 】

多重系制御基盤プログラム 3 0 0 は、他の計算機ノード 1 0 0 内で動作する多重系制御基盤プログラム 3 0 0 との間で情報を交換することで仮想化多重系システム 1 内の V M 2 0 0 の状態を判断し、V M 2 0 0 の「主系 V M」/「従系 V M」の系切替を制御する。主系ノードである計算機ノード 1 0 0 内で動作する多重系制御基盤プログラム 3 0 0 は V M 2 0 0 の状態判断と系切替の実行を指示し、従系ノードの計算機ノード 1 0 0 内で動作する多重系制御基盤プログラム 3 0 0 は主系ノードの生死状態を監視する。計算機ノード 1 0 0 の監視は、主系ノード、従系ノードの双方から行う方式も有り得るが、本実施形態では従系ノードからのみ実行する。しかし、本発明の思想がこれに限定されるものでは無い。多重系制御基盤プログラム 3 0 0 の V M 2 0 0 の状態判断の詳細は図 1 3 と共に後述する。系切替や生死状態監視処理についての詳細は、計算機ノード監視部 3 0 5 の処理説明として、図 1 4、図 1 5、図 1 6、図 1 7、図 1 8、図 1 9 と共に後述する。

20

【 0 0 2 7 】

計算機ノード構成管理テーブル 3 0 1 は多重系制御基盤プログラム 3 0 0 が計算機ノード 1 0 0 の主系ノード/従系ノードなどの状態を管理するためのテーブルである。詳細は図 3 と共に後述する。

【 0 0 2 8 】

V M 状態管理テーブル 3 0 2 は多重系制御基盤プログラム 3 0 0 が計算機ノード 1 0 0 内の V M 2 0 0 の主系 V M / 従系 V M などの状態を管理するためのテーブルである。詳細は図 4 と共に後述する。

30

【 0 0 2 9 】

リソース管理テーブル 3 0 3 は多重系制御基盤プログラム 3 0 0 が、V M 2 0 0 が使用しているメモリ 1 0 1 や実行している A P 2 1 1、他の V M 2 0 0 との関係などの「リソース」を管理するためのテーブルである。詳細は図 5 と共に後述する。

【 0 0 3 0 】

系切替制御テーブル 3 0 4 は多重系制御基盤プログラム 3 0 0 が、V M 2 0 0 が使用するリソースに異常が発生した時などにおける、障害内容とこれに対する対応動作を記載したテーブルである。詳細は図 6、図 7 と共に後述する。

【 0 0 3 1 】

計算機ノード監視部 3 0 5 は多重系制御基盤プログラム 3 0 0 が他の計算機ノード 1 0 0 内で動作する多重系制御基盤プログラム 3 0 0 との間で通信を行い、生死状態監視処理を行う処理である。計算機ノード監視部 3 0 5 は多重系制御基盤プログラム 3 0 0 を構成するプログラムではあるが、多重系制御基盤プログラム 3 0 0 の処理とは独立した契機で活動する。詳細は図 1 4、図 1 5、図 1 6、図 1 7、図 1 8、図 1 9 と共に後述する。

40

< 図 2 : 仮想化環境構成例 >

図 2 は計算機ノードの構成の例を示す図である。アプリケーションやオペレーティングシステムの動作環境である V M 2 0 0 は仮想的な計算機の実行環境であり、実際には、計算機ノード 1 0 0 内のメモリ 1 0 1 やプロセッサ 1 0 2 などの現存するリソースを割当てることで構築される。

【 0 0 3 2 】

50

図 2 は計算機ノード 100 のメモリ 101、プロセッサ 102、外部通信 I/F 104 等のリソースを複数の VM 200 やホスト OS 400 への割当てた例を示した図である。

【0033】

図 2 中において、計算機ノード A 100 - A 内には複数の VM 200 (VM 1 200 - A、VM 2 200 - B、VM 3 200 - C) が作成されている。メモリ 101 は複数のメモリ領域に分割することが可能で、それぞれのメモリ領域が個別に各 VM 200 やホスト OS 400 に割り当てられたり、共有されたりする。プロセッサ 102 は「マルチコア」のプロセッサで、コア単位で VM 200 やホスト OS 400 に割り当てられたり、共有されたりする。外部通信 I/F 104 は計算機ノード 100 内には複数存在し、それぞれの外部通信 I/F 104 が VM 200 やホスト OS 400 に割り当てられたり、共有されたりする。

10

【0034】

仮想環境 VM 1 200 - A は、プロセッサ 102 の VM 1 占有コア # 1 102 - A、VM 1 占有メモリ 101 - A と VM 1 と VM 2 間共有メモリ 101 - B のメモリ領域、VM 1 と VM 2 共有外部通信 I/F 104 - A のリソースを使用して構築される。仮想環境 VM 1 200 - A は、プロセッサ 102 の VM 2 占有コア # 2 102 - B、VM 2 占有メモリ 101 - C と VM 1 と VM 2 共有メモリ 101 - B のメモリ領域、VM 1 と VM 2 共有外部通信 I/F 104 - A のリソースを使用して構築される。

VM 1 200 - A、VM 2 200 - B に割り当てられていないリソースであるプロセッサ 102 のホスト OS 占有コア # 0 102 - C は、ホスト OS 400 のみが使用するプロセッサ・コアで、VM 200 から使用されることは無い。また外部通信 I/F 104 - B はホスト OS 400 と、VM 1 200 - A と VM 2 200 - B 以外の VM 200 が使用する。

20

【0035】

計算機ノード B 100 - B 内にも複数の VM 201 (VM 4 201 - A、VM 5 201 - B、VM 6 201 - C) が作成されており、各 VM 201 やホスト OS 400 には計算機ノード B 100 - B 内のリソースが割り当てられる。

< 図 3 : 計算機ノード構成管理テーブル >

図 3 は計算機ノード構成管理テーブル 301 の構成例を示す図である。計算機ノード構成管理テーブル 301 は多重系制御基盤プログラム 300 内に保持され、仮想化多重系システム 1 内の計算機ノード 100 の状態情報を格納する。仮想化多重系システム 1 内の各計算機ノード A 100 - A、計算機ノード B 100 - B は、同じ内容の計算機ノード構成管理テーブル 301 を保持する。

30

【0036】

計算機ノード識別子 311 には仮想化多重系システム 1 内の計算機ノード 100 を識別する一意の情報を格納する。

【0037】

主/従系状態フラグ 312 には各計算機ノード 100 が「主系ノード」、「従系ノード」のいずれの状態にあるかを識別する情報を格納する。対象の計算機ノード 100 が主系ノードの場合は「M」を格納し、従系ノードの場合は「S」を格納する。

40

【0038】

マルチキャストアドレス 313 には計算機ノード 100 にアクセスするための、内部ネットワーク 4 のマルチキャストアドレスを格納する。

【0039】

生死状態フラグ 314 には計算機ノード 100 の正常/異常状態を識別する情報を格納する。計算機ノード 100 が「正常状態」の場合には「0」値を格納し、「異常状態」の場合には「1」を格納する。本発明の実施の形態では正常/異常状態を「0/1」値で表すが、正常/異常状態の識別が可能な情報であれば格納情報は「0/1」値に限らない。

【0040】

監視周期 315 は計算機ノード 100 の状態を監視する周期時間を格納する。従系ノード

50

ドの多重系制御基盤プログラム300の計算機ノード監視部305は、監視周期315の時間周期で、主系ノードの計算機ノード100の状態を取得する。主系ノードは従系ノードの計算機ノード100の状態取得要求に応じて、状態情報を返信する。

【0041】

生死確定カウンタ316は計算機ノード100の状態取得要求に対して応答が無かった場合に、相手計算機ノード100が異常状態にあることを判断するためのしきい値を格納する。生死確定カウンタ316の回数分、状態取得要求に対する相手計算機ノード100からの返信が連続して無かった場合は、相手計算機ノード100が異常状態にあると判断する。

<図4：VM状態管理テーブル>

図4は本発明の実施の形態のVM状態管理テーブル302の構成を示す図である。VM状態管理テーブル302は多重系制御基盤プログラム300内に保持され、計算機ノード100内に作成されているVM200の状態情報を保持する。

【0042】

VM状態管理テーブル302に格納される状態情報は計算機ノード100毎に異なる。計算機ノードA 100-AのVM状態管理テーブル302はVM1 200-A、VM2 200-B、VM3 200-Cの状態情報を保持する。計算機ノードB 100-BのVM状態管理テーブル302はVM4 201-A、VM5 201-B、VM6 201-Cの状態情報を保持する。

管理VM総数321には一台の計算機ノード100内に作成されているVM200の総数を格納する。

【0043】

障害VM総数322には一台の計算機ノード100内に作成されているVM200内、異常状態にあるVM200の個数を格納する。障害VM総数322は「0」以上、管理VM総数321の値以下の値となる。

【0044】

VM識別子323には仮想化多重系システム1内でVM200を一意に識別する情報を格納する。行328はVM1 200-Aの状態情報を格納した行であることを示す。

【0045】

ホスト計算機ノード識別子324はVM200が作成されている計算機ノード100の識別情報を格納する。

【0046】

主/従系状態フラグ325にはVM200の「主系VM」、「従系VM」を識別する情報を格納する。対象のVM200が主系VMの場合は「M」を格納し、従系VMの場合は「S」を格納する。

【0047】

生死状態フラグ326にはVM200の正常/異常状態を識別する情報を格納する。VM200が「正常状態」の場合には「0」値を格納し、「異常状態」の場合には「1」を格納する。本発明の実施の形態では正常/異常状態を「0/1」値で表すが、正常/異常状態の識別が可能な情報であれば格納情報は「0/1」値に限らない。

【0048】

主/従系構成VM情報327には、VM識別子323のVM200と主系VM/従系VMの関係に対応するVM200の情報を格納する。ホスト計算機ノード識別子327-Aには対応するVM200が作成されている計算機ノード100の識別情報を格納し、VM識別子327-Bには対応するVM200を一意に特定する識別情報を格納する。

