

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5585844号
(P5585844)

(45) 発行日 平成26年9月10日(2014.9.10)

(24) 登録日 平成26年8月1日(2014.8.1)

(51) Int.Cl.		F I			
G06F	13/10	(2006.01)	G06F	13/10	330C
G06F	9/46	(2006.01)	G06F	9/46	350
G06F	9/50	(2006.01)	G06F	9/46	462Z

請求項の数 14 (全 26 頁)

(21) 出願番号	特願2011-67444 (P2011-67444)	(73) 特許権者	000005108
(22) 出願日	平成23年3月25日 (2011.3.25)		株式会社日立製作所
(65) 公開番号	特開2012-203636 (P2012-203636A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成24年10月22日 (2012.10.22)	(74) 代理人	100114236
審査請求日	平成25年6月3日 (2013.6.3)		弁理士 藤井 正弘
		(74) 代理人	100075513
			弁理士 後藤 政喜
		(74) 代理人	100120260
			弁理士 飯田 雅昭
		(72) 発明者	服部 直也
			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	澤 勇太
			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
			最終頁に続く

(54) 【発明の名称】 仮想計算機の制御方法及び計算機

(57) 【特許請求の範囲】

【請求項1】

プロセッサとメモリと仮想I/Oアダプタを仮想機能として生成する物理機能と前記仮想機能と前記物理機能のそれぞれの状態を管理するレジスタとを有するI/Oアダプタを備えた物理計算機と、前記物理計算機の計算機資源を仮想計算機に提供する仮想化部と、前記仮想計算機で実行されるOSと、を備えて、前記仮想化部が前記仮想機能を割り当てた1つ以上の仮想計算機を生成し、前記仮想計算機でOSを稼働させ、前記OSに前記仮想機能を利用させる仮想計算機の制御方法であって、

前記仮想化部が、前記I/Oアダプタのレジスタを参照し前記I/Oアダプタの障害の発生を検出する第1のステップと、

前記仮想化部が、前記障害の発生した前記I/Oアダプタの前記物理機能により提供される前記仮想機能が割り当てられた仮想計算機を特定する第2のステップと、

前記仮想化部が、前記特定された仮想計算機上で稼働する前記OSに前記仮想機能の利用中止を指示する第3のステップと、

前記仮想化部が、前記OSから前記利用中止の指示に対する応答を受領した後、前記I/Oアダプタの前記物理機能の回復処理を実施する第4のステップと、

を含むことを特徴とする仮想計算機の制御方法。

【請求項2】

請求項1に記載の仮想計算機の制御方法であって、

前記I/Oアダプタは、

障害が発生したときに、障害情報を格納する障害情報レジスタを有し、

前記第1のステップは、

前記仮想化部が、前記障害を示す信号を受信したときに前記障害情報レジスタを参照することを特徴とする仮想計算機の制御方法。

【請求項3】

請求項2に記載の仮想計算機の制御方法であって、

前記仮想計算機は、

前記I/Oアダプタの挿抜要求と通知を格納するホットプラグレジスタをさらに有し、

前記第2のステップは、

全ての前記仮想機能について前記特定を行い、

前記第4のステップは、

前記仮想化部が、前記ホットプラグレジスタに前記仮想機能の取り外し要求を設定する第5のステップと、

前記OSが、前記ホットプラグレジスタを読み込む第6のステップと、

前記特定された仮想計算機上で稼働する全ての前記OSが、前記ホットプラグレジスタに前記仮想機能の取り外し要求に対して応答するまで待機する第7のステップと、

をさらに含むことを特徴とする仮想計算機の制御方法。

10

【請求項4】

請求項3に記載の仮想計算機の制御方法であって、

前記仮想化部は、前記OSが前記仮想機能の利用を中止した後に、前記I/Oアダプタを再初期化する第8のステップと、

前記I/Oアダプタが、前記仮想機能を再度生成し、前記仮想化部が、再度生成された仮想機能を前記仮想計算機に再度割り当てる第9のステップと、

前記仮想化部が、前記ホットプラグレジスタにI/Oアダプタの取り付け通知を設定する第10のステップと、

前記OSが、前記ホットプラグレジスタを読み込んで前記I/Oアダプタの取り付けを検知する第11のステップと、

前記OSが、前記仮想機能の利用を再開する第12のステップと、

をさらに含むことを特徴とする仮想計算機の制御方法。

20

【請求項5】

請求項1に記載の仮想計算機の制御方法であって、

前記仮想化部は、

前記I/Oアダプタを操作するドライバを有し、当該ドライバには予め障害発生時の処理の方針が設定され、

前記第4のステップは、

前記仮想化部が、前記障害発生時には前記ドライバに対して、障害対処の方針を問い合わせる第13のステップと、

前記ドライバに設定された前記方針がリセット試行である場合に、前記仮想計算機で稼働する前記OSに前記仮想機能の状態を伝達する第14のステップと、

をさらに含むことを特徴とする仮想計算機の制御方法。

30

40

【請求項6】

請求項3に記載の仮想計算機の制御方法であって、

前記仮想化部は、

前記I/Oアダプタを操作するドライバを有し、当該ドライバには予め障害発生時の処理の方針が設定され、

前記仮想化部は、前記OSが前記仮想機能の利用を中止した後に、前記I/Oアダプタを再初期化する第8のステップと、

前記I/Oアダプタが、前記仮想機能を再度生成し、前記仮想化部が、再度生成された仮想機能を前記仮想計算機に再度割り当てる第9のステップと、

前記仮想化部が、前記ホットプラグレジスタにI/Oアダプタの取り付け通知を設定す

50

る第 10 のステップと、

前記 OS が、前記ホットプラグレジスタを読み込んで前記 I/Oアダプタの取り付けを検知する第 11 のステップと、

前記 OS が、前記仮想機能の利用を再開する第 12 のステップと、
をさらに含み、

前記第 4 のステップは、

前記仮想化部が、前記障害発生時には前記ドライバに対して、障害対処の方針を問い合わせる第 13 のステップと、

前記ドライバに設定された前記方針がリセット試行である場合に、前記仮想計算機で稼働する前記 OS に前記仮想機能の状態を伝達する第 14 のステップと、

をさらに含むことを特徴とする仮想計算機の制御方法。

10

【請求項 7】

請求項 1 に記載の仮想計算機の制御方法であって、

前記仮想計算機は、第 1 の OS が稼働する第 1 の仮想計算機と第 2 の OS が稼働する第 2 の仮想計算機を含み、

前記仮想化部は、前記第 1 の仮想計算機及び第 2 の仮想計算機に前記仮想機能をそれぞれ割り当てて、前記第 1 の OS 及び前記第 2 の OS に前記割り当てた前記仮想機能をそれぞれ利用させ、

前記第 4 のステップは、

前記仮想化部が、前記第 1 の OS からの前記仮想機能の利用中止の指示に対する応答と、前記第 2 の OS からの前記仮想機能の利用中止の指示に対する応答とを待ち合わせることを特徴とする仮想計算機の制御方法。

20

【請求項 8】

プロセッサとメモリと仮想 I/Oアダプタを仮想機能として生成する物理機能と前記仮想機能と前記物理機能のそれぞれの状態を管理するレジスタとを有する I/Oアダプタを備えた物理計算機と、

前記物理計算機の計算機資源を仮想計算機に提供する仮想化部と、

前記仮想計算機で実行される OS と、を備えて、前記仮想化部が、前記仮想機能を割り当てた 1 つ以上の仮想計算機を生成し、前記仮想計算機で前記 OS を稼働させ、前記 OS に前記仮想機能を利用させる計算機であって、

30

前記仮想化部は、

前記 I/Oアダプタのレジスタを参照し前記 I/Oアダプタの障害の発生を検出する障害処理部と、

前記物理機能が生成した前記仮想機能について、当該仮想機能を割り当てた前記仮想計算機を管理する割り当て情報と、

を備え、

前記障害処理部は、

前記障害の発生した前記 I/Oアダプタの前記割り当て情報を参照し、前記物理機能により提供される前記仮想機能が割り当てられた前記仮想計算機を特定し、前記特定された仮想計算機上で稼働する前記 OS に前記仮想機能の利用中止を指示し、前記 OS から前記利用中止の指示に対する応答を受領した後、前記 I/Oアダプタの前記物理機能の回復処理を実施することを特徴とする計算機。

40

【請求項 9】

請求項 8 に記載の計算機であって、

前記 I/Oアダプタは、

障害情報を格納する障害情報レジスタを有し、

前記障害処理部は、

前記 I/Oアダプタで発生した障害を示す信号を受信したときに前記障害情報レジスタを参照することを特徴とする計算機。

【請求項 10】

50

請求項 9 に記載の計算機であって、
 前記仮想計算機は、
 前記 I / O アダプタの挿抜要求と通知を格納するホットプラグレジスタをさらに有し、
 前記障害処理部は、全ての前記仮想機能について前記特定を行い、
 前記障害処理部は、
 前記ホットプラグレジスタに前記仮想機能の取り外し要求を設定し、前記特定された仮想計算機上で稼働する全ての前記 OS が、前記ホットプラグレジスタに前記仮想機能の取り外し要求に対して応答するまで待機することを特徴とする計算機。

【請求項 11】

請求項 10 に記載の計算機であって、
 前記障害処理部は、
 前記 OS が前記仮想機能の利用を中止した後に、前記 I / O アダプタを再初期化し、前記ホットプラグレジスタに I / O アダプタの取り付け通知を設定し
 前記 I / O アダプタは、前記再初期化の後に前記仮想機能を再度生成し、
 前記仮想化部が、再度生成された前記仮想機能を前記仮想計算機に再度割り当てて、
 前記 OS が、前記ホットプラグレジスタを読み込んで前記 I / O アダプタの取り付けを検知した後に、前記仮想機能の利用を再開することを特徴とする計算機。

【請求項 12】

請求項 8 に記載の計算機であって、
 前記仮想化部は、
 前記 I / O アダプタを操作するドライバを有し、当該ドライバには予め障害発生時の処理の方針が設定され、
 前記障害処理部が、前記障害発生時には前記ドライバに対して、障害対処の方針を問い合わせ、前記ドライバに設定された前記方針がリセット試行である場合に、前記仮想計算機で稼働する前記 OS に前記仮想機能の状態を伝達することを特徴とする計算機。

