

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2010-3257

(P2010-3257A)

(43) 公開日 平成22年1月7日(2010.1.7)

(51) Int.Cl.

F I

テーマコード (参考)

G 0 6 F 9/50 (2006.01)

G 0 6 F 9/46 4 6 5 Z

G 0 6 F 9/46 (2006.01)

G 0 6 F 9/46 3 5 0

審査請求 有 請求項の数 12 O L (全 31 頁)

(21) 出願番号 特願2008-163686 (P2008-163686)
(22) 出願日 平成20年6月23日 (2008. 6. 23)

(71) 出願人 000003078
株式会社東芝
東京都港区芝浦一丁目1番1号
(71) 出願人 301063496
東芝ソリューション株式会社
東京都港区芝浦一丁目1番1号
(74) 代理人 100058479
弁理士 鈴江 武彦
(74) 代理人 100108855
弁理士 蔵田 昌俊
(74) 代理人 100091351
弁理士 河野 哲
(74) 代理人 100088683
弁理士 中村 誠

最終頁に続く

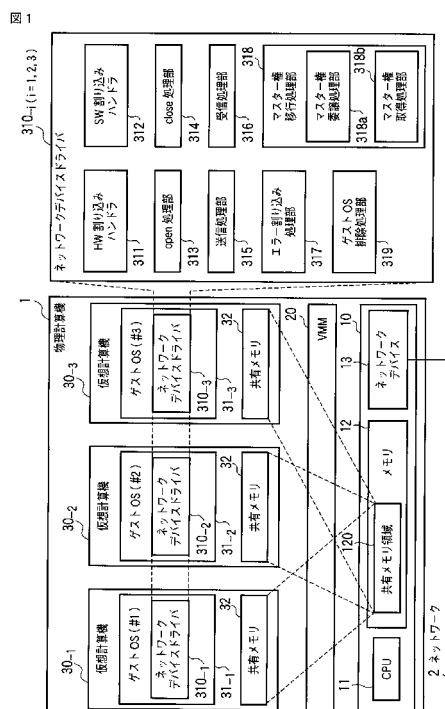
(54) 【発明の名称】 仮想計算機システム及びネットワークデバイス共有方法

(57) 【要約】

【課題】 バッファのデータのコピーによるオーバーヘッドを招くことなく、複数のゲストOSがネットワークデバイスを共有できるようにする。

【解決手段】 VMM 20は、物理計算機1のメモリ12に確保される共有メモリ領域120を用いて共有メモリ32を構築する。共有メモリ32には、ネットワークデバイス13を介してのデータ転送に用いられるバッファが配置される。ネットワークドライバ310-1~310-3は、自身が動作するゲストOSがマスターゲストOSである場合、バッファに格納された受信データを、ネットワークデバイス13からの受信完了割り込みに応じて宛先のゲストOSによって受信させ、マスターゲストOSからの送信要求または当該マスターゲストOS以外のゲストOSからの送信要求割り込みに応じて、バッファに格納された送信データをネットワークデバイス13によってネットワーク2へ送信させる。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

CPU、メモリ及びネットワークデバイスを含む物理計算機上で動作して、複数のゲストOSがそれぞれ動作可能な仮想計算機実行環境を構築する仮想計算機マネージャと、

前記複数のゲストOS上でそれぞれ動作する複数のネットワークデバイスドライバであって、自身が動作するゲストOSがマスター権を有するマスターゲストOSである場合、前記ネットワークデバイスからの割り込みを処理する複数のネットワークデバイスドライバと、

前記メモリに確保される共有メモリ領域を用いて前記仮想計算機マネージャによって構築され、前記複数のゲストOS及び前記複数のネットワークデバイスドライバからアクセス可能な共有メモリであって、前記ネットワークデバイスによってネットワークから受信された受信データを一時格納すると共に、前記複数のゲストOSのいずれかから転送された送信データを一時格納するためのバッファが確保された共有メモリとを具備し、