<図5：リソース管理テーブル>

図5はリソース管理テーブル303の構成例を示す図である。リソース管理テーブル303は多重系制御基盤プログラム300内に保持され、VM200が使用する個別リソースの状態情報を格納する。リソース管理テーブル303に格納される状態情報は計算機ノード100毎に異なる。計算機ノードA 100-Aのリソース管理テーブル303はV

10

20

30

40

50

M1 200 - A、VM2 200 - B、VM3 200 - Cが使用するリソースの状態情報を保持する。計算機ノードB 100 - Bのリソース管理テーブル303はVM4 201 - A、VM5 201 - B、VM6 201 - Cが使用するリソースの状態情報を保持する。

【0049】

VM識別子331にはVM200を一意に識別する情報を格納し、リソース識別子332にはVM200が使用するリソースを一意に特定する情報を格納する。行336はVM1 200 - AがAP1 211 - Aをリソースとして使用していることを意味する。

【0050】

監視周期333には使用リソースの正常/異常状態を監視するための周期時間を格納する。リソース状態の取得をイベントなどの契機で行い、周期的に監視しない場合は無効値である「0」を格納する。

10

【0051】

リソース生死状態フラグ334には使用リソースの正常/異常状態を識別する情報を格納する。リソースの状態が「正常状態」の場合には「0」値を格納し、「異常状態」の場合には「1」を格納する。本発明の実施の形態では正常/異常状態を「0/1」値で表すが、正常/異常状態の識別が可能な情報であれば格納情報は「0/1」値に限らない。

【0052】

リソース詳細情報335には所有状態335 - Aと個別情報335 - Bの項目が含まれる。所有状態335 - AにはVM200のリソースを使用する形態を格納し、個別情報335 - Bにはリソース状態の取得手段を示す情報を格納する。

20

【0053】

例えば図5における所有状態335 - Aの値、「占有」は使用リソースをVM1 200 - Aのみで使用することを意味し(行336)、「共有」は使用リソースをVM1 200 - Aと他VM200とで使用することを意味し(行338)、「論理/SERVER」はVM1 200 - AがVM2 200 - Bと論理的な繋がりのあるグループを形成して「サーバ」の関係になっていることを意味し(行339)、「論理/CLIENT」はVM2 200 - BがVM1 200 - Aと論理的な繋がりのあるグループを形成して「クライアント」の関係になっていることを意味する(行340)。

【0054】

例えば図5における個別情報335 - Bの値、「212」は使用リソースAP1 211 - Aの状態を多重系制御エージェントプログラム212が取得することを意味し(行336)、「EVT」は使用リソースVM1占有メモリ101 - A状態の取得はハードウェア状態を検知する計算機ノード100のハードウェア機能などが送出するシグナル信号を検出することで行う事を意味する(行337)。図5の例では、多重系制御エージェントプログラム212は監視周期333に格納された値の周期で使用リソースAP1 211 - Aの状態を取得するが(行336)、「EVT」の場合はハードウェア機能は使用リソースVM1占有メモリ101 - Aの状態に変化が生じた時にだけシグナル信号を送出するので監視周期333の格納値は使用しない(行337)。

30

<図6：系切替制御テーブル>

40

図6は系切替制御テーブル304の構成例を示す図である。系切替制御テーブル304は多重系制御基盤プログラム300内に保持され、障害発生時に多重系制御基盤プログラム300がVM200の主系VM/従系VMの系切替を行うための規則を格納する。仮想化多重系システム1内の各計算機ノード100は同じ内容の系切替制御テーブル304を保持している。

【0055】

障害内容341には計算機ノード100やVM200の障害発生部位を指定する情報を格納する。図6における「計算機ノード障害」は計算機ノード100自体の動作継続に支障が生じる部位の指定情報であることを示す。VM200の障害発生部位は、全てのVM200を対象に、使用リソースの所有状態種別毎にリソース識別子を格納する。

50

【 0 0 5 6 】

系切替対象 VM 3 4 2 には、障害内容 3 4 1 に障害発生時に系切替対象となる VM 2 0 0 の VM 識別子を格納する。図 6 における「全 VM」は障害が生じた計算機ノード 1 0 0 内に存在する全ての VM 2 0 0 を系切替対象とすることを意味する。「該当 VM のみ」は障害が発生したリソースを使用する VM 2 0 0 だけが系切替対象となることを意味し、実際には当該 VM 2 0 0 の識別子が格納される。「関連全 VM」は障害が発生したリソースを共用する複数の VM 2 0 0、または障害が発生したリソースを使用している VM 2 0 0 及び当該 VM 2 0 0 と論理的なグループ関係にある複数の VM 2 0 0 が一緒に系切替対象となることを意味し、実際には系切替対象となる複数の VM 2 0 0 の識別子が格納される。

10

【 0 0 5 7 】

過半数判断ロジック適用可否 3 4 3 には VM 2 0 0 の系切替実行時に「過半数判断ロジック」の適用の有無を指定する情報を格納する。「要」は「過半数判断ロジック」を適用することを意味し、「否」は「過半数判断ロジック」を適用しないことを意味する。「過半数判断ロジック」は主系 VM / 従系 VM の系切替対象となる VM 2 0 0 を再評価する処理であり、VM 状態管理テーブル 3 0 2 の障害 VM 総数 3 2 2 と系切替対象 VM 3 4 2 の VM 2 0 0 の合計数が、VM 状態管理テーブル 3 0 2 の管理 VM 総数 3 2 1 の半数以上であるかを判断し、系切替対象の VM 2 0 0 を決定する。過半数であれば「過半数判断ロジック」は全 VM を系切替対象と決定し、半数未満であれば系切替対象 VM 3 4 2 に格納された VM 2 0 0 のみを系切替対象とする。

20

【 0 0 5 8 】

Failback 処理適用可否 3 4 4 には系切替を行った VM 2 0 0 に「フェイルバック」処理の適用の可否を指定する情報を格納する。「可」は系切替を行った VM 2 0 0 をフェイルバックしても良いことを意味し、「否」は系切替を行った VM 2 0 0 をフェイルバックしてはいけないことを意味する。「フェイルバック」とは VM 2 0 0 の主系 VM / 従系 VM の状態を系切替実行前に戻す処理のことである。主系 VM の VM 2 0 0 に障害が発生すると多重系制御基盤プログラム 3 0 0 により系切替が実行され、VM 2 0 0 は従系 VM に切替る。「フェイルバック」処理は、系切替後の従系 VM の VM 2 0 0 を、後に障害回復した際に再度系切替を実行して、主系 VM に戻す処理を行う。

30

< 図 7 : VM 1 障害内容と系切替制御内容例 >

図 7 は計算機ノード A 1 0 0 - A に作成された VM 1 2 0 0 - A を対象とした系切替制御テーブル 3 0 4 の例を示した図である。VM 1 2 0 0 - A は、占有リソースとして AP 1 2 1 1 - A と VM 1 占有メモリ 1 0 1 - A を、共有リソースとして VM 1 と VM 2 共有メモリ 1 0 1 - B と VM 1 と VM 2 共有外部通信 I / F 1 0 4 - A を使用し、論理リソースとして VM 2 2 0 0 - B とグループ関係にある。障害内容 3 4 1 には「VM 1 占有リソース障害」として AP 1 2 1 1 - A と VM 1 占有メモリ 1 0 1 - A を格納し（行 3 5 2 ）、「VM 1 共有リソース障害」として VM 1 と VM 2 共有メモリ 1 0 1 - B と VM 1 と VM 2 共有外部通信 I / F 1 0 4 - A を格納し（行 3 5 3 ）、「VM 1 論理リソース障害」として VM 2 2 0 0 - B を格納する。

40

【 0 0 5 9 】

行 3 5 1 は VM 1 2 0 0 - A が存在する計算機ノード A 1 0 0 - A が障害により動作継続不可となった場合の系切替規則で、計算機ノード A 1 0 0 - A 内の全 VM 2 0 0 が系切替対象となり、「過半数判断ロジック」と「フェイルバック」処理は適用しない。

【 0 0 6 0 】

行 3 5 2 は占有リソースの AP 1 2 1 1 - A と VM 1 占有メモリ 1 0 1 - A が障害で使用不能となった場合の系切替規則で、VM 1 2 0 0 - A のみが系切替対象となり、系切替実行時に「過半数判断ロジック」を適用し、占有リソースの障害回復時には「フェイルバック」処理を行う。

【 0 0 6 1 】

行 3 5 3 は共有リソースの VM 1 と VM 2 共有メモリ 1 0 1 - B と VM 1 と VM 2 共有

50

外部通信 I / F 1 0 4 - A が障害で使用不能となった場合の系切替規則で、共に同じリソースを使用する V M 1 2 0 0 - A と V M 2 2 0 0 - B が系切替対象となり、系切替実行時に「過半数判断ロジック」を適用し、共有リソースの障害回復時には「フェイルバック」処理を行う。

【 0 0 6 2 】

行 3 5 4 はグループ関係にある V M 2 2 0 0 - B が障害で動作不能となった場合の系切替規則で、同一グループ内の V M 1 2 0 0 - A と V M 2 2 0 0 - B が系切替対象となり、系切替実行時に「過半数判断ロジック」を適用し、V M 2 2 0 0 - B の障害回復時には「フェイルバック」処理を行う。

< 図 8 : 計算機ノード間通信フォーマット >

図 8 は計算機ノード間通信フォーマット 5 0 0 を示す図である。計算機ノード間通信フォーマット 5 0 0 は異なる計算機ノード 1 0 0 の多重系制御基盤プログラム 3 0 0 の間で送受信される通信データのフォーマットで、計算機ノード 1 0 0 の動作状況を取得する「生死監視」、及び主系 V M / 従系 V M の系切替要求である「フェイルバック」、「異常 V M 系切替」、「全 V M 系切替」の用途で用いられる。

【 0 0 6 3 】

計算機ノード識別子 5 0 1 には通信データを生成した計算機ノード 1 0 0 を特定する識別子を格納する。通信データ生成時には、送受信される当該通信データを識別する情報をシーケンス番号 5 0 2 に格納し、また当該通信データ生成時の時刻情報を送信タイムスタンプ 5 0 3 に格納する。