【請求項 13】

請求項 10 に記載の計算機であって、
 前記仮想化部は、
 前記 I / O アダプタを操作するドライバを有し、当該ドライバには予め障害発生時の処理の方針が設定され、
 前記障害処理部は、前記 OS が前記仮想機能の利用を中止した後に、前記 I / O アダプタを再初期化し、
 前記 I / O アダプタが、前記仮想機能を再度生成し、
 前記仮想化部が、前記再度生成した仮想機能を前記仮想計算機に再度割り当てて、
 前記障害処理部が、前記ホットプラグレジスタに I / O アダプタの取り付け通知を設定し、

前記 OS が、前記ホットプラグレジスタを読み込んで前記 I / O アダプタの取り付けを検知して、前記仮想機能の利用を再開し、
 前記障害処理部が、前記障害発生時には前記ドライバに対して、障害対処の方針を問い合わせ、前記ドライバに設定された前記方針がリセット試行である場合に、前記仮想計算機で稼働する前記 OS に前記仮想機能の状態を伝達することを特徴とする計算機。

【請求項 14】

請求項 8 に記載の計算機であって、
前記仮想計算機は、第 1 の OS が稼働する第 1 の仮想計算機と第 2 の OS が稼働する第 2 の仮想計算機を含み、
前記仮想化部は、前記第 1 の仮想計算機及び第 2 の仮想計算機に前記仮想機能をそれぞれ割り当てて、前記第 1 の OS 及び前記第 2 の OS に前記割り当てた前記仮想機能をそれぞれ利用させ、
前記障害処理部は、
前記第 1 の OS からの前記仮想機能の利用中止の指示に対する応答と、前記第 2 の OS

からの前記仮想機能の利用中止の指示に対する応答とを待ち合わせることを特徴とする計算機。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、仮想計算機に関し、仮想計算機上で稼動するゲストOSを停止させずにI/Oデバイスの障害回復を処理する技術に関する。

【背景技術】

【0002】

オープン系サーバの性能向上および機能拡充に伴い、物理サーバに搭載されるプロセッサコアを有効に活用する方法として、サーバ仮想化ソフトウェア（ハイパバイザやVMM）が広く用いられている。ハイパバイザは、1台の物理サーバ上で複数の仮想計算機を生成し、それぞれの仮想計算機でOSやアプリケーションを稼働させる。プロセッサの性能向上と相まって、近年では物理サーバ上で10以上の仮想計算機（LPAR）を稼働させるケースも珍しくない。しかし、物理サーバ上で稼動する仮想計算機の数が増えるにつれ、I/Oデバイスに関する2つの問題が浮上していた。

【0003】

問題1（性能問題）：物理I/Oアダプタを、複数の仮想I/Oアダプタとして振る舞わせる「I/O共有」を実現するためにハイパバイザの介入が不可欠だった。この場合、ハイパバイザの介入オーバーヘッドが原因で、仮想計算機の利用できるI/O性能が限られていた。

【0004】

問題2（信頼性問題）：オープン系サーバでは従来、I/Oアダプタ障害の内容をOS等のソフトウェアに伝達する仕組みが存在しなかった。そのため、物理I/Oアダプタに障害が発生すると、障害の規模や種類が判断できないために必ず物理サーバ全体がダウンし、物理サーバ上で稼働する全ての仮想計算機に被害が及んでいた。

【0005】

これらの問題を解決するためにPCI Sigから、IOV(I/O Virtualization)とAER(Advanced Error Reporting)という2つの規格が制定された。

【0006】

IOVは、上記の性能問題を解決するために、「I/O共有」の主要部をハードウェアで提供する仕組みである。IOVを用いると、ハイパバイザの介入が初期化等の低頻度処理のみとなるため、仮想計算機で高いI/O性能を利用可能になる。

【0007】

一方、AERは、非特許文献1で示すように、上記信頼性問題を改善するために、I/Oアダプタの障害情報をOS等のソフトウェアに伝達する仕組みである。AERによって、OS等のソフトウェアは障害の深刻度を判断可能となり、軽度な障害であればリセット等の方法でI/Oアダプタを回復させ、物理サーバや仮想計算機の稼働を継続可能となる。

【0008】

なお、IOVを利用する技術として、特許文献1にハイパバイザの基本的な動作が開示されている。特許文献1ではIOVに対応するI/Oアダプタの物理機能（以下、PF: Physical Function）の障害に対処する方法として、PFのHot Plug（サーバを稼働させた状態でPFを抜き差しする技術）を利用する技術が含まれている。特許文献1を使用すれば、PFに障害が発生した場合には、保守要員や管理者等がI/Oアダプタを交換することができる。

【先行技術文献】

【特許文献】

【0009】

【特許文献1】米国特許出願公開第2009/133016号明細書

10

20

30

40

50

【非特許文献】

【0010】

【非特許文献1】PCI Express(r) Base Specification Revision 2.1 (§7.10. Advanced Error Reporting Capability)

【発明の概要】

【発明が解決しようとする課題】

【0011】

しかしながら、上記特許文献1には軽度な障害への対処方法が言及されていない。I/Oアダプタの電源投入直後、及びリセット直後は、従来との互換性の観点からIOV機能が無効化される。このため、IOVに対応したI/Oアダプタは、Hot Plugの直後ではI/O共有の為の仮想I/Oアダプタ機能(以下、VF: Virtual Functionとする)が無効化されているため、特許文献1の技術を用いると以下の問題が生じる。

10

【0012】

軽度な障害に対して、AERの仕組みに則ってI/Oアダプタ(PF)をリセットすると、IOV機能が無効化されるため、仮想計算機からIOVによって生成されていた仮想I/Oアダプタ機能(VF)が消失する。仮想計算機上で動くOS(以下、ゲストOS)は、一般に物理I/Oアダプタ(PF)向けに設計されているため、I/Oアダプタが突然消失する事態を想定していない。そのためVFが突然消失するとゲストOSがクラッシュするケースがある。また、上記特許文献1では、軽度な障害に対してもPFをリセットしなければならず、障害から回復するために保守員による交換作業が発生し、保守の手間がかかる、という問題があった。

20

【0013】

そこで本発明は、上記問題点に鑑みてなされたもので、IOVに対応するI/Oアダプタの高性能、高信頼性、保守容易性を充足することを目的とする。特に、I/OアダプタのPFリセットが引き起こすVFの無断消失に起因するゲストOSのクラッシュを回避し、I/Oアダプタの軽度な障害から自動的に回復させることを目的とする。

【課題を解決するための手段】

【0014】

本発明は、プロセッサとメモリと仮想I/Oアダプタを仮想機能として生成する物理機能と前記仮想機能と物理機能のそれぞれの状態を管理するレジスタとを有するI/Oアダプタを備えた物理計算機と、前記物理計算機の計算機資源を仮想計算機に提供する仮想化部と、前記仮想計算機で実行されるOSと、を備えて、前記仮想化部が前記仮想機能を割り当てた1つ以上の仮想計算機を生成し、前記仮想計算機でOSを稼働させ、前記OSに前記仮想機能を利用させる仮想計算機の制御方法であって、前記仮想化部が、前記I/Oアダプタのレジスタを参照し前記I/Oアダプタの障害の発生を検出する第1のステップと、前記仮想化部が、前記障害の発生した前記I/Oアダプタの前記物理機能により提供される前記仮想機能が割り当てられた仮想計算機を特定する第2のステップと、前記仮想化部が、前記特定された仮想計算機上で稼働する前記OSに前記仮想機能の利用中止を指示する第3のステップと、前記仮想化部が、前記OSから前記利用中止の指示に対する応答を受領した後、前記I/Oアダプタの前記物理機能の回復処理を実施する第4のステップと、を含む。

30

40

【発明の効果】

【0015】

したがって、本発明は、I/O共有機能を有する物理I/Oアダプタを備えて、仮想化部は、仮想I/Oアダプタ(VF)を仮想計算機に割り当てる。障害発生時などの特定のケースを除けば仮想化部の動作頻度が低いため、高性能を実現できる。次に、VF消失の契機となる物理I/Oアダプタの機能(PF)のリセットに先だって、仮想計算機のOSからVFの取り外し許可を受ける。これにより、VF消失によるOSクラッシュを回避できる。従って高信頼である。PFドライバが軽度な障害と判断した場合に、保守要員の手

50

を借りずに P F をリセットし、物理 I / O アダプタの回復を試みる。物理 I / O アダプタの回復に成功した場合は仮想計算機に V F を割り当て直し、仮想計算機を障害発生前の状態に戻す。従って保守が容易である。以上より、高性能、高信頼、保守容易性が充足される。

【図面の簡単な説明】

【 0 0 1 6 】

【図 1】本発明の第 1 の実施形態を示し、計算機システムの一例を示すブロック図である。

【図 2】本発明の第 1 の実施形態を示し、物理計算機及びソフトウェアの関係を示すブロック図である。

【図 3】本発明の第 1 の実施形態を示し、物理計算機のメモリマップの一例である。

【図 4】本発明の第 1 の実施形態を示し、アダプタ割り当て表の一例である。

【図 5】本発明の第 1 の実施形態を示し、P F 状態表の一例である。

【図 6】本発明の第 1 の実施形態を示し、V F 状態表の一例である。

【図 7】本発明の第 1 の実施形態を示し、ハイパバイザの全体的な処理を説明するフローチャートの例である。

【図 8】本発明の第 1 の実施形態を示し、ハイパバイザの I / O 障害処理を説明するフローチャートの例である。

【図 9】本発明の第 1 の実施形態を示し、ハイパバイザの P F 障害処理を説明するフローチャートの例である。

【図 10】本発明の第 1 の実施形態を示し、ハイパバイザの V F 障害処理を説明するフローチャートの例である。

【図 11】本発明の第 1 の実施形態を示し、ハイパバイザの I / O 仮想化処理を説明するフローチャートの例である。

【図 12】本発明の第 1 の実施形態を示し、ハイパバイザの P F 回復処理を説明するフローチャートの例である。

【図 13】本発明の第 1 の実施形態を示し、ハイパバイザの V F 交換処理を説明するフローチャートの例である。

【図 14】本発明の第 2 の実施形態を示し、物理計算機及びソフトウェアの関係を示すブロック図である。

【図 15】本発明の第 2 の実施形態を示し、物理計算機のメモリマップの一例である。

【図 16】本発明の第 2 の実施形態を示し、ハイパバイザの P F 障害処理を説明するフローチャートの例である。

【図 17】本発明の第 2 の実施形態を示し、ハイパバイザの V F 障害処理を説明するフローチャートの例である。

【図 18】本発明の第 2 の実施形態を示し、ハイパバイザの I / O 仮想化処理を説明するフローチャートの例である。

【図 19】本発明の第 2 の実施形態を示し、ハイパバイザの P F 回復処理を説明するフローチャートの例である。

【図 20】本発明の第 2 の実施形態を示し、ハイパバイザの V F 回復処理を説明するフローチャートの例である。

【図 21】本発明の第 2 の実施形態を示し、ハイパバイザの V F 交換処理を説明するフローチャートの例である。

【発明を実施するための形態】

【 0 0 1 7 】

以下、本発明の一実施形態を添付図面に基づいて説明する。

【 0 0 1 8 】

図 1 は、本発明の第 1 の実施形態を示し、I O V (I/O Virtualization) と A E R (Advanced Error Reporting) に対応する I / O アダプタを用いた仮想計算機システムの一例を示すブロック図である。本実施形態の I / O アダプタ 6 0 は、物理的な I / O アダプ