前記複数のネットワークデバイスドライバの各々は、

自身が動作するゲストOSが前記マスターゲストOSである場合、前記受信データが前記バッファに格納されたことを通知するための前記ネットワークデバイスからの受信完了割り込みに応じて、当該受信データを当該受信データの宛先のゲストOSによって受信させる受信処理手段と、

自身が動作するゲストOSが前記マスターゲストOSである場合、当該マスターゲストOSからの送信要求または当該マスターゲストOS以外のゲストOSからの送信要求割り込みに応じて、前記バッファに格納された前記送信データを前記ネットワークデバイスによって前記ネットワークへ送信させる送信処理手段とを含む

ことを特徴とする仮想計算機システム。

【請求項 2】

前記ネットワークデバイスドライバの前記送信処理手段は、当該ネットワークデバイスドライバが動作するゲストOSが前記マスターゲストOSである場合には、当該マスターゲストOSからの送信要求に応じて、当該マスターゲストOSからの送信データを前記バッファに格納し、しかる後に当該バッファに格納された前記送信データを前記ネットワークデバイスによって前記ネットワークへ送信させ、当該ネットワークデバイスドライバが動作するゲストOSが前記マスターゲストOSとは別のゲストOSである場合、当該別のゲストOSからの送信要求に応じて、当該別のゲストOSからの送信データを前記バッファに格納し、しかる後に前記マスターゲストOSに前記送信要求割り込みを発行することを特徴とする請求項 1 記載の仮想計算機システム。

【請求項 3】

前記複数のネットワークデバイスドライバの各々は、自身が動作するゲストOSが前記マスターゲストOSであって、且つ当該ゲストOSから他のゲストOSへのマスター権の委譲が必要な場合、前記複数のゲストOSの中からマスター権が委譲されるべきゲストOSを選択して、当該選択されたゲストOSにマスター権を委譲するためのマスター権委譲処理を実行するマスター権委譲手段を含む請求項 1 または 2 記載の仮想計算機システム。

【請求項 4】

前記マスター権委譲手段は、前記マスターゲストOSが前記ネットワークデバイスの使用を終了した場合に前記マスター権委譲処理を実行することを特徴とする請求項 3 記載の仮想計算機システム。

【請求項 5】

前記複数のネットワークデバイスドライバの各々は、前記複数のゲストOSのうち自身が動作しているゲストOS以外のゲストOSの稼働状況を監視することによって異常なゲストOSを検出し、当該検出された異常なゲストOSを前記ネットワークデバイスの共有から排除するための排除手段を含み、

前記排除手段は、前記検出された異常なゲストOSが前記マスターゲストOSである場合、当該排除手段を含む前記ネットワークデバイスドライバが動作するゲストOSに前記

10

20

30

40

50

マスター権を取得させるためのマスター権取得処理を実行するマスター権取得手段を含むことを特徴とする請求項 1 乃至 4 のいずれかに記載の仮想計算機システム。

【請求項 6】

前記複数のネットワークデバイスドライバの各々は、前記複数のゲスト OS の通信回数、または前記複数のゲスト OS がディスパッチされる時間もしくは回数、またはディスパッチされる時間もしくは回数の指標となる前記複数のゲスト OS の優先度に基づき、マスターゲスト OS として最適なゲスト OS を決定し、決定されたゲスト OS が現在のマスターゲスト OS と異なる場合、当該決定されたゲスト OS に前記マスター権を移行させるためのマスター権移行処理を実行するマスター権移行手段を含む請求項 1 または 2 記載の仮想計算機システム。

10

【請求項 7】

CPU、メモリ及びネットワークデバイスを含む物理計算機上で動作する仮想計算機マネージャによって提供される仮想計算機実行環境で複数のゲスト OS が動作する仮想計算機システムにおいて、前記ネットワークデバイスを前記複数のゲスト OS で共有させるネットワークデバイス共有方法であって、