【 0 0 6 4 】

通信種類 5 0 4 は通信データの用途を識別する情報を格納する。通信用途を識別する情報として、生死 5 0 4 - A、F a i l B a c k 5 0 4 - B、異常 V M 5 0 4 - C、全 V M 5 0 4 - D のうちのの一つを使用する。従系ノードから主系ノードへ「生死監視」を行う場合には生死 5 0 4 - A を使用し、主系ノードから従系ノードへ異常状態になった特定の主系 V M の系切替を指示する「異常 V M 系切替」の場合には異常 V M 5 0 4 - C を使用し、全ての主系 V M の系切替を指示する「全 V M 系切替」の場合には全 V M 5 0 4 - D を使用し、V M 2 0 0 の状態が回復した際に切替っていた主系 V M / 従系 V M の状態を元に戻す「フェイルバック」を指示する場合には F a i l B a c k 5 0 4 - B を使用する。

【 0 0 6 5 】

計算機ノード 1 0 0 の間でデータ通信を行う際には、相手計算機ノード 1 0 0 に自計算機ノード 1 0 0 内の V M 2 0 0 の状態を通知するために、V M 状態管理テーブル内容 5 0 5 には V M 状態管理テーブル 3 0 2 の内容を格納し、リソース管理テーブル内容 5 0 6 にはリソース管理テーブル 3 0 3 の内容を格納する。

「生死監視」などの用途での計算機ノード 1 0 0 間でデータ通信では、相手計算機ノード 1 0 0 状態の異常発生などにより、発信したデータ通信に対する応答が返ってこないなどの応答異常が発生する場合がある。応答異常カウント 5 0 7 には応答異常の発生回数を格納する。

< 図 9 : 多重系制御エージェントプログラムと多重系制御基盤プログラム間の通信フォーマット >

図 9 は多重系制御エージェントプログラムと多重系制御基盤プログラム間の通信フォーマット 6 0 0 を示す図である（以降、プログラム間通信フォーマット 6 0 0 と称する）。プログラム間通信フォーマット 6 0 0 は計算機ノード 1 0 0 内の多重系制御基盤プログラム 3 0 0 と多重系制御エージェントプログラム 2 1 2 の間で送受信されるデータ通信のフォーマットである。

【 0 0 6 6 】

V M 識別子 6 0 1 には通信元 / 先の V M 2 0 0 を特定する情報を格納し、シーケンス番号 6 0 2 には送受信される通信データを識別する情報を格納し、送信タイムスタンプ 6 0 3 には通信データ生成時の時刻情報を格納する。

【 0 0 6 7 】

10

20

30

40

50

多重系制御基盤プログラム300と多重系制御エージェントプログラム212の間でVM200が使用するリソース情報の通知を行う場合には、リソース管理テーブル内容604にリソース管理テーブル303の内容を格納する。

【0068】

多重系制御基盤プログラム300と多重系制御エージェントプログラム212の間で主系VM/従系VMの系切替要求を通知したい場合には、系切替要求605に値を設定する。本発明の形態では、系切替要求がある場合には「1」を設定し、系切替要求が無い場合には「0」を設定するが、系切替要求の有無が識別可能であれば「0/1」値に限られるものではない。

<図10：仮想化多重系システムの全体フローチャート>

10

図10は仮想化多重系システム1で実行される処理の全体フローチャートである。仮想化多重系システム1の全体処理は、多重系制御基盤プログラム300の起動処理及び系切替通知処理S116、多重系制御エージェントプログラム212の起動処理及びリソース監視処理S106と系切替処理、計算機ノード監視部305のメッセージ送信処理S112及びノード監視処理S123より成り立つ。

【0069】

ホストOS400及びプロセッサ102は多重系制御基盤プログラム300及びVM200を計算機ノード100のメモリ101内に展開すると、多重系制御基盤プログラム300の処理を開始する。計算機ノード監視部305及びVM200は多重系制御基盤プログラム300の処理の中から処理を開始される。

20

【0070】

多重系制御基盤プログラム300は処理が開始されると、初めに計算機ノード監視部305の起動を行う(S111)。

【0071】

次に多重系制御基盤プログラム300は記憶装置103からVM状態管理テーブル302とリソース管理テーブル303をメモリ101内に読み込み(S112)、系切替制御テーブル304もメモリ101内に読み込む(S113)。

【0072】

多重系制御基盤プログラム300は制御下にある各VM200のメモリ領域にリソース管理テーブル303をコピーし(S114)、VM200の処理を起動する(S115)

30

【0073】

多重系制御基盤プログラム300はS115の後、系切替通知処理S116を実行する。

【0074】

計算機ノード監視部305は多重系制御基盤プログラム300によって立ち上げられ処理を開始すると、多重系制御基盤プログラム300からの「VM切替」または「フェイルバック」実行の送信要求の有無を判定する(S121)。

【0075】

多重系制御基盤プログラム300からの送信要求が存在する場合(S121の結果が「Y」)はメッセージ送信処理S112を実行し、その後ノード監視処理S123を実行する。多重系制御基盤プログラム300からの送信要求が存在しない場合(S121の結果が「N」)はノード監視処理S123を実行する。

40

【0076】

計算機ノード監視部305はノード監視処理S123の実行終了後は、再び多重系制御基盤プログラム300からの送信要求の有無判定S121から処理を実行する。

【0077】

多重系制御エージェントプログラム212は処理を開始すると、自VM200のメモリ領域内のリソース管理テーブル303を読み込む(S101)。

【0078】

50

多重系制御エージェントプログラム 2 1 2 は読込んだリソース管理テーブル 3 0 3 から、VM 識別子 3 3 1 の格納値が自 VM 2 0 0 の識別子と一致し、かつ個別情報 3 3 5 - B の格納値が「多重系制御エージェントプログラム 2 1 2」である個別リソース管理情報を選別する。そして選別した個別リソース管理情報のリソース識別子 3 3 2 と監視周期 3 3 3 から状態監視対象のリソースと監視周期を取得する。(S 1 0 2)

次に多重系制御エージェントプログラム 2 1 2 は自 VM 2 0 0 の状態を多重系制御基盤プログラム 3 0 0 に通知する。多重系制御エージェントプログラム 2 1 2 は状態監視対象リソースの状態を取得、自 VM 2 0 0 のメモリ領域内のリソース管理テーブル 3 0 3 を更新し、プログラム間通信フォーマット 6 0 0 の通信データを作成して多重系制御基盤プログラム 3 0 0 に送信する。(S 1 0 3)

次に多重系制御エージェントプログラム 2 1 2 は多重系制御基盤プログラム 3 0 0 からの主系 VM / 従系 VM の系切替要求の有無を確認する (S 1 0 4)。

【0079】

系切替要求が有る場合 (S 1 0 4 の結果が「Y」)、多重系制御エージェントプログラム 2 1 2 は系切替処理 S 1 0 5 を実行の後、再び S 1 0 4 を実行する。系切替処理 S 1 0 5 は自 VM 2 0 0 内で動作する AP 2 1 1 に対して系切替を要求する処理である。

【0080】

系切替要求が無い場合 (S 1 0 4 の結果が「N」)、多重系制御エージェントプログラム 2 1 2 はリソース監視処理 S 1 0 6 を実行の後、再び S 1 0 4 を実行する。

< 図 1 1 : リソース監視処理概要 >

図 1 1 はリソース監視処理 S 1 0 6 の概要を示す図である。本実施形態では、多重系制御エージェントプログラム 2 1 2 のリソース状態は、AP 2 1 1 の実行状態を取得することで判断する。

【0081】

AP 2 1 1 が動作する時には、VM 2 0 0 内の任意の場所 (例えば / var / run) に AP 2 1 1 の実行状態やエラー情報を出力する APdog 7 0 0 ファイルが確保される。APdog 7 0 0 ファイルは AP 2 1 1 以外のプログラムからも参照が可能なファイルである。AP 2 1 1 は自身の実行状態を定期的に APdog 7 0 0 ファイルに出力し、多重系制御エージェントプログラム 2 1 2 は APdog 7 0 0 ファイルを参照することで AP 2 1 1 の実行状態の監視を行う。

【0082】

多重系制御エージェントプログラム 2 1 2 のリソース状態の監視は、例えば AP 2 1 1 リソースの稼動状態は APdog 7 0 0 が定期的に更新されていることを確認することで認識し、外部通信 I / F 1 0 4 リソースの異常は AP 2 1 1 が APdog 7 0 0 に出力した端末 2 との外部通信異常情報を取得することで行い、論理グループ関係にある他の VM 2 0 0 リソースの異常は AP 2 1 1 が APdog 7 0 0 に出力した他の AP 2 1 1 との内部通信異常情報を取得することで行う。

リソース状態の変化を検知した多重系制御エージェントプログラム 2 1 2 は、リソース状態を多重系制御基盤プログラム 3 0 0 に通知する。

< 図 1 2 : リソース監視処理のフローチャート >

図 1 2 はリソース監視処理 S 1 0 6 の一例を示すフローチャートである。多重系制御エージェントプログラム 2 1 2 は監視対象リソースに異常を検出すると、異常状態を多重系制御基盤プログラム 3 0 0 に通知する。

【0083】

多重系制御エージェントプログラム 2 1 2 は起動すると、初めに APdog 7 0 0 を参照してリソース状態の監視を開始し (S 2 0 1)、監視対象リソースの異常の有無を判断する (S 2 0 2)。

【0084】

監視対象リソースの内のいずれかに異常が認められた場合 (S 2 0 2 の結果が「Y」)、VM 2 0 0 内に保持されているリソース管理テーブル 3 0 3 の該当リソースのリソース

10

20

30

40

50

生死状態フラグ 3 3 4 に死状態である「1」を設定し (S 2 0 3)、リソース管理テーブル 3 0 3 をプログラム間通信フォーマット 6 0 0 のリソース管理テーブル内容 6 0 4 に設定して通信データを作成して多重系制御基盤プログラム 3 0 0 に送信して (S 2 0 4)、リソース監視処理 S 1 0 6 を終了する (S 2 0 5)。

【0085】

監視対象リソースの内のいずれにも異常が認められない場合 (S 2 0 2 の結果が「N」) は、リソース監視処理 S 1 0 6 を終了する (S 2 0 5)。