10

20

30

40

50

タの物理機能（以下、PF：Physical Function）と、「I/O共有」を実現する仮想I/Oアダプタ（以下、VF：Virtual Function）を生成する機能を有し、VFを仮想計算機に割り当てる場合の構成及び動作を説明する。

【0019】

< 1. ハードウェア構成 >

図1は、本発明の第1の実施形態となる仮想計算機システムを動作させる物理計算機の構成例を示す。物理計算機は、障害検出機能を有するCPU70（またはプロセッサコア）を1つ以上有し、これらのCPU70はQPI（Quick Path Interconnect）やSMI（Scalable Memory Interconnect）等のインターコネクタ110を介してChipSet100やメモリ90に接続される。

10

【0020】

ChipSet100には、PCIexpress等のバス120を介してI/Oアダプタ60が接続される。I/Oアダプタ60はLAN130に接続されるネットワークアダプタ、ディスク装置140等に接続されるSCSIアダプタ、SAN150（Storage Area Network）に接続されるファイバーチャネルアダプタ、コンソール80に接続されるグラフィックコントローラなどで構成される。

【0021】

CPU70はインターコネクタ110を介してメモリ90にアクセスし、ChipSet100からI/Oアダプタ60にアクセスして所定の処理を行う。

【0022】

メモリ90には、仮想化部としてハイパバイザ20がロードされ、ハイパバイザ20が制御する仮想計算機30-1~30-nが、ゲストOS40-1~40-nを実行する。ハイパバイザ20は、物理計算機10の計算器資源を分割または仮想化して仮想計算機に割り当てる。なお、以下では仮想計算機30-1~3-nの総称を符号30で示し、ゲストOS40-1~40-nの総称を符号40で示す。

20

【0023】

< 2. ソフトウェア構成 >

次に、物理計算機10上で仮想計算機30を実現するソフトウェアの構成の主要部と、制御対象となるハードウェア要素について、図2を参照しながら詳述する。図2は、物理計算機及びソフトウェアの関係を示すブロック図である。

30

【0024】

物理計算機10上では、1つ以上の仮想計算機30を制御するハイパバイザ20が稼働する。物理計算機10は、ChipSet100と、1つ以上のI/Oアダプタ60を含む。

【0025】

ChipSet100は、I/Oアダプタ60の取り外し（挿抜）要求（Hot Removeの要求）を保持するレジスタと、取り付け告知（Hot Addの通知）を保持するレジスタから成るHotPlugレジスタ（以下、HPレジスタ）185を含む。またChipSet100は、I/Oアダプタ60の取り外し要求または取り付け告知あるいは障害を認識すると割り込みを発生して、物理計算機10上で稼働するハイパバイザ20に要求や障害の発生を通知する機能を有する。

40

【0026】

1つ以上のI/Oアダプタ60はIOV機能を有しており、物理的な機能（PF：Physical Function）160（160-1~160-m）と仮想的な機能（VF：Virtual Function）190（190-1-1~190-m-k）から成る。PF160はいつでも利用できる機能だが、VF190はIOV機能が有効な場合のみ利用できる機能である。なお、仮想計算機10は、IOVに非対応のI/Oアダプタ60を含んでも構わない。

【0027】

PF160は、物理計算機10の外部とデータを送受信する機能を有し、IOV機能を制御するIOVレジスタ170と、PFに関する障害の有無や種類といった障害情報を保

50

持するA E Rレジスタ1 8 0を含む。

【0 0 2 8】

V F 1 9 0は、I O V機能が有効な場合のみ、物理計算機1 0の外部とデータを送受信する機能を有し、V Fに関する障害の有無や種類といった障害情報を保持するA E Rレジスタ1 8 0を含む。

【0 0 2 9】

ここで、A E Rレジスタ1 8 0は、2種類存在し、P F 1 6 0のA E Rレジスタ1 8 0とV F 1 9 0のA E Rレジスタ1 8 0が存在する。P F 1 6 0のA E Rレジスタ1 8 0は物理機能、つまり物理I / Oアダプタ6 0の障害に関する情報を保持し、V F 1 9 0のA E Rレジスタ1 8 0は仮想機能、つまり仮想I / Oアダプタの障害に関する情報を保持する。なお、I / Oアダプタ6 0は、障害発生時にA E Rレジスタ1 8 0に障害情報を設定するコントローラ(図示省略)を有する。

10

【0 0 3 0】

ハイバイザ2 0は、仮想計算機3 0を生成し、仮想計算機3 0上に対してC h i p S e t 1 0 0に相当する機能(仮想C h i p S e t 3 0 0)を提供する。またハイバイザ2 0は、任意のV F 1 9 0を任意の仮想計算機3 0に占有的に割り当て、当該仮想計算機3 0上で稼働するゲストOS 4 0に当該V F 1 9 0の直接操作を許す機能(パススルー機能)を備える。このパススルー機能によって、複数の仮想計算機3 0でひとつのI / Oアダプタ6 0をV F 1 9 0から共有することができる。

【0 0 3 1】

20

ハイバイザ2 0は、仮想計算機3 0に対するV F 1 9 0の割り当て関係を保持するアダプタ割り当て表2 0 0と、物理計算機1 0または物理計算機資源を監視して障害発生時に呼び出される障害処理部2 1 0と、仮想計算機3 0の状態を保持する仮想計算機のエミュレーションデータと、P F 1 6 0の種類に応じた制御手順を備えたP Fドライバ2 5 0(2 5 0 - 1 ~ 2 5 0 - m)を含む。

【0 0 3 2】

障害処理部2 1 0は、I O V機能に対応したI / Oアダプタ6 0の障害を受け持つI O V障害処理部2 1 2を含み、I O V障害処理部2 1 2はP F 1 6 0の障害処理状況保持するP F状態表2 1 4と、V F 1 9 0の障害処理状況保持するV F状態表2 1 6を含む。

30

【0 0 3 3】

仮想計算機のエミュレーションデータ2 2 0は、仮想計算機3 0に提供する仮想C h i p S e t 3 0 0などの仮想的な部品またはデバイスに関する状態を保持するために、仮想C h i p S e tデータ2 2 2を含む。仮想C h i p S e tデータ2 2 2は、仮想C h i p S e t 3 0 0が保持すべきH Pレジスタ1 8 5などの状態を保持する。

【0 0 3 4】

P Fドライバ2 5 0は、I O Vレジスタ1 7 0を操作する機能と、P F 1 6 0で障害が発生した場合に軽度障害か重度障害かを判断する機能、及びP Fリセット後に障害回復の成否を判断する機能を有する。

【0 0 3 5】

40

仮想計算機3 0は、ハイバイザ2 0によって提供される仮想C h i p S e t 3 0 0等の仮想的な部品と、占有的に割り当てられたV F 1 9 0を含む。仮想計算機3 0の上ではゲストOS 4 0が動作する。ゲストOS 4 0は、V F 1 9 0の種類に応じたV Fドライバ2 6 0を用いてV F 1 9 0を操作する。

【0 0 3 6】

V Fドライバ2 6 0は、V F 1 9 0で障害が発生した場合に軽度障害か重度障害かを判断する障害の程度を判定する機能、及びV Fリセット後に障害回復の成否を判定する機能を有する。

【0 0 3 7】

図3はハイバイザ2 0が管理するメモリ9 0のアドレスマップの一例を示す。

50

【 0 0 3 8 】

ハイパバイザ 2 0 は、メモリ 9 0 上に自身を配置する領域と、仮想計算機 3 0 が使用する領域を割り当てる。例えば図 3 のように、ハイパバイザ 2 0 は、自身にアドレス A D 0 ~ A D 1 を割り当て、仮想計算機 3 0 - 1 にアドレス A D 1 + 1 ~ A D 2 を、仮想計算機 3 0 - n にアドレス A D 3 ~ A D 4 を割り当てる。各仮想計算機 3 0 の領域には、ゲスト O S 4 0 と V F ドライバ 2 6 0 (2 6 0 - 1 ~ 2 6 0 - n) が格納される。

【 0 0 3 9 】

ハイパバイザ 2 0 が使用する領域には、アダプタ割り当て表 2 0 0、障害処理部 2 1 0、仮想計算機のエミュレーションデータ 2 2 0 と、P F ドライバ 2 5 0 が格納される。

【 0 0 4 0 】

図 4 は、アダプタ割り当て表 2 0 0 の構成例である。アダプタ割り当て表 2 0 0 は、V F 1 9 0 と仮想計算機 3 0 の割り当て関係を保持する表である。本表では V F 1 9 0 を識別する手段として、V F 1 9 0 を生成した P F 1 6 0 の通し番号 (P F #) 4 0 0 と、V F 1 9 0 の通し番号 (V F #) 4 1 0 を用い、当該 V F 1 9 0 を割り当てた仮想計算機 3 0 の通し番号 (仮想計算機番号) 4 2 0 を格納する。未割り当ての場合は仮想計算機番号 4 2 0 を "未割り当て" とする。

【 0 0 4 1 】

図 5 は、P F 状態表 2 1 4 の構成例である。P F 状態表 2 1 4 は、P F 1 6 0 の通し番号 (P F #) 4 0 0 と、P F 1 6 0 の障害処理状況として、P F 1 6 0 の状態 5 1 0 と、当該 P F 1 6 0 から生成された各 V F 1 9 0 の取り外し許可の有無を格納する待ち合わせ b i t m a p 5 2 0 からひとつのエントリが構成される。P F 1 6 0 の識別には P F # 4 0 0 を使用する。状態 5 1 0 は、障害処理が不要な状態 (正常)、軽度障害に伴って P F 1 6 0 のリセットが必要となった状態 (リセット保留)、重度障害に伴って P F 1 6 0 の交換が必要となった状態 (交換保留) のいずれかを保持する。待ち合わせ b i t m a p 5 2 0 は、各 b i t が当該 P F 1 6 0 の生成した各 V F 1 9 0 に対応しており、V F 1 9 0 の取り外し許可をゲスト O S 4 0 に要求中の場合に「 1」、未要求または許可済みの場合に「 0」が格納される。待ち合わせ b i t m a p 5 2 0 の全 b i t が 0 になると P F のリセットが可能になる。