前記複数のゲスト OS 上でそれぞれ動作する複数のネットワークデバイスドライバのうちのいずれかのネットワークデバイスドライバが、当該ネットワークデバイスドライバが動作するゲスト OS からのオープン要求に応じて、前記メモリに確保される共有メモリ領域が仮想化された、前記複数のゲスト OS 及び前記複数のネットワークデバイスドライバからアクセス可能な共有メモリであって、前記ネットワークデバイスによってネットワークから受信された受信データを一時格納すると共に、前記複数のゲスト OS のいずれかから転送された送信データを一時格納するためのバッファが確保された共有メモリを前記仮想計算機マネージャによって構築させるためのステップと、

20

前記複数のゲスト OS のうちのマスター権を有するマスターゲスト OS 上で動作する前記ネットワークデバイスドライバが、前記受信データが前記バッファに格納されたことを通知するための前記ネットワークデバイスからの受信完了割り込みに応じて、当該受信データを当該受信データの宛先のゲスト OS によって受信させるステップと、

前記マスターゲスト OS 上で動作する前記ネットワークデバイスドライバが、当該マスターゲスト OS からの送信要求または当該マスターゲスト OS 以外のゲスト OS からの送信要求割り込みに応じて、前記バッファに格納された前記送信データを前記ネットワークデバイスによって前記ネットワークへ送信させるステップと

30

を具備することを特徴とするネットワークデバイス共有方法。

【請求項 8】

前記複数のネットワークデバイスドライバのうち、自身が動作するゲスト OS から前記送信要求を受けたネットワークデバイスドライバが、当該ゲスト OS からの送信データを前記バッファに格納するステップと、

前記送信要求を受けたネットワークデバイスドライバが、送信要求元が前記マスターゲスト OS であるかを判定するステップと、

前記送信要求元が前記マスターゲスト OS でない場合、前記送信要求を受けたネットワークデバイスドライバが、前記マスターゲスト OS に前記送信要求割り込みを発行するステップと

40

を更に具備することを特徴とする請求項 7 記載のネットワークデバイス共有方法。

【請求項 9】

前記複数のネットワークデバイスドライバのうち、前記マスターゲスト OS 上で動作するネットワークデバイスドライバが、前記複数のゲスト OS の中からマスター権が委譲されるべきゲスト OS を選択して、当該選択されたゲスト OS にマスター権を委譲するステップを更に具備することを特徴とする請求項 8 記載のネットワークデバイス共有方法。

【請求項 10】

前記複数のネットワークデバイスドライバの各々が、前記複数のゲスト OS のうち自身が動作しているゲスト OS 以外のゲスト OS の稼働状況を監視することによって異常なゲ

50

ストOSを検出するステップと、

前記検出された異常なゲストOSが前記ゲストOSである場合、前記異常なゲストOSを検出したネットワークデバイスドライバが、自身が動作しているゲストOSに前記マスター権を取得させるステップと、

前記異常なゲストOSを検出したネットワークデバイスドライバが、当該異常なゲストOSを前記ネットワークデバイスの共有から排除するステップと

を更に具備することを特徴とする請求項7乃至9のいずれかに記載のネットワークデバイス共有方法。

【請求項11】

前記複数のネットワークデバイスドライバの各々が、前記複数のゲストOSの通信回数、または前記複数のゲストOSがディスパッチされる時間もしくは回数、またはディスパッチされる時間もしくは回数の指標となる前記複数のゲストOSの優先度に基づき、マスターゲストOSとして最適なゲストOSを決定し、決定されたゲストOSが現在のマスターゲストOSと異なる場合、当該決定されたゲストOSに前記マスター権を移行させるステップを更に具備することを特徴とする請求項7または8記載のネットワークデバイス共有方法。