【0086】

リソース監視処理 S 1 0 6 を終了した後、多重系制御エージェントプログラム 2 1 2 は、多重系制御基盤プログラム 3 0 0 からの主系 VM / 従系 VM の系切替要求の有無を確認するために図 1 0 に示した S 1 0 4 を起動する。

10

< 図 1 3 : 系切替通知処理のフローチャート >

図 1 3 は本発明の実施の形態の系切替通知処理 S 1 1 6 のフローチャートである。多重系制御基盤プログラム 3 0 0 は系切替通知処理 S 1 1 6 で多重系制御エージェントプログラム 2 1 2 からの VM 2 0 0 状態の通知を取得し、計算機ノード監視部 3 0 5 に主系 VM / 従系 VM の系切替要求を行う。

【0087】

多重系制御基盤プログラム 3 0 0 は系切替通知処理 S 1 1 6 を開始すると、初めに多重系制御エージェントプログラム 2 1 2 からの VM 2 0 0 状態の通知、及び計算機ノード 1 0 0 やホスト OS 4 0 0 からのハードウェア状態検出の通知を待つ (S 3 0 1)。

20

【0088】

VM 2 0 0 状態やハードウェア状態の通知内容を解釈し、計算機ノード 1 0 0 内のある VM 2 0 0 に新たに異常が認められた場合 (S 3 0 2 の結果が「Y」)、VM 状態管理テーブル 3 0 2 の該当 VM 2 0 0 の主 / 従系状態フラグ 3 2 5 には従系 VM 状態を意味する「S」を設定して、生死状態フラグ 3 2 6 には死状態を意味する「1」を設定する (S 3 1 1)。

【0089】

そして VM 状態管理テーブル 3 0 2 の障害 VM 総数 3 2 2 の値を「1」増加させる (S 3 1 2)。

【0090】

次に多重系制御基盤プログラム 3 0 0 は計算機ノード構成管理テーブル 3 0 1 の自計算機ノード 1 0 0 に該当する主 / 従系状態フラグ 3 1 2 を参照し、自計算機ノード 1 0 0 の主系ノード / 従系ノード状態を判別する (S 3 1 3)。

30

【0091】

主 / 従系状態フラグ 3 1 2 の値が「S」で自計算機ノード 1 0 0 が従系ノード状態である場合 (S 3 1 3 の結果が「N」)、多重系制御基盤プログラム 3 0 0 は何も処理することなく、再び VM 2 0 0 状態やハードウェア状態の通知を待つために S 3 0 1 を起動する。

【0092】

主 / 従系状態フラグ 3 1 2 の値が「M」で自計算機ノード 1 0 0 が主系ノード状態である場合 (S 3 1 3 の結果が「Y」)、多重系制御基盤プログラム 3 0 0 は異常 VM 2 0 0 の主系 VM / 従系 VM の系切替要求を行うために計算機ノード間通信フォーマット 5 0 0 の通信パケットを作成する (S 3 3 1)。

40

【0093】

多重系制御基盤プログラム 3 0 0 は、異常 VM 2 0 0 のリソース障害内容が「過半数判断ロジック」の適用を要するものである場合には、VM 状態管理テーブル 3 0 2 の障害 VM 総数 3 2 2 と異常 VM 2 0 0 の数量の和が管理 VM 総数 3 2 1 の半数以上であるかの判断を行う (S 3 3 2)。

【0094】

障害 VM 2 0 0 の総数が過半数となる場合には (S 3 3 2 の結果が「Y」)、正常な V

50

M 2 0 0 も含めて全ての V M 2 0 0 を系切替対象とするために、計算機ノード間通信フォーマット 5 0 0 の通信パケットの通信種類 5 0 4 には全 V M 5 0 4 - D の識別値を設定する (S 3 4 1) 。

【 0 0 9 5 】

障害 V M 2 0 0 の総数が過半数とならない場合には (S 3 3 2 の結果が「 N 」)、異常 V M 2 0 0 のみを系切替対象とするために、計算機ノード間通信フォーマット 5 0 0 の通信パケットの通信種類 5 0 4 には異常 V M 5 0 4 - C の識別値を設定する (S 3 5 1) 。

【 0 0 9 6 】

計算機ノード間通信フォーマット 5 0 0 の通信パケットの作成が完了すると、多重系制御基盤プログラム 3 0 0 は計算機ノード監視部 3 0 5 に送信要求を行う (S 3 4 2) 。

【 0 0 9 7 】

次に異常 V M 2 0 0 を停止させるためにリセットを行い (S 3 4 3)、リソース管理テーブル 3 0 3 を参照して異常 V M 2 0 0 が使用しているリソースに対してもリセットを行い (S 3 4 4)、停止させた異常 V M 2 0 0 を再起動させるために S 1 1 5 を起動する。

【 0 0 9 8 】

V M 2 0 0 状態やハードウェア状態の通知内容から計算機ノード 1 0 0 内のある V M 2 0 0 に異常が認められない場合 (S 3 0 2 の結果が「 N 」)、当該 V M 2 0 0 は「死状態」から回復したものであるかを確認するために、V M 状態管理テーブル 3 0 2 に記録された当該 V M 2 0 0 の生死状態フラグ 3 2 6 の値を参照する (S 3 2 1) 。

【 0 0 9 9 】

生死状態フラグ 3 2 6 の値が「 0 」で当該 V M 2 0 0 が「正状態」であると判断されていた場合には (S 3 2 1 の結果が「 N 」)、当該 V M 2 0 0 の状態に変化は起きていないため、多重系制御基盤プログラム 3 0 0 は再び多重系制御エージェントプログラム 2 1 2 からの通知や計算機ノード 1 0 0 / ホスト O S 4 0 0 からのハードウェア状態検出の通知を待つために、S 3 0 1 を起動する。

【 0 1 0 0 】

生死状態フラグ 3 2 6 の値が「 1 」で当該 V M 2 0 0 が「死状態」であると判断されていた場合には (S 3 2 1 の結果が「 Y 」)、当該 V M 2 0 0 は状態が正常に回復したものと判断されるため、多重系制御エージェントプログラム 2 1 2 から受信したプログラム間通信フォーマット 6 0 0 の通信データのリソース管理テーブル内容 6 0 4 の該当 V M 2 0 0 の状態を、多重系制御基盤プログラム 3 0 0 内の V M 状態管理テーブル 3 0 2 に反映する (S 3 2 2) 。

【 0 1 0 1 】

次に当該 V M 2 0 0 の主系 V M / 従系 V M の状態のフェイルバック要求を行うために、計算機ノード構成管理テーブル 3 0 1 の主 / 従系状態フラグ 3 1 2 を参照して、自計算機ノード 1 0 0 の主系ノード / 従系ノード状態を判別する (S 3 2 3) 。

【 0 1 0 2 】

主 / 従系状態フラグ 3 1 2 の値が「 S 」で自計算機ノード 1 0 0 が従系ノードであった場合には (S 3 2 3 の結果が「 N 」)、従系ノードが自らフェイルバック要求は行わないため、再び多重系制御エージェントプログラム 2 1 2 からの通知や計算機ノード 1 0 0 / ホスト O S 4 0 0 からのハードウェア状態検出の通知を待つために S 3 0 1 を起動する。

【 0 1 0 3 】

主 / 従系状態フラグ 3 1 2 の値が「 M 」で自計算機ノード 1 0 0 が主系ノードであった場合には (S 3 2 3 の結果が「 Y 」)、計算機ノード間通信フォーマット 5 0 0 の通信パケットを作成し (S 3 2 4)、フェイルバック要求であるため通信種類 5 0 4 には F a i l B a c k 5 0 4 - B の識別値を設定し (S 3 2 5)、計算機ノード監視部 3 0 5 に作成した通信パケットの送信要求を行い (S 3 2 6)、再び多重系制御エージェントプログラム 2 1 2 からの通知や計算機ノード 1 0 0 / ホスト O S 4 0 0 からのハードウェア状態検出の通知を待つために S 3 0 1 を起動する。

< 図 1 4 : メッセージ送信処理のフローチャート >

10

20

30

40

50

図14はメッセージ送信処理S112のフローチャートである。計算機ノード監視部305は多重系制御基盤プログラム300から受けた他計算機ノード100への「フェイルバック」要求、「異常VM切替」要求、「全VM切替」要求を、メッセージ送信処理S112で発信する。

【0104】

主系VM/従系VMの「全VM切替」が実行される場合、切替対象となるVM200の中には本来切替を必要としない正常状態のVM200が含まれることもある。「全VM切替」を実行するために、正常VM200の系切替を行うための事前処理はメッセージ送信処理S112内で実行する。

【0105】

計算機ノード監視部305はメッセージ送信処理S112を開始すると、初めに通信要求の種類が「全VM切替」とその他を特定するために、計算機ノード間通信フォーマット500の通信要求パケットの通信種類504を確認する(S401)。

【0106】

通信種類504の値が全VM504-Dの識別値ではなく通信要求が「全VM切替」以外の場合(S401の結果が「N」)、系切替対象は状態が異常に転じた、または状態が回復したVM200に限られるため、S411を起動して通信要求パケットを仮想化多重系システム1内の他の計算機ノード100にマルチキャスト送信する。

【0107】

通信種類504の値が全VM504-Dの識別値であり通信要求が「全VM切替」の場合は(S401の結果が「Y」)、正常状態のVM200を系切替するための処理を起動する。

【0108】

計算機ノード監視部305はVM状態管理テーブル302の主/従系状態フラグ325の値が「M」の主系VMのVM200を検索する(S402)。

【0109】

検索された該当VM200に系切替を要求するために、プログラム間通信フォーマット600の系切替要求605に「1」を設定した通信データを作成し(S403)、該当VM200の多重系制御エージェントプログラム212に対して送信する(S404)。

【0110】

次にVM状態管理テーブル302の全てのVM200の主/従系状態フラグ325に従系VMを意味する「S」を設定し、また生死状態フラグ326を死状態である「1」に変更する(S405)。

【0111】

そしてVM状態管理テーブル302の障害VM総数322の値を「0」にリセットして(S406)、「全VM切替」の通信要求パケットを仮想化多重系システム1内の他の計算機ノード100にマルチキャスト送信する(S411)。