【 0 0 4 2 】

図 6 は、V F 状態表 2 1 6 の構成例である。V F 状態表 2 1 6 は、P F 1 6 0 の通し番号 (P F #) 4 0 0 と、V F 1 9 0 の通し番号 (V F #) 4 1 0 と、V F 1 9 0 の障害処理状況として、V F 1 9 0 の状態 6 2 0 を保持する。V F 1 9 0 の識別には、P F # 4 0 0 と、V F # 4 1 0 を使用する。状態 6 2 0 は、障害処理が不要な状態 (正常)、軽度障害に伴って V F 1 9 0 のリセットが必要となった状態 (リセット保留)、重度障害に伴って V F 1 9 0 の交換が必要となった状態 (交換保留) のいずれかを保持する。

【 0 0 4 3 】

< 3 . ハイパバイザによる処理 >

次に、ハイパバイザ 2 0 が行う処理の一例について、以下、フローチャートを参照しながら説明する。

【 0 0 4 4 】

< 3 . 1 . ハイパバイザによる処理の概要 >

図 7 は、ハイパバイザ 2 0 が行う処理の全体像を示すフローチャートの例である。物理計算機 1 0 の電源を投入するとハイパバイザ 2 0 がメモリ 9 0 にロードされ、ハイパバイザ 2 0 自身と物理計算機 1 0 を初期化し、I / O アダプタ 6 0 の I O V 機能を有効化する (7 0 0)。ステップ 7 0 0 では、P F 状態表 2 1 4 は全 P F が正常状態に初期化され、V F 状態表 2 1 6 は全 V F が正常状態に初期化され、アダプタ割り当て表 2 0 0 は全 V F が未割り当てに初期化される。

【 0 0 4 5 】

続いてハイパバイザ 2 0 は、コンソール 8 0 からの入力、或いは前回起動時の割り当て指示に基づいて、V F 1 9 0 を割り当てた仮想計算機 3 0 を生成する (7 1 0)。ステッ

10

20

30

40

50

プ710では、アダプタ割り当て表200の、VF190の行に割り当てられた仮想計算機30の通し番号(識別子)を登録する。

【0046】

続いてハイバイザ20は、仮想計算機30の稼働を開始し、仮想計算機30上でゲストOS40及びアプリケーション(以下、両者を総称してゲストと記載)を実行する(720)。ステップ720ではゲストが任意の命令コードを実行できるが、1命令を実行する毎にCPU70がハイバイザ20の介入が必要なイベントが発生しているかをチェックする(730)。ハイバイザ20の介入が必要なイベントが発生していれば、CPU70がハイバイザ20を呼び出してステップ740に進む。ハイバイザ20の介入が必要なイベントが発生していなければ、ゲストの次の命令を実行するためにステップ720に進む。

10

【0047】

ハイバイザ20が呼び出されると、CPU70から発生したイベントの情報を受け取って解析し、ハイバイザ20はイベントの内容を取得する(740)。

【0048】

続いてハイバイザ20は、I/Oアダプタ60に障害が発生したか否かを判定し(750)、I/Oアダプタ60に障害が発生していればステップ790に進んでI/O障害処理を行う。本実施形態のI/Oアダプタ60は、障害発生時に割込信号を発生する。ハイバイザ20は、この割込信号を障害発生イベントとして受信する。

20

一方、ハイバイザ20は、I/Oアダプタ60に障害が発生していなければ、ステップ760に進んでゲストが仮想チップセット300等に対してI/O操作を行ったか否かを判定する。ゲストがI/O操作を行っていれば、ステップ780に進んでI/O仮想化処理を行う。ゲストがI/O操作を行っていなければステップ770に進み、ハイバイザ20の従来の処理(解析したイベントに応じた処理)を行う。ハイバイザ20の従来の処理については、公知または周知の処理を実施すればよいので、説明を省略する。

【0049】

ステップ770、780、790の何れかで発生したイベントに対する処理が完了すると、ハイバイザ20はゲストの実行を再開させる(795)。

【0050】

<3.2.I/O障害処理>

図8は、図7のステップ790で行われるI/O障害処理のフローチャートの例である。I/Oアダプタ60に障害の発生を検出したハイバイザ20は、PF160又はVF190が備えるAERレジスタ180を参照し、どのI/Oアダプタ60で障害が発生したのかを特定する(800)。

30

【0051】

続いてハイバイザ20は、障害が発生したI/Oアダプタ60の種類によって処理を分岐させる(810)。ハイバイザ20は発生した障害が、IOVに非対応のI/Oアダプタ(従来アダプタ)60で発生した場合は、ステップ820に進んで従来の障害処理を行う。従来の障害処理については、公知または周知の処理を実施すればよいので、説明を省略する。

40

【0052】

ハイバイザ20は発生した障害が、IOV機能を有するI/OアダプタのPF160部分で発生した場合は、ステップ830に進んでPF障害処理を後述のように行う。一方、ハイバイザ20は発生した障害が、IOV機能を有するI/OアダプタのVF190部分で発生した場合は、ステップ840に進んでVF障害処理を後述のように行う。

【0053】

<3.2.1.PF障害処理>

図9は、図8のステップ830で行われるPF障害処理のフローチャートの例である。障害が発生したPF160を特定したハイバイザ20は、当該PF160の種類に応じ

50

た P F ドライバ 2 5 0 に、障害からの回復が可能か否かを問い合わせ、障害の重度を取得する (9 0 0)。 I / O アダプタ 6 0 の障害の程度には 3 通りの重度があり、重い障害から順に、

(1) 回復不能で I / O アダプタ 6 0 の交換が必要となる重度障害 (P C I _ E R S _ R E S U L T _ D I S C O N N E C T)、

(2) P F 1 6 0 をリセットすれば回復できる可能性のある軽度障害 (P C I _ E R S _ R E S U L T _ N E E D _ R E S E T)、

(3) P F 1 6 0 のレジスタの一部を再設定する等の処理で回復が完了する障害 (P C I _ E R S _ R E S U L T _ R E C O V E R E D)、

に分類される。

10

【 0 0 5 4 】

ハイパバイザ 2 0 は、障害の重度に応じて処理を分岐させる (9 1 0)。

【 0 0 5 5 】

(1) 回復不能な重度障害であれば、 P F 状態表 2 1 4 の該当する P F 1 6 0 のエントリの状態 5 1 0 を、「交換保留」に変更する (9 3 0)。

【 0 0 5 6 】

(2) リセットで回復できる可能性がある軽度障害であれば、ハイパバイザ 2 0 は P F 状態表 2 1 4 の該当する P F 1 6 0 のエントリの状態 5 1 0 を、「リセット保留」に変更する (9 6 0)。

【 0 0 5 7 】

(3) 回復完了した障害であれば、ハイパバイザ 2 0 の回復処理は省略する。

20

【 0 0 5 8 】

上記 (1) 及び (2) のケースでは更に、ハイパバイザ 2 0 がアダプタ割り当て表 2 0 0 を参照し、障害が起きた P F 1 6 0 から生成された全ての V F 1 9 0 と、各 V F 1 9 0 を割り当てた仮想計算機 3 0 を特定する (9 4 0)。

【 0 0 5 9 】

続いてハイパバイザ 2 0 は、各仮想計算機 3 0 に含まれる仮想 C h i p S e t 3 0 0 の H P レジスタ 1 8 5 に V F 1 9 0 の稼働時取り外し (H o t R e m o v e) 要求を設定する。そして、ハイパバイザ 2 0 は、仮想計算機 3 0 に H o t P l u g 割り込みを通知して、ゲスト O S 4 0 に V F 1 9 0 の利用中止と安全な取り外し (S a f e l y R e m o v e) を指令する (9 5 0)。

30

【 0 0 6 0 】

最後にハイパバイザ 2 0 は、障害内容と対処結果をコンソール 8 0 に出力する (9 2 0)。コンソール 8 0 への出力により保守要員に I / O アダプタ 6 0 の保守を促す。 I / O アダプタ 6 0 の障害の種類が (2) または (3) の場合は、直ちに I / O アダプタ 6 0 の交換を行う必要は無いが、 I / O アダプタ 6 0 が経年劣化している可能性を考慮し、予防保守の観点から保守員に通報する。一方、障害の種類が (1) の場合は、直ちに I / O アダプタ 6 0 の交換を行う必要がある旨をコンソール 8 0 に通知する。

【 0 0 6 1 】

< 3 . 2 . 2 . V F 障害処理 >

40

図 1 0 は、図 8 の 8 4 0 で行われる V F 障害処理のフローチャートの例である。障害が発生した V F 1 9 0 を特定したハイパバイザ 2 0 は、 V F 状態表 2 1 6 の、当該 V F 1 9 0 のエントリを“交換保留”に変更する (1 0 6 0)。更に、ハイパバイザ 2 0 がアダプタ割り当て表 2 0 0 を参照し、障害が起きた V F 1 9 0 を割り当てた仮想計算機 3 0 を特定する (1 0 4 0)。

【 0 0 6 2 】

続いてハイパバイザ 2 0 は、当該仮想計算機 3 0 に含まれる仮想 C h i p S e t 3 0 0 の H P レジスタ 1 8 5 に V F 1 9 0 の稼働時取り外し (H o t R e m o v e) 要求を設定する。そして、ハイパバイザ 2 0 は、仮想計算機 3 0 に H o t P l u g 割り込みを通知して、ゲスト O S 4 0 に V F 1 9 0 の利用中止と安全な取り外し (S a f e l y R e m

50

ove)を指令する(1050)。

【0063】

最後にハイバイザ20は、障害内容と対処結果をコンソール80に出力通報し、I/Oアダプタ60の交換を促す(1020)。VF障害に対しては直ちに交換を行う必要は無いが、I/Oアダプタ60が経年劣化している可能性を考慮し、予防保守の観点から保守員に通報する。

【0064】

<3.3.I/O仮想化処理>

図11は、図7のステップ780で行われるI/O仮想化処理のフローチャートの例である。ゲストによるI/O操作を検出したハイバイザ20は、操作対象(仮想ChipSet300等のレジスタ)及び操作内容(読み出し、書き込み)に応じて従来のI/Oエミュレーション処理を行う(1100)。従来のI/Oエミュレーション処理については、公知または周知の処理を実施すればよいので、説明を省略する。