【請求項12】

前記ネットワークデバイスドライバが前記共有メモリを構築させるためのステップは、当該ネットワークデバイスドライバが動作するゲストOSからのオープン要求に応じて、前記共有メモリが既に構築されているかを当該ネットワークデバイスドライバが判定するステップと、

前記共有メモリが構築されていない場合、前記ネットワークデバイスドライバが前記共有メモリを前記仮想計算機マネージャによって構築させるステップと、

前記共有メモリを構築させた前記ネットワークデバイスドライバが、当該ネットワークデバイスドライバが動作するゲストOSが前記マスターゲストOSであることを示すマスターゲストOS情報を前記共有メモリに登録するステップと

を含むことを特徴とする請求項7記載のネットワークデバイス共有方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、CPU及びメモリを含む物理計算機上で動作する仮想計算機マネージャによって提供される仮想計算機実行環境で複数のゲストOSが動作する仮想計算機システムに係り、特に、単一のネットワークデバイスを当該複数のゲストOSから共有するのに好適な、仮想計算機システム及びネットワークデバイス共有方法に関する

【背景技術】

【0002】

近年、パーソナルコンピュータのような計算機システムにおいて、仮想計算機(Virtual Machine: VM)の技術を応用した研究や開発が行われている。また、仮想計算機システムを実現するための商用のアプリケーションプログラムが実際に広く利用されている。また、メインフレームにおいて古くからハードウェア(HW)で構成された仮想計算機支援機構を用いた仮想計算機システムも存在する。

【0003】

一般に、パーソナルコンピュータのような計算機(物理計算機、実計算機)は、CPU(実CPU)及びメモリ(実メモリ)を含むHW(実HW)で構成されている。この実HWで構成される物理計算機(物理環境)上でハイパバイザOS(オペレーティングシステム)である仮想計算機マネージャ(Virtual Machine Manager: VMM)が動作する。そして、以下に述べるように、VMMの提供する仮想計算機環境のもとで複数のゲストOSが動作をする(特許文献1参照)。

【0004】

VMMは、実CPU、実メモリ及び物理計算機に付属した各種の実入出力デバイス(デ

10

20

30

40

50

ィスドライブ等)のような物理資源の管理を行い、それらを仮想化して仮想計算機(VMM)を構築する。つまりVMMは、物理計算機が有する実CPU、実メモリ及び実入出力デバイスを時間的、空間(領域)的に、複数の仮想計算機のそれぞれ仮想CPU、仮想メモリ及び仮想入出力デバイスとして割り当てることによって、当該複数の仮想計算機が同時に動作可能な環境(仮想計算機環境)を構築する。

【0005】

複数の仮想計算機(が動作する仮想計算機環境)にはそれぞれゲストOSがロードされる。各仮想計算機にロードされたゲストOSは、当該仮想計算機内の仮想CPUによって実行される。つまり、各仮想計算機(が動作する仮想計算機環境)では、ゲストOSが動作する。このようにVMMは、実CPU、実メモリ及び実入出力装置のような物理資源を各ゲストOSに割り当てることによって、当該各ゲストOSが同時に動作できる環境を提供し、仮想計算機システムの管理を行う。

10

【0006】

さて、上記実入出力装置の1つとしてネットワークデバイス(実ネットワークデバイス)が知られている。また、ゲストOSがネットワークデバイスを使用する仕組みとして、以下に挙げる仕組みが従来から知られている。

【0007】

第1は、VMMが仮想ネットワークデバイスを構築して仮想計算機環境に提供することにより、実ネットワークデバイスの管理をVMMが完全に行う仕組み(第1の仕組み)である

20

第2は、1つのゲストOSがマスターとなって実ネットワークデバイスを管理し、他のゲストOSと外部との通信が、マスターとなるゲストOSを経由して行われる仕組み(第2の仕組み)である。