【0112】

S411で通信要求パケットを送信し終えた後は、メッセージ送信処理S112を終了し(S421)、次のノード監視処理S123を起動する。

<図15：ノード監視処理のフローチャート>

図15はノード監視処理S123のフローチャートである。主系ノード及び従系ノードの計算機ノード100内で動作する多重系制御基盤プログラム300の計算機ノード監視部305の間では、「生死監視」及び「フェイルバック」、「異常VM切替」、「全VM切替」要求の種別の計算機ノード間通信が行われる。計算機ノード監視部305は受信した計算機ノード間通信データの種別を判別し、種別に応じた処理の起動を行う。

【0113】

計算機ノード監視部305はノード監視処理S123の処理を開始すると、初めに他の計算機ノード100からの計算機ノード間通信パケットを受信しているかを判別する。従系ノードの計算機ノード監視部305は、「生死監視」に対する応答、および「異常VM

10

20

30

40

50

切替」、「全VM切替」、「フェイルバック」要求の計算機ノード間通信データを受信する。主系ノードの計算機ノード監視部305は、「生死監視」の計算機ノード100状態取得要求、および「異常VM切替」、「全VM切替」、「フェイルバック」に対する応答の計算機ノード間通信データを受信する。(S501)。

【0114】

計算機ノード間通信パケットを受信している場合(S501の結果が「Y」)、計算機ノード監視部305は計算機ノード構成管理テーブル301の主/従系状態フラグ312を参照し、自計算機ノード100の主系ノード/従系ノード状態を判別する。主/従系状態フラグ312が「M」で自計算機ノード100が主系ノードの場合(S502の結果が「Y」)は主系処理S503を起動し、主/従系状態フラグ312が「S」で自計算機ノード100が従系ノードの場合(S502の結果が「N」)は従系処理S504を起動する。(S502)

主系処理S503または従系処理S504が終了すると、計算機ノード監視部305はノード監視処理S123を終了する(S531)。

【0115】

計算機ノード間通信データを受信していない場合(S501の結果が「N」)、計算機ノード監視部305は計算機ノード構成管理テーブル301の主/従系状態フラグ312を参照し、自計算機ノード100の主系ノード/従系ノード状態を判別する。

【0116】

主/従系状態フラグ312が「M」で自計算機ノード100が主系ノードの場合(S511の結果が「N」)はノード監視処理S123を終了し、主/従系状態フラグ312が「S」で自計算機ノード100が従系ノードの場合(S511の結果が「Y」)は「生死監視」要求に対する応答が返ってこないことであるため次のS512を起動する。(S511)

S512では、「生死監視」要求に対する応答の無い計算機ノード100の計算機ノード構成管理テーブル301の生死確定カウンタ316の値と、計算機ノード間通信フォーマット500の応答異常カウンタ507の値を比較し、当該計算機ノード100の生死状態の判別を行う。

【0117】

応答異常カウンタ507の値が生死確定カウンタ316の値未満の場合(S512の結果が「N」)は当該計算機ノード100がまだ正常状態であると判断して生死監視処理S513を起動し、応答異常カウンタ507の値が生死確定カウンタ316の値以上の場合(S512の結果が「Y」)は当該計算機ノード100が異常状態にあると判断して全VM切替処理S521を起動する。

【0118】

生死監視処理S513または全VM切替処理S521が終了すると、計算機ノード監視部305はノード監視処理S123を終了する(S531)。

【0119】

ノード監視処理S123の終了後、計算機ノード監視部305は多重系制御基盤プログラム300からの送信要求の有無判定S121を起動する。

<図16：生死監視処理のフローチャート>

図16は本発明の実施の形態の生死監視処理S513のフローチャートである。従系ノードである計算機ノード100の多重系制御基盤プログラム300の計算機ノード監視部305は生死監視処理S513で「生死監視」要求を実行する。

【0120】

従系ノードの計算機ノード監視部305は生死監視処理S513を起動すると、計算機ノード構成管理テーブル301の主系ノードである計算機ノード100の監視周期315を参照する(S601)。

【0121】

前回の「生死監視」要求からの経過時間が監視周期315を超えていない場合(S60

10

20

30

40

50

1の結果が「N」)は生死監視処理S513を終了する。

【0122】

経過時間が監視周期315を超えている場合(S601の結果が「Y」)は「生死監視」要求を行うために次のS602を起動する。

【0123】

主系ノードである計算機ノード100に送る計算機ノード間通信フォーマット500のパケットを作成するために、自計算機ノード100のVM状態管理テーブル302をVM状態管理テーブル内容505に、リソース管理テーブル303の内容をリソース管理テーブル内容506に設定する(S602)。そして通信種類504には生死504-Aの識別値を設定する(S603)。

【0124】

次に作成したパケットを、仮想化多重系システム1内の主系ノードである計算機ノード100に対してマルチキャストで送信する(S604)。

【0125】

「生死監視」要求のパケット送信が完了すると生死監視処理S513を終了する。

<図17：全VM切替処理のフローチャート>

図17は全VM切替処理S521のフローチャートである。従系ノードである計算機ノード100で動作する多重系制御基盤プログラム300は、全VM切替処理S521を実行して、当該計算機ノード100内の全てのVM200を主系VM状態に、当該計算機ノード100を主系ノード状態に変更する。

【0126】

多重系制御基盤プログラム300は全VM切替処理S521を起動すると、初めに計算機ノード構成管理テーブル301の自計算機ノード100に該当する主/従系状態フラグ312に「M」を設定し、主系ノード状態に変更する(S701)。

【0127】

次に自計算機ノード100内の全てのVM200の系切替を行うために、プログラム間通信フォーマット600の系切替要求605に系切替の実行を指定する「1」を設定した、通信データを作成する(S702)。

【0128】

次にS703で自計算機ノード100内の全てのVM200の多重系制御エージェントプログラム212に対して作成したプログラム間通信フォーマット600の通信データを送信した後、S704で全VM切替処理S521を終了する。

<図18：主系処理のフローチャート>

図18は主系処理S503のフローチャートである。ノード監視処理S123のS502から起動された主系処理S503は、従系ノードである計算機ノード100からの「生死監視」要求、および「フェイルバック」、「異常VM切替」、「全VM切替」の応答に対する処理を行う。

【0129】

処理を開始した主系処理S503は、初めに計算機ノード間通信フォーマット500の受信パケットの通信種類504を参照し(S801)、「生死監視」要求の受信パケットであるかを判別する。

【0130】

受信パケットの通信種類504が生死504-Aの識別値である場合(S802の結果が「Y」)、S811を起動して自計算機ノード100の状態の返信処理を開始する。受信パケットの通信種類504が生死504-Aの識別値でなかった場合(S802の結果が「N」)、受信パケットの種別を更に特定するためにS821を起動する。(S802)

S811では計算機ノード間通信フォーマット500の返信パケットに自計算機ノード100のVM状態管理テーブル302の内容をVM状態管理テーブル内容505に設定し、リソース管理テーブル303の内容をリソース管理テーブル内容506に設定する。

10

20

30

40

50

【 0 1 3 1 】

次に作成した「生死監視」要求に対する返信パケットを、S 8 1 2 で従系ノードである計算機ノード 1 0 0 にマルチキャストで送信した後、主系処理 S 5 0 3 を終了する。

【 0 1 3 2 】

S 8 2 1 では受信パケットが「フェイルバック」の応答であるかを判別する。受信パケットの通信種類 5 0 4 が F a i l B a c k 5 0 4 - B の識別値である場合 (S 8 2 1 の結果が「Y」)、主系ノードである計算機ノード 1 0 0 での「フェイルバック」処理を行うために S 8 3 1 を起動する。受信パケットの通信種類 5 0 4 が F a i l B a c k 5 0 4 - B の識別値でない場合 (S 8 2 1 の結果が「N」)、受信パケットの種別を更に特定するために S 8 2 2 を起動する。

10

【 0 1 3 3 】

S 8 3 1 では V M 状態管理テーブル 3 0 2 の主 / 従系状態フラグ 3 2 5 を参照して、値が「S」である従系 V M の V M 2 0 0 を探索し、これを主系 V M への系切替対象とする。

【 0 1 3 4 】

次にプログラム間通信フォーマット 6 0 0 の系切替要求 6 0 5 に系切替要求である「1」を設定した通信データを作成する (S 8 3 2) 。

【 0 1 3 5 】

次の S 8 3 3 において、S 8 3 1 で探索した系切替対象の V M 2 0 0 の多重系制御エージェントプログラム 2 1 2 に S 8 3 2 で作成した通信データを送信した後、主系処理 S 5 0 3 を終了する。

20

【 0 1 3 6 】

S 8 2 2 では受信パケットが「全 V M 切替」の応答であるかを判別する。受信パケットの通信種類 5 0 4 が全 V M 5 0 4 - D の識別値である場合 (S 8 2 2 の結果が「Y」)、S 8 2 3 を起動して計算機ノード構成管理テーブル 3 0 1 の自計算機ノード 1 0 0 に該当する主 / 従系状態フラグ 3 1 2 に「S」を設定することで自計算機ノード 1 0 0 を従系ノード状態に変更した後、主系処理 S 5 0 3 を終了する。

【 0 1 3 7 】

受信パケットの通信種類 5 0 4 が全 V M 5 0 4 - D の識別値でない場合 (S 8 2 2 の結果が「N」) は主系処理 S 5 0 3 を終了する。

< 図 1 9 : 従系処理のフローチャート >

30

図 1 9 は従系処理 S 5 0 4 のフローチャートである。ノード監視処理 S 1 2 3 の S 5 0 2 から起動された従系処理 S 5 0 4 は、主系ノードである計算機ノード 1 0 0 からの「生死監視」への応答、及び「フェイルバック」、「異常 V M 切替」、「全 V M 切替」要求に対する処理を行う。