10

【0065】

続いてハイバイザ20は、ゲストによるI/O操作の内容がVF190の取り外し(HotRemove)許可であったか否かを判定する(1110)。つまり、ハイバイザ20は、ゲストがVF190の取り外しを許可するまで待機する。なお、VF190の取り外し許可の通知は、ゲストが仮想ChipSet300のホットプラグレジスタ185に許可の応答を書き込むことで行う。ハイバイザ20はホットプラグレジスタ185の値が所定の許可応答であれば、ゲストによるI/O操作の内容がVF190の取り外し許可であった場合となる。この場合のみ、ハイバイザ20は、PF状態表214を参照し、当該VF190を生成したPF160がリセット保留中であることを判定する(1120)。

20

【0066】

ステップ1120に於いてPF160がリセット保留中だった場合、ハイバイザ20は、PF状態表214の当該PF160のエントリについて、待ち合わせbitmap520の取り外し許可が得られたVF190に対応するbitを「0」に変更する(1130)。

【0067】

続いてハイバイザ20は、PF状態表214の当該PF160の待ち合わせbitmap520で全bitが「0」に揃ったか否かを判定する(1140)。全bitが「0」に揃った場合のみ、1150に進んでPF回復処理を行う。

30

【0068】

ステップ1120に於いてPF160がリセット保留中でなかった場合は、ハイバイザ20がVF状態表216を参照し、取り外し許可が得られたVF190が交換保留中であるか否かを判定する(1180)。ハイバイザ20は、VF190が交換保留中であった場合のみ、ステップ1190に進んでVF交換処理を行う。

【0069】

<3.3.1.PF回復処理>

図12は、図11のステップ1150で行われるPF回復処理のフローチャートの例である。リセット保留中のPF160から生成された全てのVF190について、ゲストOS40からの取り外し許可を検出したハイバイザ20は、当該PF160をリセットする(1200)。これによって、当該PF160=I/Oアダプタ60はリセットされ再起動する。

40

【0070】

続いてハイバイザ20は、リセットしたPF160の種類に応じたPFドライバ250を呼び出し、PF160の回復に成功したか否かをI/Oアダプタ60に問い合わせる(1210)。ハイバイザ20は、PF160の回復の成否に応じて処理を分岐する(1220)。

【0071】

50

PF160の回復に成功した場合は、ハイバイザ20がPFドライバ250にPF160の再初期化を依頼し、IOV機能を再度有効化してPF160にVF190を再生成させる(1230)。続いてハイバイザ20は、アダプタ割り当て表200を参照して、当該PF160から生成されていた全てのVF190について、障害発生前にどの仮想計算機30に割り当てられていたかを取得し、当該仮想計算機30に含まれる仮想ChipSet300のHPレジスタ185にVF190の稼働時取り付け(Hot Add)告知(または通知)を設定する。そして、ハイバイザ20は、仮想計算機30にHot Plug割り込みを通知して、ゲストOS40にVF190の増設認識(再認識)と初期化及び利用開始を促す(1240)。これにより、Hot Plug割り込みを受信したゲストOS40は、仮想ChipSet300のHPレジスタ185を読み込んで、Hot Addが発生したことを取得し、PF160が再生成したVF190の再検知を行う。ゲストOS40は既にロードしていたVFドライバ260で再検知したVF190の利用を再開する。

10

【0072】

続いてハイバイザ20は、PF状態表214の当該PF160に対応するエントリの状態510を正常状態に戻す(1250)。

【0073】

一方、ステップ1220でPF160の回復に失敗した場合は、ハイバイザ20がPF状態表214の障害PFのエントリの状態510を交換保留に変更する(1270)。ハイバイザ20によるPF回復処理の最後に、当該処理の結果をコンソール80に出力して保守要員によるI/Oアダプタ60の交換を促す(1260)。回復に成功した場合は直ちに交換を行う必要は無いが、I/Oアダプタ60が経年劣化している可能性を考慮し、予防保守の観点から保守要員に通報する。

20

【0074】**<3.3.2.VF交換処理>**

図13は、図11のステップ1190で行うVF回復処理のフローチャートの例である。交換保留中のVF190について、ゲストOS40からの取り外し許可を検出したハイバイザ20は、アダプタ割り当て表200を参照し、当該VF190と同じPF160から生成された別のVF190を代替VFとして割り当て可能かを判断する(1400)。

30

【0075】

代替VFを割り当て可能な場合のみ、ハイバイザ20はアダプタ割り当て表200を更新して、代替VFを当該ゲストOS40が動いている仮想計算機30に割り当て直す(1420)。続いてハイバイザ20は、当該仮想計算機30に含まれる仮想ChipSet300のHPレジスタ185に代替VFの稼働時取り付け(Hot Add)の告知を設定させる。そして、ハイバイザ20は仮想計算機30にHot Plug割り込みを通知して、ゲストOS40に代替VFの増設認識と初期化及び利用開始を促す(1430)。これにより、Hot Plug割り込みを受信したゲストOS40は、仮想ChipSet300のHPレジスタ185を読み込んで、Hot Addが発生したことを取得し、PF160が再生成した他のVF190の追加を行う。ゲストOS40は追加したVF190の種類に応じたVFドライバ260で代替のVF190の利用を開始する。

40

【0076】

ハイバイザ20によるVF交換処理の最後に、当該処理の結果をコンソール80に出力し、保守要員によるI/Oアダプタ60の交換を促す(1410)。交換に成功した場合は直ちに交換を行う必要は無いが、I/Oアダプタ60が経年劣化している可能性を考慮し、予防保守の観点から保守員に通報する。

【0077】**<4.まとめ>**

以上の構成及び処理により、障害発生時を除いてVF190をゲストOS40に直接操作させるため、高い性能を実現できる。

50

【0078】

P F 1 6 0 で軽度な障害が発生した場合は、当該 P F 1 6 0 のリセットによって消失する全ての V F 1 9 0 を予め仮想計算機 3 0 から H o t R e m o v e して、ゲスト O S 4 0 のクラッシュを回避できるため高信頼を実現できる。また、保守要員の手を借りずに P F 1 6 0 を回復させ、再生成された V F 1 9 0 を自動的に仮想計算機 3 0 に H o t A d d するため、I / O アダプタ 6 0 の保守を容易にすることができる。

【0079】

また、V F 1 9 0 で軽度及び重度な障害が発生した場合についても、保守要員の手を借りずにハイバイザ 2 0 が当該 V F 1 9 0 を仮想計算機 3 0 から H o t R e m o v e する。そして、ハイバイザ 2 0 は、代替 V F 1 9 0 を割り当て直して仮想計算機 3 0 に H o t A d d する。これにより仮想計算機 3 0 は代替 V F 1 9 0 で回復することができるため、V F 1 9 0 の障害に対しても高信頼と保守容易性を実現できる。

10

【0080】

< 第 2 実施形態 >

本発明の第 2 実施形態では、ハイバイザ 2 0 が V F 1 9 0 を利用して V F 1 9 0 の機能を提供する Emulated Function 3 1 0 (以下、E F とする) を生成し、仮想計算機 3 0 に E F 3 1 0 を割り当てる場合の構成及び動作を説明する。本第 2 実施形態では、ゲスト O S 4 0 では無くハイバイザ 2 0 が V F 1 9 0 を直接操作するため、V F 障害に対してより適切な処理が可能となる。なお、E F 3 1 0 は、V F 1 9 0 と同様にゲストに対して仮想 I / O アダプタとして振る舞う。

20

【0081】

< 5 . ハードウェア構成 >

本第 2 の実施形態の計算機システムのハードウェア構成は、前記第 1 実施形態の図 1 と同様であるため説明を省略する。

【0082】

< 6 . ソフトウェア構成 >

次に、物理計算機 1 0 上で仮想計算機 3 0 を実現するソフトウェアの構成の主要部と、制御対象となるハードウェア要素について、図 1 4 を参照しながら詳述する。実施形態 1 に於ける図 2 と同一の符号の同一構成要素については説明を省略する。

【0083】

物理計算機 1 0 は、前記第 1 の実施形態と同一である。ハイバイザ 2 0 は、任意の V F 1 9 0 を用いて E F 3 1 0 を生成し、任意の仮想計算機 3 0 に E F 3 1 0 を割り当てる。E F 3 1 0 は V F 1 9 0 と 1 対 1 に対応し、V F 1 9 0 と同じ機能を提供する。E F 3 1 0 は機能を提供するインタフェースが V F 1 9 0 と違っていても構わない。ハイバイザ 2 0 は、V F 1 9 0 を使用するために V F 1 9 0 の種類に応じた V F ドライバ 2 6 0 (V F 2 6 0 - 1 ~ 2 6 0 - m) を保持する。またハイバイザ 2 0 は、E F 3 1 0 を実現するために、仮想計算機 3 0 のエミュレーションデータ 2 2 0 の内部に E F データ 2 2 4 を保持する。

30

【0084】

仮想計算機 3 0 は、ハイバイザ 2 0 によって提供される仮想 C h i p S e t 3 0 0 等の仮想的な部品と、E F 3 1 0 を含む。仮想計算機 3 0 の上ではゲスト O S 4 0 が動作する。ゲスト O S 4 0 は、E F 3 1 0 に対応する仮想アダプタドライバ 2 7 0 を用いて E F 3 1 0 を操作する。

40

【0085】

図 1 5 はハイバイザ 2 0 が管理するメモリ 9 0 の一例を示す。前記第 1 実施形態に於ける図 3 と同一の構成については説明を省略する。各仮想計算機 3 0 の領域には、ゲスト O S 4 0 と仮想アダプタドライバ 2 7 0 が格納される。

【0086】

ハイバイザ 2 0 が使用する領域には、アダプタ割り当て表 2 0 0、障害処理部 2 1 0、仮想計算機のエミュレーションデータ 2 2 0 と、P F ドライバ 2 5 0 に加えて、V F ド

50

ライバ 260 が格納される。

【0087】

アダプタ割り当て表 200 の構成は、前記第 1 実施形態の図 4 と同一であるため説明を省略する。また、PF 状態表 214 の構成は、前記第 1 実施形態の図 5 と同一であるため説明を省略する。

【0088】

VF 状態表 216 の構成は、VF 190 と EF 310 は 1 対 1 で対応するため、EF 310 には VF 190 と同じ通し番号を付与する。よって本第 2 実施形態では、前記第 1 実施形態の図 6 に示した VF # 410 は仮想アダプタ # の意味を兼ねる。