【0008】

第3は、ゲストOSが自身のデバイスドライバによって直接に実ネットワークデバイスを管理する仕組み(第3の仕組み)である。

【特許文献1】特開2006-039763号公報

【発明の開示】

【発明が解決しようとする課題】

【0009】

30

しかしながら、上記第1乃至第3の仕組みは、それぞれ次のような問題を有している。

【0010】

第1の仕組みの問題は、仮想ネットワークデバイスを提供するためにVMM内にドライバを実装しなければならない点である。第1の仕組みの問題はまた、VMMが複雑になり実装のコストも増大する点にもある。また第1の仕組みの問題は更に、状況によっては、ネットワークデバイスの故障に引きずられてVMMの動作環境が機能しなくなり、加えて、本来は影響を受けないはずの仮想計算機環境までが機能しなくなり、結果としてVMM自体の信頼性が低下する点にもある。

【0011】

40

第2の仕組みの問題は、ネットワークデバイスを介してのデータ転送に用いられるバッファ(DMAバッファ領域)のデータを、マスターゲストOS及びその他のゲストOSの仮想メモリ空間相互間でコピーする必要があるため、オーバーヘッドを招いて十分なスループットが得られない点にある。また、第2の仕組みの問題は更に、マスターゲストOSが動作不能に陥れば、他のゲストOSが通信することができなくなる点にもある。

【0012】

第3の仕組みの問題は、デバイスドライバが、他のゲストOSが同じネットワークデバイスを使用することを考慮して作られていないことに起因する。つまり第3の仕組みの問題は、複数のゲストOS(で動作する第3の仕組み)のそれぞれでデバイスドライバが動作するために、複数のゲストOS間の排他的仕組みがなく複数のゲストOSでネットワークデバイスを共有するのが困難である点にある。そのため、一般には1つのネットワーク

50

デバイスは１つのゲストＯＳしか使用することができない。また、ＶＭＭは、どのゲストＯＳにネットワークデバイスを管理させるか（見せるか）を決定する機能を持たなければならない。

【００１３】

本発明は上記事情を考慮してなされたものでその目的は、ネットワークデバイスを介してのデータ転送に用いられるバッファのデータのコピーによるオーバーヘッドを招くことなく、複数のゲストＯＳがネットワークデバイスを共有できる仮想計算機システム及びネットワークデバイス共有方法を提供することにある。

【課題を解決するための手段】

【００１４】

本発明の１つの観点によれば、ＣＰＵ、メモリ及びネットワークデバイスを含む物理計算機上で動作する仮想計算機マネージャによって提供される仮想計算機実行環境で複数のゲストＯＳが動作する仮想計算機システムが提供される。この仮想計算機システムは、前記複数のゲストＯＳ上でそれぞれ動作する複数のネットワークデバイスドライバであって、自身が動作するゲストＯＳがマスター権を有するマスターゲストＯＳである場合、前記ネットワークデバイスからの割り込みを処理する複数のネットワークデバイスドライバと、前記メモリに確保される共有メモリ領域を用いて前記仮想計算機マネージャによって構築され、前記複数のゲストＯＳ及び前記複数のネットワークデバイスドライバからアクセス可能な共有メモリであって、前記ネットワークデバイスによってネットワークから受信された受信データを一時格納すると共に、前記複数のゲストＯＳのいずれかから転送された送信データを一時格納するためのバッファが確保された共有メモリとを具備し、前記複数のネットワークデバイスドライバの各々は、自身が動作するゲストＯＳが前記マスターゲストＯＳである場合、前記受信データが前記バッファに格納されたことを通知するための前記ネットワークデバイスからの受信完了割り込みに応じて、当該受信データを当該受信データの宛先のゲストＯＳによって受信させる受信処理手段と、自身が動作するゲストＯＳが前記マスターゲストＯＳである場合、当該マスターゲストＯＳからの送信要求または当該マスターゲストＯＳ以外のゲストＯＳからの送信要求割り込みに応じて、前記バッファに格納された前記送信データを前記ネットワークデバイスによって前記ネットワークへ送信させる送信処理手段とを含むことを特徴とする。