【 0 1 3 8 】

処理を開始した従系処理 S 5 0 4 は、初めに計算機ノード間通信フォーマット 5 0 0 の受信パケットの通信種類 5 0 4 を参照し (S 9 0 1)、「生死監視」への応答の受信パケットであるかを判別する。受信パケットの通信種類 5 0 4 が生死 5 0 4 - A の識別値である場合 (S 9 0 2 の結果が「Y」) は S 9 1 1 を起動して「生死監視」の処理を行い、通信種類 5 0 4 が生死 5 0 4 - A の識別値で無い場合 (S 9 0 2 の結果が「N」) は受信パケットの種別を更に特定するために S 9 2 1 を起動する。 (S 9 0 2)

40

S 9 1 1 において、「生死監視」要求に対する返答が合ったことから主系ノードの計算機ノード 1 0 0 は正常状態にあると判断できるため、計算機ノード間通信フォーマット 5 0 0 の通信パケットの応答異常カウント 5 0 7 の値を「0」にリセットした後、従系処理 S 5 0 4 を終了する。

【 0 1 3 9 】

S 9 2 1 では受信パケットが「フェイルバック」要求であるかを判別する。受信パケットの通信種類 5 0 4 が F a i l B a c k 5 0 4 - B の識別値である場合 (S 9 2 1 の結果が「Y」)、従系ノードである計算機ノード 1 0 0 での「フェイルバック」処理を行うために S 9 4 1 を起動する。受信パケットの通信種類 5 0 4 が F a i l B a c k 5 0 4 - B

50

の識別値で無い場合（S 9 2 1の結果が「N」）、受信パケットの種別を更に特定するためにS 9 2 2を起動する。

【0140】

S 9 4 1では計算機ノード間通信フォーマット500の受信パケットのVM状態管理テーブル内容505に格納されているVM状態管理テーブル302情報を参照して、主/従系状態フラグ325の値が「S」である従系VM状態にある主系ノードの計算機ノード100内のVM200を検索する。

【0141】

次に、自計算機ノード100の識別子及び自計算機ノード100内のVM200の識別子を、S 9 4 1で検索したVM200の主/従系構成VM情報327のホスト計算機ノード識別子327-A及びVM識別子327-Bと比較し、系切替に該当する自計算機ノード100内のVM200を検索する（S 9 4 2）。

10

【0142】

次に該当VM200に系切替要求を行うために、S 9 4 3でプログラム間通信フォーマット600の系切替要求605に系切替要求である「1」を設定した通信データを作成し、S 9 4 4で該当VM200の多重系制御エージェントプログラム212に作成した通信データを送信する。

【0143】

次に主系ノードの計算機ノード100に「フェイルバック」の応答を行うために、S 9 4 5で自計算機ノード100のVM状態管理テーブル302とリソース管理テーブル303の内容を計算機ノード間通信フォーマット500のVM状態管理テーブル内容505とリソース管理テーブル内容506に格納した通信パケットを作成し、次のS 9 4 6では通信パケットの通信種類504にFail Back 504-Bの識別値を設定する。

20

【0144】

次のS 9 4 7で、作成した通信パケットを仮想化多重系システム1内の主系ノードである計算機ノード100にマルチキャストで送信した後、従系処理S 5 0 4を終了する。

【0145】

S 9 2 2では受信パケットが「全VM切替」要求であるかを判別する。受信パケットの通信種類504が全VM504-Dの識別値である場合（S 9 2 2の結果が「Y」）、従系ノードである計算機ノード100での「全VM切替」処理を行うために全VM切替処理S 5 2 1を起動する。受信パケットの通信種類504が全VM504-Dの識別値で無い場合（S 9 2 2の結果が「N」）、従系ノードである計算機ノード100での「異常VM切替」処理を行うためにS 9 2 3を起動する。

30

【0146】

全VM切替処理S 5 2 1の終了後、主系ノードの計算機ノード100に「全VM切替」の応答を行うために、S 9 3 2で自計算機ノード100のVM状態管理テーブル302とリソース管理テーブル303の内容を計算機ノード間通信フォーマット500のVM状態管理テーブル内容505とリソース管理テーブル内容506に格納した通信パケットを作成し、次のS 9 3 3では通信パケットの通信種類504に全VM504-Dの識別値を設定する。

40

【0147】

次のS 9 3 4で、作成した通信パケットを仮想化多重系システム1内の主系ノードである計算機ノード100にマルチキャストで送信した後、従系処理S 5 0 4を終了する。

【0148】

S 9 2 3では、「異常VM切替」処理を行うために、計算機ノード間通信フォーマット500の受信パケットのVM状態管理テーブル内容505に格納されているVM状態管理テーブル302情報を参照して、主/従系状態フラグ325の値が「S」である従系VM状態にある主系ノードの計算機ノード100内のVM200を検索する。

【0149】

次に、自計算機ノード100の識別子及び自計算機ノード100内のVM200の識別

50

子を、S 9 2 3 で検索した V M 2 0 0 の主 / 従系構成 V M 情報 3 2 7 のホスト計算機ノード識別子 3 2 7 - A 及び V M 識別子 3 2 7 - B と比較し、系切替に該当する自計算機ノード 1 0 0 内の V M 2 0 0 を検索する (S 9 2 4)。

【 0 1 5 0 】

次に該当 V M 2 0 0 に系切替要求を行うために、S 9 2 5 でプログラム間通信フォーマット 6 0 0 の系切替要求 6 0 5 に系切替要求である「 1 」を設定した通信データを作成し、次の S 9 2 6 で該当 V M 2 0 0 の多重系制御エージェントプログラム 2 1 2 に作成した通信データを送信し、従系処理 S 5 0 4 を終了する。

【 0 1 5 1 】

以上の実施形態によれば、物理計算機上で障害検知した主系モードの仮想マシンが搭載済みの仮想マシンよりも半数未満の場合には、仮想マシン毎に対して系切替え（フェイルオーバー）を行った後にフェイルバックする。一方、物理計算機上で障害検知した主系モードの仮想マシンが搭載済みの仮想マシンよりも半数以上の場合には、全ての仮想マシンに対して系切替えを行う。以上によって、単一の主系モードの仮想マシンが障害を起こしたとしても、複数の主系モードの仮想マシンが障害を起こしたとしても、つまり、いずれの場合にも、最終的には物理計算機の片方に全ての仮想マシンが主系モードで稼働するようになり、利用者にとって物理計算機の全体の動作モードが管理しやすくなる。その結果、物理計算機上の全ての仮想マシンを主系モードで稼働させながら、もう一方の物理計算機に対してハードウェアやソフトウェアの保守を行うことが可能となる。

10

【 符号の説明 】

20

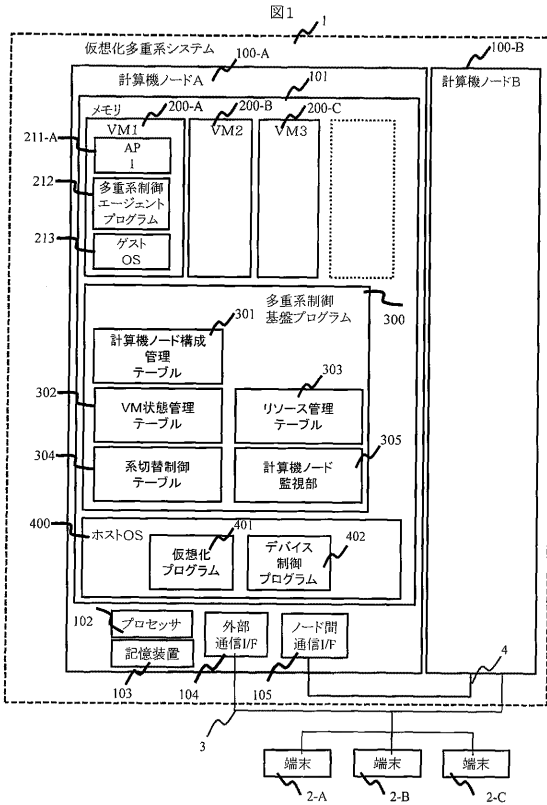
【 0 1 5 2 】

- 1 ... 仮想化多重系システム
- 2 ... 端末
- 3 ... 外部ネットワーク
- 4 ... 内部ネットワーク
- 1 0 0 ... 計算機ノード
- 1 0 1 ... メモリ
- 1 0 2 ... プロセッサ
- 1 0 3 ... 記憶装置、
- 1 0 4 ... 外部通信 I / F
- 1 0 5 ... ノード間通信 I / F
- 2 0 0 ... V M
- 2 0 1 ... V M
- 2 1 1 ... A P
- 2 1 2 ... 多重系制御エージェントプログラム
- 2 1 3 ... ゲスト O S
- 3 0 0 ... 多重系制御基盤プログラム
- 3 0 1 ... 計算機ノード構成管理テーブル
- 3 0 2 ... V M 状態管理テーブル
- 3 0 3 ... リソース管理テーブル
- 3 0 4 ... 系切替制御テーブル
- 3 0 5 ... 計算機ノード監視部
- 4 0 0 ... ホスト O S
- 4 0 1 ... 仮想化プログラム
- 4 0 2 ... デバイス制御プログラム
- 7 0 0 ... A P d o g

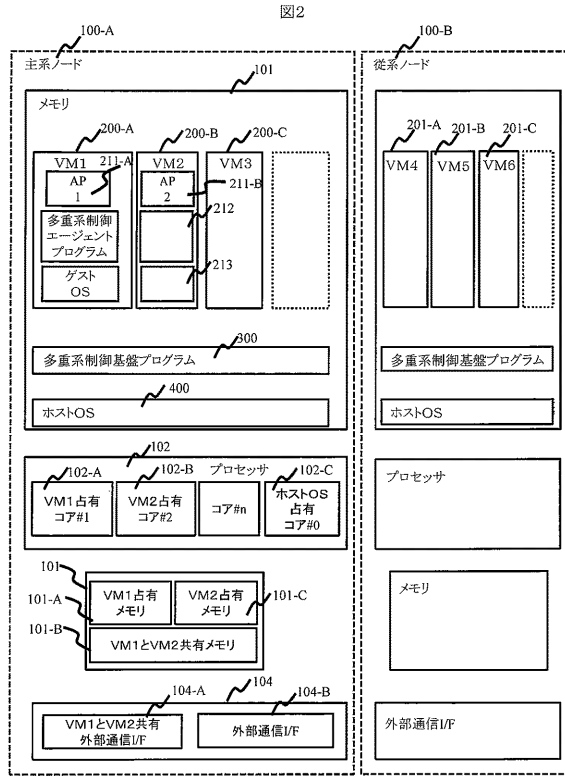
30

40

【 図 1 】



【 図 2 】



【 図 3 】

図3

計算機ノード構成管理テーブル 301

計算機ノード識別子	主/従系状態フラグ	マルチキャストアドレス	生死状態フラグ	監視周期 [ms]	生死確定カウンタ
100-A	M	226.94.1.1	0	1000	3
100-B	S	226.94.1.1	0	1000	3