【0089】

< 7 . ハイパバイザによる処理 >

次に、ハイパバイザ 20 が行う処理の一例について、以下、フローチャートを参照しながら説明する。

【0090】

< 7 . 1 . ハイパバイザによる処理の概要 >

ハイパバイザ 20 が行う処理の全体像を示すフローチャートは、前記第 1 実施形態の図 7 と同一であるため説明を省略する。

【0091】

< 7 . 2 . I / O 障害処理 >

図 7 のステップ 790 で行われる I / O 障害処理は、前記第 1 実施形態の図 8 の処理と同一であるため説明を省略する。

【0092】

< 7 . 2 . 1 . PF 障害処理 >

図 16 は、図 8 のステップ 830 で行われる PF 障害処理のフローチャートの例である。前記第 1 実施形態の図 9 と同一符号の同一処理については説明を省略する。

【0093】

本第 2 の実施形態では前記第 1 実施形態で示した (1) 重度障害及び (2) 軽度障害のケースで更に、ハイパバイザ 20 がアダプタ割り当て表 200 を参照し、障害が起きた PF 160 から生成された全ての VF 190 と、各 VF 190 と 1 対 1 に対応する EF 310 を割り当てた仮想計算機 30 を特定する (1740) 。

【0094】

続いてハイパバイザ 20 は、各仮想計算機 30 に含まれる仮想 Chip Set 300 の HP レジスタ 185 に EF 310 の稼働時取り外し (Hot Remove) 要求を設定させる。ハイパバイザ 20 は、仮想計算機 30 に Hot Plug 割り込みを通知して、ゲスト OS 40 に EF 310 の安全な取り外し (Safely Remove) を指令する (1750) 。

【0095】

< 7 . 2 . 2 . VF 障害処理 >

図 17 は、図 8 のステップ 840 で行われる VF 障害処理のフローチャートの例である。なお、以下では前記第 1 実施形態の図 10 と同一符号の同一処理については説明を省略する。

【0096】

障害が発生した VF 190 を特定したハイパバイザ 20 は、当該 VF 190 の種類に応じた VF ドライバ 260 に、障害からの回復が可能か否かを問い合わせ、障害の重度を取得する (1800) 。 I / O アダプタ 60 の障害の程度には前記第 1 実施形態で述べたとおり 3 通りの重度があり、重い障害から順に、(1) 回復不能で VF 190 の交換が必要となる重度障害、(2) VF 190 をリセットすれば回復できる可能性のある軽度障害、(3) VF 190 のレジスタの一部を再設定する等の手段で回復を完了した障害に分類される。

【0097】

10

20

30

40

50

ハイバイザ20は、I/Oアダプタ60の障害の重度に応じて処理を分岐させる(1810)。(1)回復不能な重度障害であれば、VF状態表216の該当するVF190のエントリの状態620を、「交換保留」に変更する(1830)。(2)リセットで回復できる可能性がある軽度障害であれば、VF状態表216の該当するVF190のエントリの状態620を、「リセット保留」に変更する(1860)。(3)回復完了した障害であれば、ハイバイザ20に回復処理は省略する。

【0098】

上記障害の程度が(1)及び(2)のケースでは更に、ハイバイザ20がアダプタ割り当て表200を参照し、障害が起きたVF190と1対1に対応するEF310と、当該EF310を割り当てた仮想計算機30を特定する(1840)。

10

【0099】

続いてハイバイザ20は、各仮想計算機30に含まれる仮想ChipSet300のHPレジスタ185にEF310の稼働時取り外し(Hot Remove)要求を設定させる。そして、ハイバイザ20は、仮想計算機30にHotPlug割り込みを通知して、ゲストOS40にEF310の利用中止と安全な取り外し(Safely Remove)を指令する(1850)。

【0100】

<7.3.I/O仮想化処理>

図18は、前記第1実施形態の図7のステップ780で行われるI/O仮想化処理のフローチャートの例である。以下では、前記第1実施形態の図11と同一符号の同一処理については説明を省略する。

20

【0101】

ゲストによるI/O操作を検出したハイバイザ20は、操作対象(仮想ChipSet300等のレジスタ)及び操作内容(読み出し、書き込み)に応じて従来のI/Oエミュレーション処理を行う(1900)。この際、EF310のインタフェースが操作された場合は、ハイバイザ20が、VF190の同一機能のインタフェースを代行操作する。

【0102】

続いてハイバイザ20は、ゲストによるI/O操作の内容がEF310の取り外し(Hot Remove)許可であったか否かを判定する(1910)。ゲストによるI/O操作の内容がVF190の取り外し許可であった場合のみ、ハイバイザ20は、PF状態表214を参照し、当該VF190を生成したPF160がリセット保留中であるか否かを判断する(1120)。

30

【0103】

ステップ1120に於いてPF160がリセット保留中でなかった場合は、ハイバイザ20がVF状態表216を参照し、取り外し許可が得られたEF310と1対1に対応するVF190がリセット保留中であるか否かを判定する(1960)。当該VF190がリセット保留中であつた場合のみ、1970に進んでVF回復処理を行う。

【0104】

<7.3.1.PF回復処理>

図19は、図18のステップ1150で行われるPF回復処理のフローチャートの例である。以下では、前記第1実施形態の図12と同一符号、同一処理については説明を省略する。

40

【0105】

PF160の回復に成功した場合は、ハイバイザ20がPFドライバ250にPF160の再初期化を依頼し、IOV機能を再度有効化してVF190を再生成させる(1230)。続いてハイバイザ20は、アダプタ割り当て表200を参照して、当該PF160から生成されていた全てのVF190について、各VF190と1対1に対応するEF310が障害発生前にどの仮想計算機30に割り当てられていたかを取得する。そして、ハイバイザ20は、EF310に対応する仮想計算機30にEF310を再初期化さ

50

せる(2035)。更にハイパバイザ20は、当該仮想計算機30に含まれる仮想ChipSet300のHPレジスタ185にEF310の稼働時取り付け(Hot Add)告知を設定する。そして、ハイパバイザ20は、仮想計算機30にHot Plug割り込みを通知して、ゲストOS40にEF310の増設認識と初期化及び利用開始を促す(2040)。

【0106】

<7.3.2.VF回復処理>

図20は、図18のステップ1970で行われるVF回復処理のフローチャートの例である。リセット保留中のVF190と1対1に対応するEF310について、ゲストOS40からの取り外し許可を検出したハイパバイザ20は、当該VF190をリセットする(2100)。

10

【0107】

続いてハイパバイザ20は、VF190の種類に応じたVFドライバ260を呼び出し、VF190の回復に成功したか否かをVFドライバ260に問い合わせる(2110)。ハイパバイザ20は、VF190の回復の成否に応じて処理を分岐する(2120)。

【0108】

VF190の回復に成功した場合は、ハイパバイザ20が当該EF310を再初期化するように仮想計算機30に指令する(2130)。続いてハイパバイザ20は、当該EF310を有する仮想計算機30に含まれる仮想ChipSet300のHPレジスタ185にEF310の稼働時取り付け(Hot Add)告知を設定する。そして、ハイパバイザ20は仮想計算機30にHot Plug割り込みを通知して、ゲストOS40にEF310の増設認識と初期化及び利用開始を促す(2140)。これにより、Hot Plug割り込みを受信したゲストOS40は、仮想ChipSet300のHPレジスタ185を読み込んで、Hot Addが発生したことを取得し、初期化されたEF310を加え、仮想アダプタドライバ270でEF310の利用を再開する。

20

【0109】

続いてハイパバイザ20は、VF状態表216の当該VF190に対応するエントリの状態620を正常状態に戻す(2150)。最後にハイパバイザ20は、当該処理の結果をコンソール80に出力し、保守要員によるI/Oアダプタ60の交換を促す(2160)。回復に成功した場合は直ちに交換を行う必要は無いが、I/Oアダプタ60が経年劣化している可能性を考慮し、予防保守の観点から保守員に通報する。

30

【0110】

一方、ステップ2120でVF190の回復に失敗した場合は、ハイパバイザ20がVF状態表216の障害VFのエントリの状態620を交換保留に変更する(2170)。続いてハイパバイザはVF交換処理を行う(2180)。

【0111】

<7.3.3.VF交換処理>

図21は、図18のステップ1190及び図20のステップ2180で行うVF回復処理のフローチャートの例である。以下では、前記第1実施形態の図13と同一符号、同一処理については説明を省略する。

40

【0112】

ステップ1420で代替VFを仮想計算機30に割り当て直したのち、ハイパバイザ20は、代替VFと1対1に対応する代替EF310を初期化する(2225)。更にハイパバイザは、当該仮想計算機30に含まれる仮想ChipSet300のHPレジスタ185に代替EF310の稼働時取り付け(Hot Add)告知を設定する。そして、ハイパバイザ20は、仮想計算機30にHot Plug割り込みを通知して、ゲストOS40に代替EF310の増設認識と初期化及び利用開始を促す(1430)。

【0113】

<8.まとめ>

以上の構成及び処理により、ハイパバイザ20の処理がEF310とVF190のイン

50

タフェース変換のみとなるため、前記第1実施形態には劣るが、比較的高い性能を実現できる。

【0114】

P F 1 6 0 で軽度な障害が発生した場合は、当該 P F 1 6 0 のリセットによって消失する全ての V F 1 9 0 を予め仮想計算機 3 0 から H o t R e m o v e しておき、ゲスト O S 4 0 のクラッシュを回避できるため高信頼を実現できる。また、保守要員の手を借りずに P F 1 6 0 を回復させ、再生成された V F 1 9 0 を自動的に仮想計算機 3 0 へ H o t A d d するため、保守容易性を実現できる。

【0115】

また、V F 1 9 0 で軽度及び重度な障害が発生した場合についても、保守要員の手を借りずに当該 V F 1 9 0 を仮想計算機 3 0 から H o t R e m o v e し、障害の重度に応じて V F 1 9 0 をリセットまたは交換する対応をとることができる。そして、E F 3 1 0 を仮想計算機 3 0 に H o t A d d して仮想計算機を回復させるため、V F 1 9 0 の障害に対しても高信頼性と保守容易性を実現できる。

10

【0116】

なお、上記各実施形態では、物理計算機 1 0 上に仮想計算機を生成する仮想化部としてハイパバイザを適用した例を示したが、V M M (Virtual Machine Monitor) を用いることができる。

【0117】

また、上記各実施形態では、物理計算機 1 0 にチップセット 1 0 0 が C P U 7 0 から独立して存在する例を示したが、C P U 7 0 内にチップセット 1 0 0 の機能を含めるようにしてもよい。