【発明の効果】

【００１５】

本発明によれば、マスターゲストＯＳとその他のゲストＯＳとが協調することによってネットワークデバイスを共有する仕組みを適用しながら、メモリに確保される共有メモリ領域を用いて構築される、各ゲストＯＳ及び各ネットワークデバイスドライバからアクセス可能な共有メモリに、ネットワークデバイスを介して送受信されるデータを一時格納するためのバッファを配置することにより、当該バッファのデータのコピーによるオーバーヘッドを招くことなく、各ゲストＯＳがネットワークデバイスを共有できる。

【発明を実施するための最良の形態】

【００１６】

以下、本発明の実施の形態につき図面を参照して説明する。

図１は本発明の一実施形態に係る仮想計算機システムの構成を示すブロック図である。図１に示す仮想計算機システムは、物理計算機（実計算機）１を用いて実現される。物理計算機１は、仮想計算機環境を提供するのに用いられるハードウェア（ＨＷ）１０を備えている。ＨＷ１０は、ＨＷ資源（物理資源）であるＣＰＵ（実ＣＰＵ）１１、メモリ（実メモリ）１２及びネットワークデバイス（実ネットワークデバイス）１３を含む。なお、図１では、ディスクドライブのような、物理計算機１に付属する、ネットワークデバイス１３以外の入出力デバイスは省略されている。

【００１７】

ネットワークデバイス１３は、物理計算機１の外部のネットワーク２に接続されている。ネットワークデバイス１３は、ネットワーク２を介して物理計算機１と外部の物理デバ

10

20

30

40

50

イスとの間の通信を行う。

【 0 0 1 8 】

物理計算機 1 (の HW 1 0) 上では、仮想計算機マネージャ (VMM) 2 0 が動作する。VMM 2 0 は、仮想計算機システムの管理を行い、仮想計算機が動作する環境 (仮想計算機環境) を構築する。更に具体的に述べるならば、VMM 2 0 は、物理計算機 1 の HW 1 0 を構成する CPU 1 1、メモリ 1 2 及び、ネットワークデバイス 1 3 を含む各種の入出力デバイスのような物理資源 (ハードウェアリソース) を制御すると共に当該物理資源の時間的、領域的な配分を管理することにより、複数の仮想計算機、例えば 3 台の仮想計算機 3 0 -1, 3 0 -2, 3 0 -3 が動作する仮想計算機環境 (仮想計算機実行環境) を構築する。

10

【 0 0 1 9 】

仮想計算機 3 0 -1, 3 0 -2, 3 0 -3 上では、それぞれゲスト OS 3 1 -1 (# 1) , 3 1 -2 (# 2) , 3 1 -3 (# 3) が動作する。更に具体的に述べるならば、ゲスト OS 3 1 -1 ~ 3 1 -3 は、それぞれ仮想計算機 3 0 -1 ~ 3 0 -3 にロードされて、当該仮想計算機 3 0 -1 ~ 3 0 -3 内の仮想 CPU によって実行される。ゲスト OS 3 1 -1, 3 1 -2, 3 1 -3 は固有の識別子 (ゲスト OS 識別子) を有する。本実施形態においてゲスト OS 3 1 -1, 3 1 -2, 3 1 -3 のゲスト OS 識別子は、それぞれ、1, 2, 3 であるものとする。

【 0 0 2 0 】

VMM 2 0 はまた、ゲスト OS 3 1 -1 ~ 3 1 -3 (の後述するドライバ 3 1 0 -1 ~ 3 1 0 -3) からアクセス (参照 / 書き込み) 可能な共有メモリ 3 2 を構築する。共有メモリ 3 2 は、HW 1 0 に含まれているメモリ 1 2 の記憶領域の一部 (共有メモリ領域) 3 2 を用いて構築される仮想化されたメモリである。VMM 2 0 は、ゲスト OS 3 1 -1 ~ 3 1 -3 の間の通知のために割り込みを発生させるインタフェース (ゲスト OS 間割り込み配信インタフェース) を当該ゲスト OS 3 1 -1 ~ 3 1 -3 に提供する。