【 図 4 】

図4

VM状態管理テーブル 302

管理VM総数	障害VM総数			主/従系構成VM情報	
	VM識別子	VM識別子	VM識別子	327-A	327-B
321	323	324	325	326	327
VM1	100-A	M	0	100-B	VM4
VM2	100-A	M	0	100-B	VM5
VM3	100-A	M	0	100-B	VM6

【 図 5 】

図5

リソース管理テーブル 303

VM識別子	リソース識別子	監視周期 [ms]	リソース生死状態フラグ	335-Aリソース詳細情報		
				335-B 所有状態	335-B 個別情報	
336	211	100	0	占有	212	
337	VM1	101-A	0	占有	EVT	
338	VM1	101-B	0	共有	EVT	
	VM1	104-A	50	共有	212	
339	VM1	VM2	100	0	論理/SERVER	212
	VM2	AP2	100	0	占有	212
	VM2	101-C	0	0	占有	EVT
	VM2	101-B	0	0	共有	EVT
	VM2	104-A	50	0	共有	212
340	VM2	VM1	100	0	論理/CLIENT	212

【 図 6 】

図6

系切替制御テーブル 304

障害内容	系切替対象VM	過半数判断ロジック適用要否	Failback処理適用可否
計算機ノード障害	全VM	否	否
VM占有リソース障害	該当VMのみ	要	可
VM共有リソース障害	関連全VM	要	可
VM論理リソース障害	関連全VM	要	可

【 図 7 】

図7

VM1障害内容と系切替制御内容例

障害内容	系切替対象VM	過半数判断 ロジック 適用可否	Failback処理 適用可否
計算機ノード障害	全VM切替	否	否
VM1占有 リソース障害	AP1	VM1切替	要
	101-A		
VM1共有 リソース障害	101-B	VM1、VM2切替	要
	104-A		
VM1論理 リソース障害	VM2	VM1、VM2切替	要

【 図 9 】

図9

プログラム間通信フォーマット

VM 識別子	シーケンス 番号	送信タイム スタンプ	リソース管理 テーブル内容	系切替 要求
323	001	HHMMSS (例:081203)	303	0

【 図 8 】

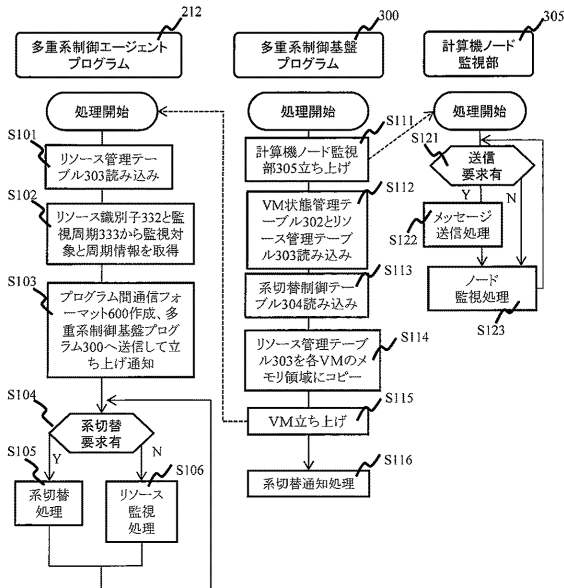
図8

計算機ノード間通信フォーマット

計算機 ノード 識別子	シーケンス 番号	送信 タイム スタンプ	通信 種類	VM 状態 管理 テーブル 内容	リソース 管理 テーブル 内容	応答 異常 カウン ト
200-A	001	HHMMSS	生死	302	303	0
			系切替			
			Fail Back			
			異常 VM			
			全 VM			