20

【産業上の利用可能性】

【0118】

以上のように、本発明は、I O V 機能と A E R 機能を備えた I / O アダプタを備えた仮想計算機や仮想計算機の管理システムに適用することができる。

【符号の説明】

【0119】

- 1 0 物理計算機
- 2 0 ハイパバイザ
- 3 0 仮想計算機
- 4 0 ゲスト O S
- 6 0 I / O アダプタ
- 7 0 C P U
- 8 0 コンソール
- 9 0 メモリ
- 1 0 0 チップセット
- 1 1 0 インターコネクト
- 1 2 0 バス
- 1 3 0 L A N
- 1 4 0 ディスク
- 1 5 0 S A N
- 1 6 0 P F
- 1 7 0 I O V レジスタ
- 1 8 0 A E R レジスタ
- 1 8 5 H P レジスタ (H o t P l u g レジスタ)
- 1 9 0 V F
- 2 0 0 アダプタ割り当て表
- 2 1 0 障害処理部
- 2 1 2 I O V 障害処理部

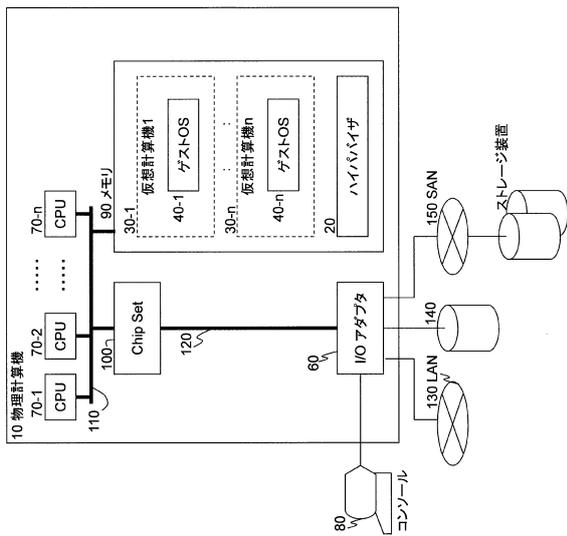
30

40

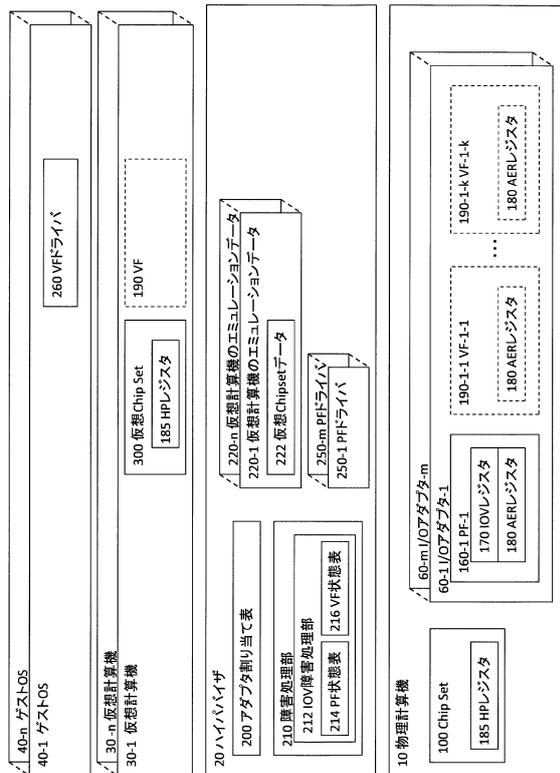
50

- 2 1 4 P F 状態表
- 2 1 6 V F 状態表
- 2 2 0 仮想計算機のエミュレーションデータ
- 2 2 2 仮想Chipsetデータ
- 2 2 4 仮想アダプタデータ
- 2 5 0 P F ドライバ
- 2 6 0 V F ドライバ
- 2 7 0 仮想アダプタドライバ
- 3 0 0 仮想ChipSet
- 3 1 0 仮想アダプタ

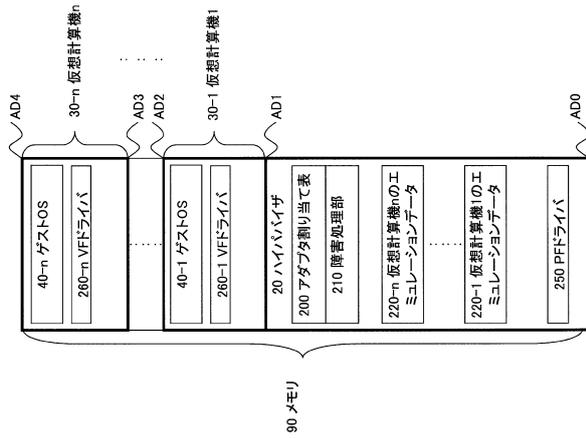
【 図 1 】



【 図 2 】



【図3】



【図4】

PF#	VF#	仮想計算機番号
PF1	VF1-1	0
PF1	VF1-2	1
:	:	:
PF8	VF8-8	未割り当て

400 410 420
 ~~~~~  
 ~~~~~

【図5】

PF#	状態	520 待ち合わせbitmap
PF1	リセット保留	00000011
PF2	正常	00000000
PF3	交換保留	00010000
:	:	:
PF8	正常	00000000

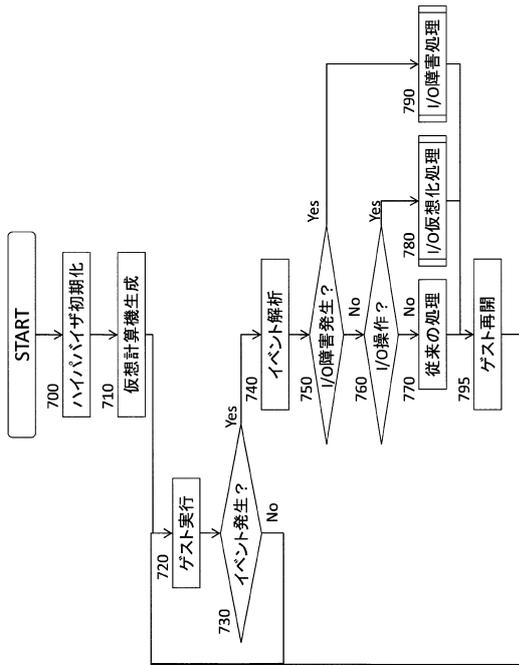
400 510 520
 ~~~~~  
 ~~~~~

【図6】

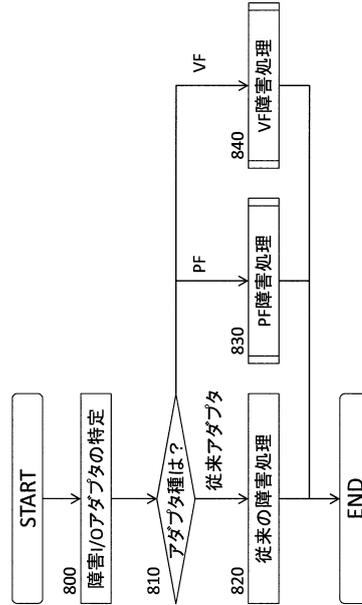
PF#	VF#	状態
PF1	VF1-1	リセット保留
PF1	VF1-2	正常
PF1	VF1-3	交換保留
:	:	:
PF8	VF8-8	正常