20

【 0 0 2 1 】

ゲスト OS 3 1 -1, 3 1 -2, 3 1 -3 には、それぞれネットワークデバイスドライバ (以下、ドライバと略称する) 3 1 0 -1, 3 1 0 -2, 3 1 0 -3 が付加されている。ドライバ 3 1 0 -1 ~ 3 1 0 -3 は、それぞれ、ネットワークデバイス 1 3 をゲスト OS 3 1 -1 ~ 3 1 -3 が使用可能とするためのデバイスドライバであり、当該ゲスト OS 3 1 -1 ~ 3 1 -3 の一部として当該ゲスト OS 3 1 -1 ~ 3 1 -3 上で動作する。つまりゲスト OS 3 1 -1 ~ 3 1 -3 は、それぞれドライバ 3 1 0 -1 ~ 3 1 0 -3 を含む。

30

【 0 0 2 2 】

ゲスト OS 3 1 -1 ~ 3 1 -3 は、それぞれがマスター (マスターゲスト OS) となる機能を持つ。ドライバ 3 1 0 -1 ~ 3 1 0 -3 は、それぞれゲスト OS 3 1 -1 ~ 3 1 -3 がマスターとして動作する場合、自身もマスターとして動作する。本実施形態において、ドライバ 3 1 0 -1, 3 1 0 -2, 3 1 0 -3 は同等の機能を持つ。但し、ドライバ 3 1 0 -1 ~ 3 1 0 -3 の間で本発明に直接関係しない一部の機能が異なっても構わない。例えば、ドライバ 3 1 0 -1 ~ 3 1 0 -3 のうちのいずれかで、マスターとなる機能が制限されていても構わない。

【 0 0 2 3 】

ドライバ 3 1 0 -i (i = 1, 2, 3) は、HW (ハードウェア) 割り込みハンドラ 3 1 1 及び SW (ソフトウェア) 割り込みハンドラ 3 1 2 の 2 種類の割り込みハンドラを有する。HW 割り込みハンドラ 3 1 1 は、ネットワークデバイス 1 3 によって発生される割り込み (HW 割り込み) を処理する。SW 割り込みハンドラ 3 1 2 は、ゲスト OS 間割り込み配信インタフェースによって発生される割り込み (ゲスト OS 3 1 -1 ~ 3 1 -3 の間の通知のための割り込み) を処理する。

40

【 0 0 2 4 】

ドライバ 3 1 0 -i は更に、open (オープン) 処理部 3 1 3、close (クローズ) 処理部 3 1 4、送信処理部 3 1 5、受信処理部 3 1 6、エラー割り込み処理部 3 1 7、マスター権移行処理部 3 1 8、ゲスト OS 排除処理部 3 1 9 を有する。マスター権移行処

50

理部 3 1 8 は、マスター権委譲処理部 3 1 8 a 及びマスター権取得処理部 3 1 8 b を含む。これらの各処理部の機能については後述する。

【 0 0 2 5 】

図 2 は、図 1 に示される共有メモリ 3 2 における領域の割り当て例を示す。共有メモリ 3 2 は、物理レジスタマップ領域 3 2 1、特殊処理ステータス領域 3 2 2、ゲスト OS 情報領域 3 2 3 及び DMA バッファ領域 3 2 4 を含む。共有メモリ 3 2 はまた、送信ディスクリプタチェーン領域 3 2 5、受信ディスクリプタチェーン領域 3 2 6 及び他共有情報領域 3 2 7 を含む。