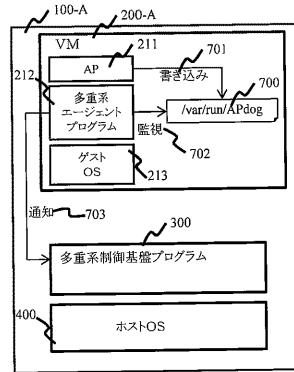
【 図 1 0 】

図10



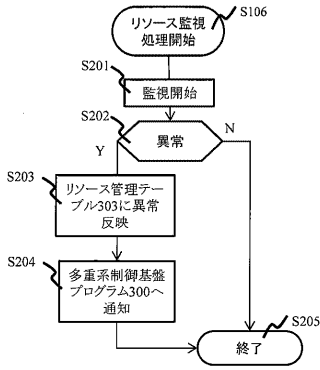
【 図 1 1 】

図11



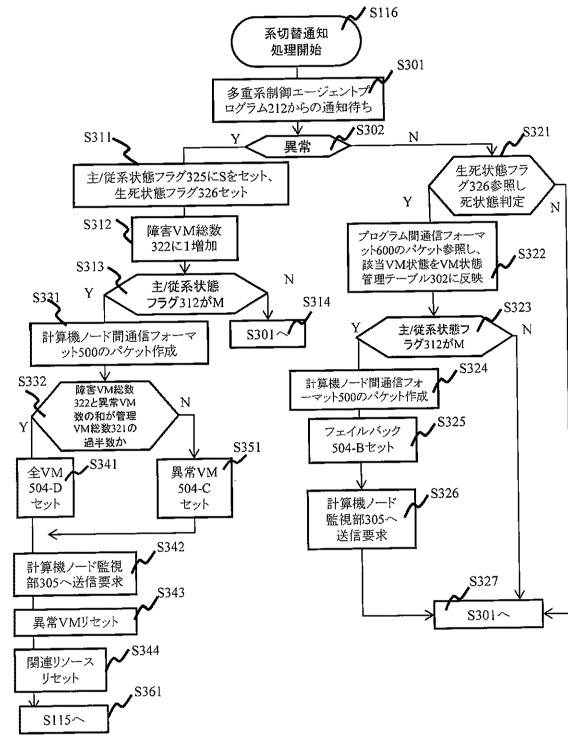
【 図 1 2 】

図12



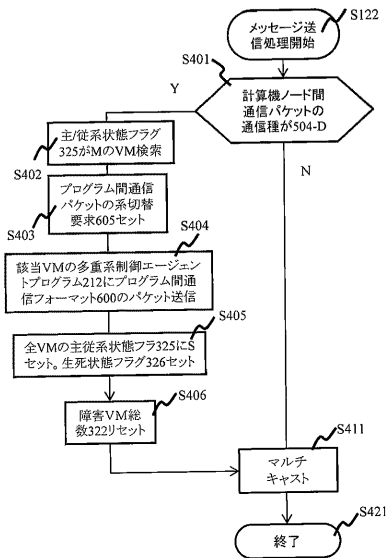
【 図 1 3 】

図13



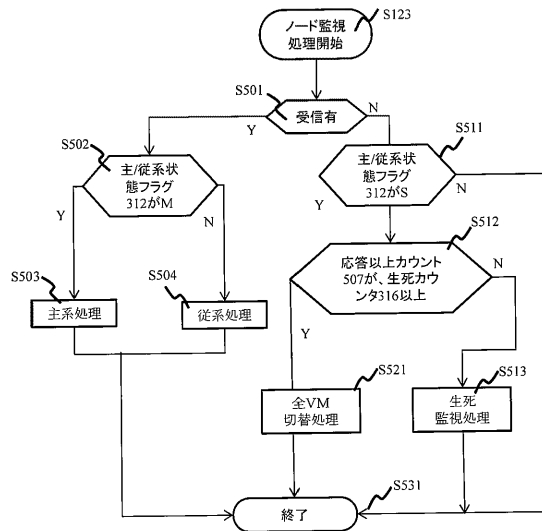
【 図 1 4 】

図14



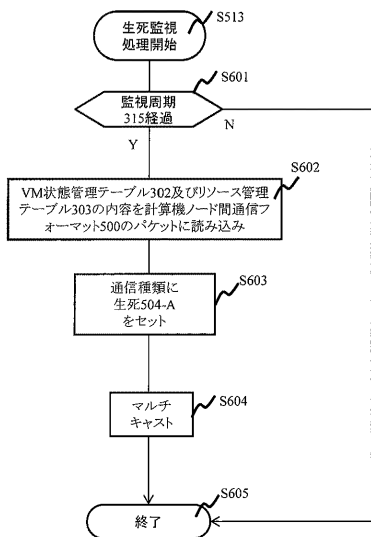
【 図 1 5 】

図15



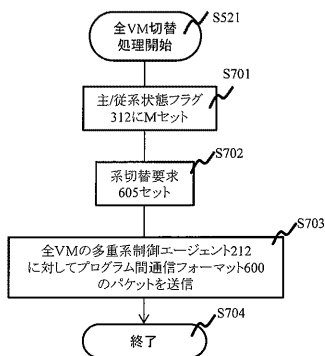
【 図 1 6 】

図16



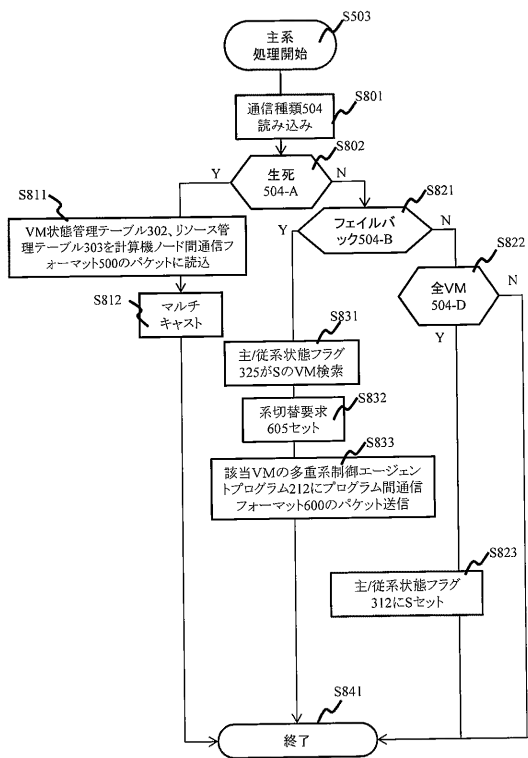
【 図 1 7 】

図17



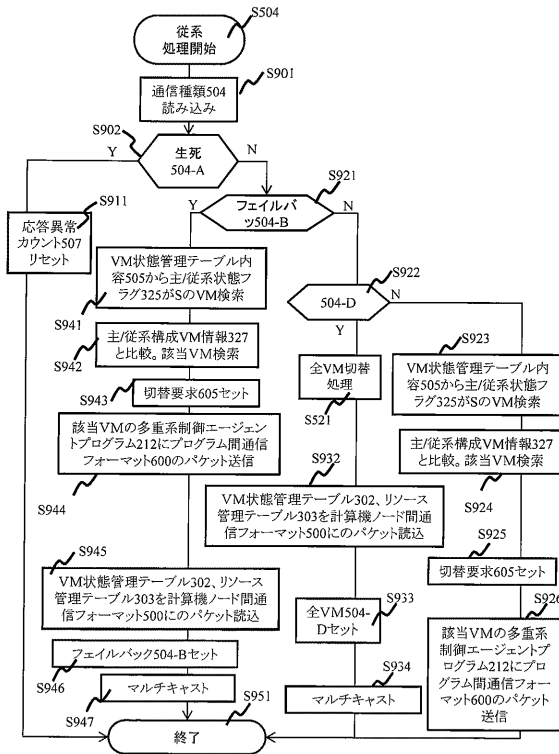
【 図 1 8 】

図18



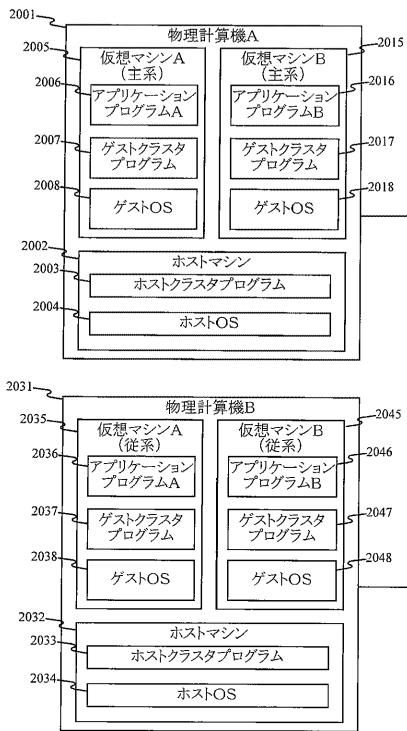
【 図 1 9 】

図19



【 図 20 】

図20



【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2010/006654

A. CLASSIFICATION OF SUBJECT MATTER G06F11/20(2006.01)i, G06F9/46(2006.01)i		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) G06F11/16-11/20, G06F9/46-9/54		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Jitsuyo Shinan Koho 1922-1996 Jitsuyo Shinan Toroku Koho 1996-2011 Kokai Jitsuyo Shinan Koho 1971-2011 Toroku Jitsuyo Shinan Koho 1994-2011		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y A	JP 2008-269332 A (Hitachi, Ltd.), 06 November 2008 (06.11.2008), paragraphs [0023] to [0036]; fig. 1, 2 & US 2008/0263390 A1	1, 5-8, 12-14 2-4, 9-11
Y A	JP 2009-211517 A (NEC Corp.), 17 September 2009 (17.09.2009), paragraphs [0016] to [0021], [0026] to [0031]; fig. 1, 2 (Family: none)	1, 5, 6, 8, 12, 13 2-4, 9-11
Y A	JP 2005-327279 A (International Business Machines Corp.), 24 November 2005 (24.11.2005), paragraph [0036] & US 2005/0268298 A1 & CN 1696902 A	1, 5, 7, 8, 12, 14 2-4, 9-11
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed		"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family
Date of the actual completion of the international search 03 February, 2011 (03.02.11)		Date of mailing of the international search report 15 February, 2011 (15.02.11)
Name and mailing address of the ISA/ Japanese Patent Office		Authorized officer
Facsimile No.		Telephone No.

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2010/006654

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	JP 2009-187090 A (NEC Corp.), 20 August 2009 (20.08.2009), entire text; all drawings (Family: none)	1-14

特許協力条約

PCT

国際調査報告

(法8条、法施行規則第40、41条)
〔PCT18条、PCT規則43、44〕

出願人又は代理人 の書類記号 341050897PCT	今後の手続きについては、様式PCT/ISA/220 及び下記5を参照すること。	
国際出願番号 PCT/JP2010/006654	国際出願日 (日.月.年) 12. 11. 2010	優先日 (日.月.年)
出願人 (氏名又は名称) 株式会社日立製作所		

国際調査機関が作成したこの国際調査報告を法施行規則第41条（PCT18条）の規定に従い出願人に送付する。
この写しは国際事務局にも送付される。

この国際調査報告は、全部で 3 ページである。

この調査報告に引用された先行技術文献の写しも添付されている。

1. 国際調査報告の基礎

a. 言語に関し、この国際調査は以下のものに基づき行った。

出願時の言語による国際出願

出願時の言語から国際調査のための言語である _____ 語に翻訳された、
この国際出願の翻訳文（PCT規則12.3(a)及び23.1(b)）

b. この国際調査報告は、PCT規則91の規定により国際調査機関が認めた又は国際調査機関に通知された明らかな誤りの訂正を考慮して作成した（PCT規則43.6の2(a)）。

c. この国際出願は、ヌクレオチド又はアミノ酸配列を含んでいる（第I欄参照）。

2. 請求の範囲の一部の調査ができない（第II欄参照）。

3. 発明の単一性が欠如している（第III欄参照）。

4. 発明の名称は 出願人が提出したものを承認する。

次に示すように国際調査機関が作成した。

5. 要約は 出願人が提出したものを承認する。

第IV欄に示されているように、法施行規則第47条第1項（PCT規則38.2）の規定により国際調査機関が作成した。出願人は、この国際調査報告の発送の日から1月以内にこの国際調査機関に意見を提出することができる。

6. 図面に関して

a. 要約書とともに公表される図は、

第 1 図とする。 出願人が示したとおりである。

出願人は図を示さなかったため、国際調査機関が選択した。

本図は発明の特徴を一層よく表しているため、国際調査機関が選択した。

b. 要約とともに公表される図はない。

様式PCT/ISA/210（第1ページ）（2009年7月）（改訂）

国際調査報告		国際出願番号 PCT/JP2010/006654									
A. 発明の属する分野の分類 (国際特許分類 (IPC)) Int.Cl. G06F11/20(2006.01)i, G06F9/46(2006.01)i											
B. 調査を行った分野 調査を行った最小限資料 (国際特許分類 (IPC)) Int.Cl. G06F11/16-11/20, G06F9/46-9/54											
最小限資料以外の資料で調査を行った分野に含まれるもの <table border="0"> <tr> <td>日本国実用新案公報</td> <td>1922-1996年</td> </tr> <tr> <td>日本国公開実用新案公報</td> <td>1971-2011年</td> </tr> <tr> <td>日本国実用新案登録公報</td> <td>1996-2011年</td> </tr> <tr> <td>日本国登録実用新案公報</td> <td>1994-2011年</td> </tr> </table>				日本国実用新案公報	1922-1996年	日本国公開実用新案公報	1971-2011年	日本国実用新案登録公報	1996-2011年	日本国登録実用新案公報	1994-2011年
日本国実用新案公報	1922-1996年										
日本国公開実用新案公報	1971-2011年										
日本国実用新案登録公報	1996-2011年										
日本国登録実用新案公報	1994-2011年										
国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)											
C. 関連すると認められる文献											
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号									
Y A	JP 2008-269332 A (株式会社日立製作所) 2008.11.06, 段落[0023]-[0036], 第1,2図 & US 2008/0263390 A1	1, 5-8, 12-14 2-4, 9-11									
Y A	JP 2009-211517 A (日本電気株式会社) 2009.09.17, 段落[0016]-[0021], [0026]-[0031], 第1,2図 (ファミリーなし)	1, 5, 6, 8, 12, 13 2-4, 9-11									
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <input type="checkbox"/> パテントファミリーに関する別紙を参照。											
* 引用文献のカテゴリー		の日の後に公表された文献									
「A」特に関連のある文献ではなく、一般的技術水準を示すもの		「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの									
「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの		「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの									
「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)		「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの									
「O」口頭による開示、使用、展示等に言及する文献		「&」同一パテントファミリー文献									
「P」国際出願日前で、かつ優先権の主張の基礎となる出願											
国際調査を完了した日 03.02.2011		国際調査報告の発送日 15.02.2011									
国際調査機関の名称及びあて先 日本国特許庁 (ISA/JP) 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号		特許庁審査官 (権限のある職員) 北元 健太	5B 3856								
		電話番号 03-3581-1101	内線 3545								

国際調査報告		国際出願番号 PCT/JP2010/006654
C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
Y A	JP 2005-327279 A (インターナショナル・ビジネス・マシーンズ・ コーポレーション) 2005.11.24, 段落[0036] & US 2005/0268298 A1 & CN 1696902 A	1, 5, 7, 8, 12, 14 2-4, 9-11
A	JP 2009-187090 A (日本電気株式会社) 2009.08.20, 全文, 全図 (ファミリーなし)	1-14

(注)この公表は、国際事務局(WIPO)により国際公開された公報を基に作成したものである。なおこの公表に係る日本語特許出願(日本語実用新案登録出願)の国際公開の効果は、特許法第184条の10第1項(実用新案法第48条の13第2項)により生ずるものであり、本掲載とは関係ありません。