400 410 620
 ~~~~~  
 ~~~~~

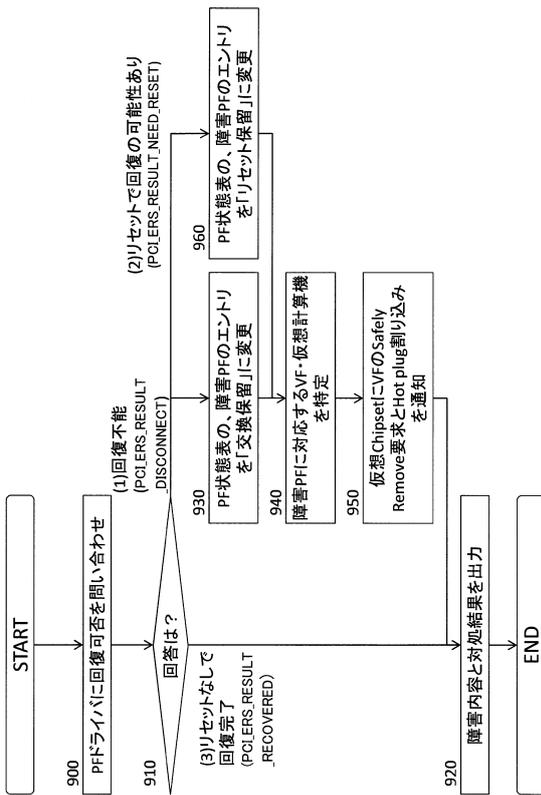
【 図 7 】



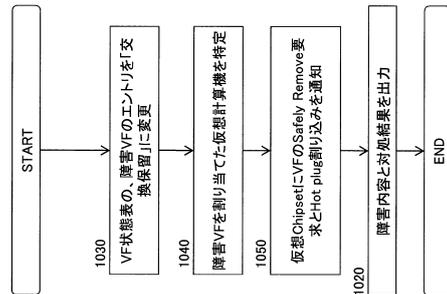
【 図 8 】



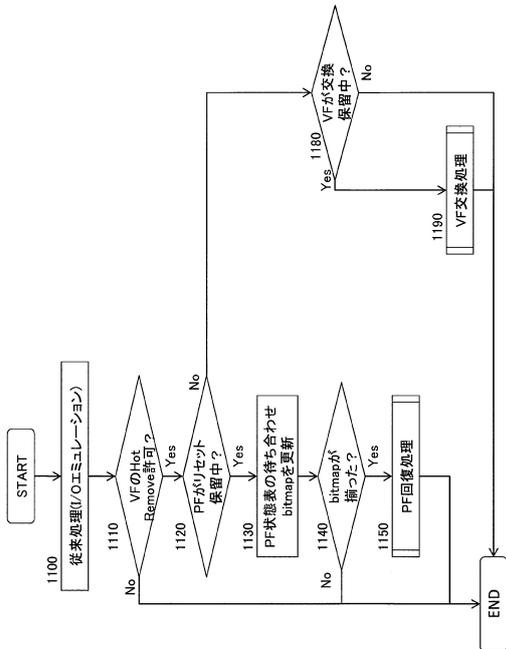
【 図 9 】



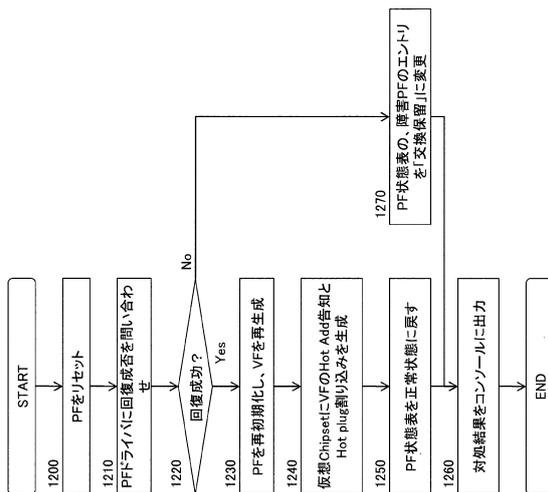
【 図 10 】



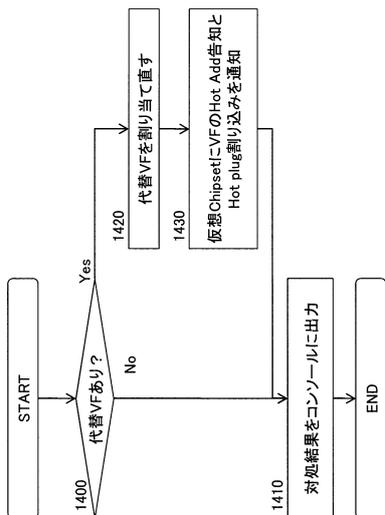
【 図 1 1 】



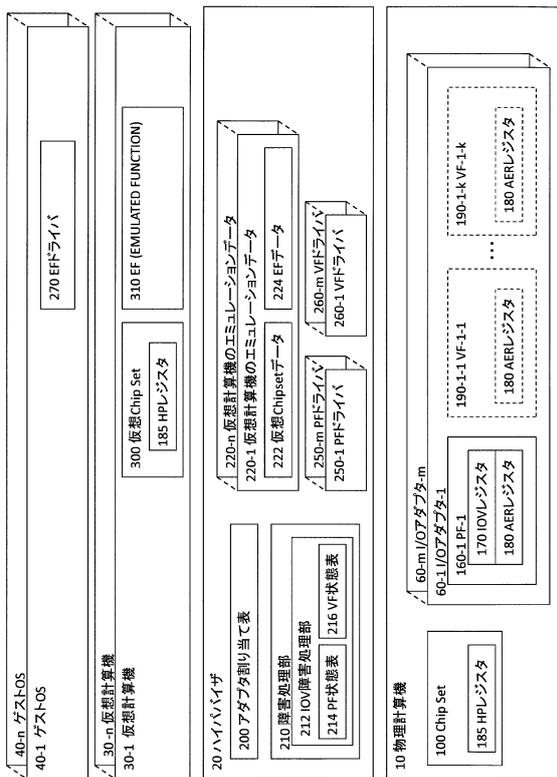
【 図 1 2 】



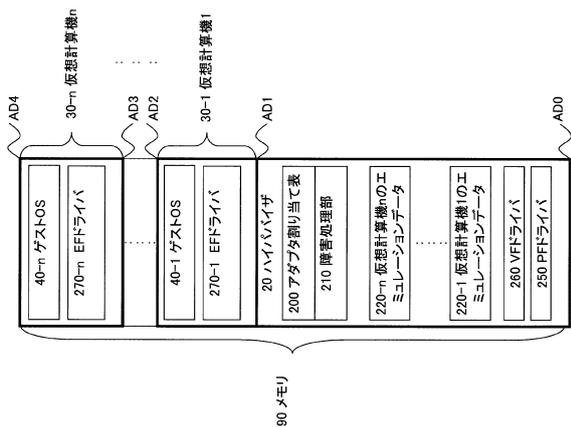
【 図 1 3 】



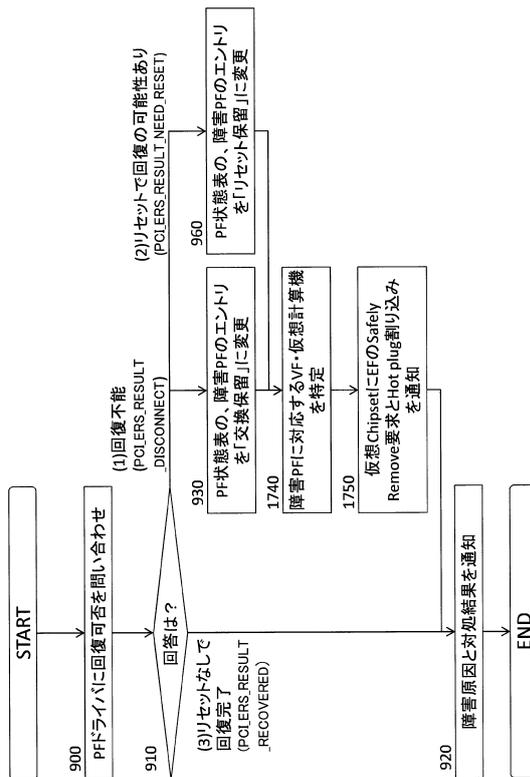
【 図 1 4 】



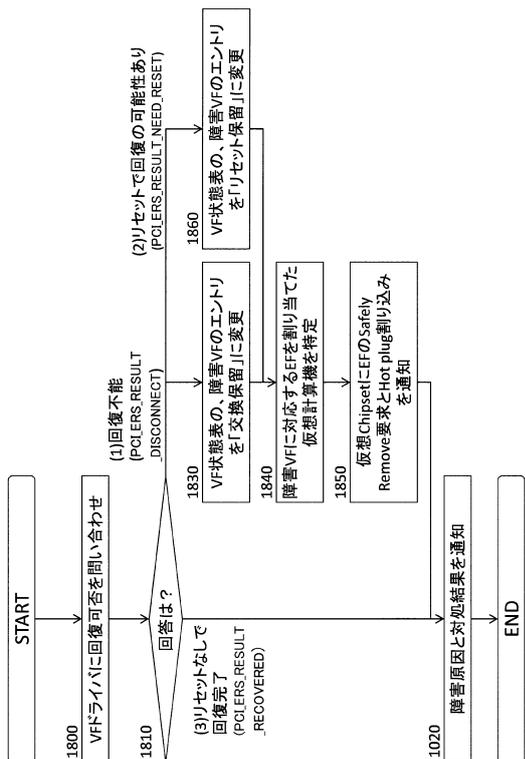
【 図 1 5 】



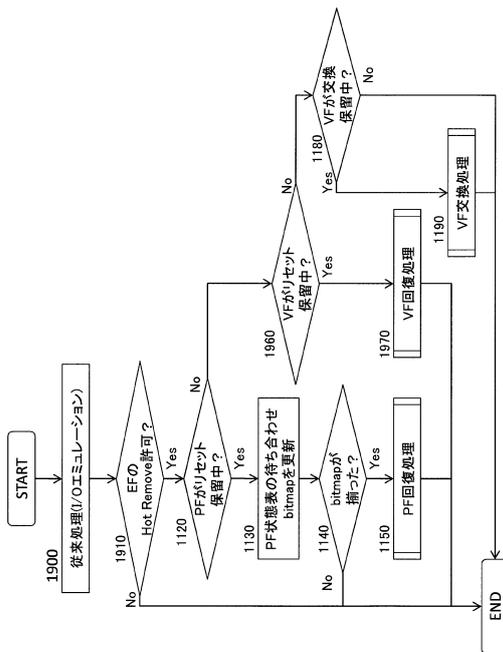
【 図 1 6 】



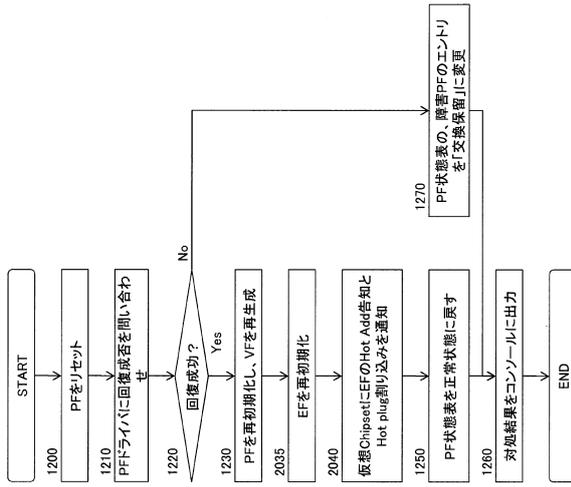
【 図 1 7 】



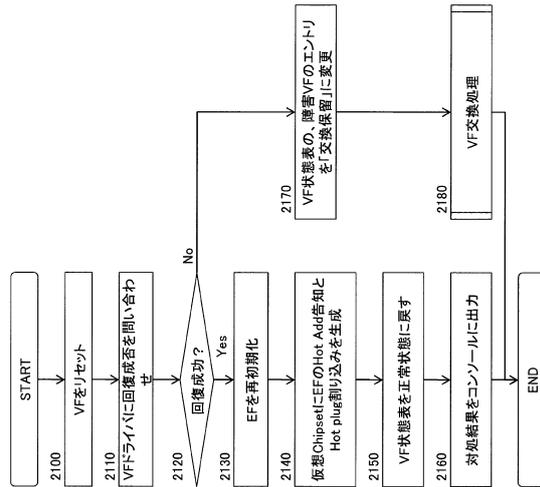
【 図 1 8 】



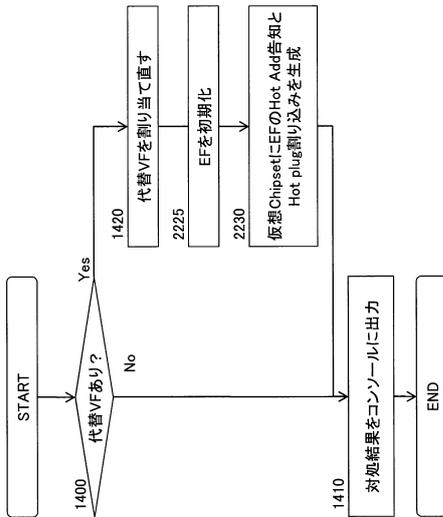
【 図 19 】



【 図 20 】



【 図 21 】



フロントページの続き

(72)発明者 上原 敬太郎

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

審査官 井上 宏一

(56)参考文献 特開2009-301162(JP,A)

特開2010-128911(JP,A)

特開2008-146566(JP,A)

特開2010-39729(JP,A)

特表2013-546111(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/10

G06F 9/46 - 9/54