【 0 0 2 6 】

物理レジスタマップ領域 3 2 1 は、ネットワークデバイス 1 3 に所属する物理レジスタの群をマップした領域である。物理レジスタマップ領域 3 2 1 によってマップされる物理レジスタの群は、初期化処理で設定が必要なレジスタ、及び割り込み要因レジスタを含む。物理レジスタマップ領域 3 2 1 にアクセスすることによって物理レジスタがアクセスされる。

【 0 0 2 7 】

特殊処理ステータス領域 3 2 2 は、ドライバ 3 1 0 -i が行うべき特殊な処理（特殊処理）が発生しているかを示すステータスを格納するのに用いられる。本実施形態において、特殊処理ステータス領域 3 2 2 を用いてステータスが管理される特殊処理は、共有メモリ 3 2 を初期化する処理、ネットワークの物理的切断（ネットワークケーブル抜けなどへの対応）処理、ネットワークの物理的接続処理を含む。これらの特殊処理は、共有メモリ 3 2 内の情報を使用して各ドライバ 3 1 0 -i で進めることができる処理であり、各ドライバ 3 1 0 -i で重複して実行されても問題の生じない処理である。本実施形態において、特殊処理ステータス領域 3 2 2 は、予め定められた特殊処理毎に、その処理が発生しているかをビットの状態を示すビットマップを格納する。ドライバ 3 1 0 -i は、ある特殊処理が発生しているかを確認するには、特殊処理ステータス領域 3 2 2（に格納されているビットマップ）内の当該特殊処理に対応するビットを参照すればよい。

【 0 0 2 8 】

ゲスト OS 情報領域 3 2 3 は、ネットワークデバイス 1 3 を使用している全てのゲスト OS の情報と、マスターゲスト OS 情報とを格納するのに用いられる。ゲスト OS の情報（ゲスト OS 情報）は、ゲスト OS 識別子、SW 割り込み番号及び SW 割り込み要因情報を含む。

【 0 0 2 9 】

ゲスト OS 識別子は、ネットワークデバイス 1 3 を現在使用しているゲスト OS の識別子を示す。

SW 割り込み番号は、ゲスト OS 識別子で示されるゲスト OS に対して他のゲスト OS から、ゲスト OS 間の通知機能（ゲスト OS 間割り込み配信インタフェース）によって割り込みを発生させる場合の割り込み配信先のゲスト OS の SW 割り込みハンドラ 3 1 2 を示す。

【 0 0 3 0 】

SW 割り込み要因情報は、上記割り込みを発生させる場合の通知内容を表す情報。SW 割り込み要因情報は例えばビットマップから構成される。配信先のゲスト OS は、配信元のゲスト OS（ゲスト OS 識別子で示されるゲスト OS）のゲスト OS 間割り込み配信インタフェースから送信されたビットマップ形式の SW 割り込み要因情報（の各ビットの状態）を参照することにより通知内容を判断する。

【 0 0 3 1 】

ゲスト OS 情報は更に、付加的な情報として、タイムスタンプ、通信回数統計情報、ゲスト OS 優先度及びネットワークアドレスを含む。

タイムスタンプは、ゲスト OS 識別子で示されるゲスト OS の最終動作時刻を示す。

【 0 0 3 2 】

通信回数統計情報は、直近の通信回数を示す統計情報である。本実施形態において、通

10

20

30

40

50

信回数統計情報には、1 ディスクリプタチェーンサイクルの期間にゲストOSの使用したディスクリプタの数が用いられる。ディスクリプタチェーンサイクルについては後述する。

【0033】

ゲストOS優先度は、ゲストOS識別子で示されるゲストOSの実行優先度を示す。つまりゲストOS優先度は、VMM20が、ディスパッチされるべきゲストOSを選択する際に選定の基準（指標）とする値である。本実施形態では、ゲストOS優先度の値が大きいゲストOSほど、優先的にCPU11（CPU時間）が割り当てられる。

【0034】

ネットワークアドレスは、ネットワーク2を介しての通信のためにゲストOS識別子で示されるゲストOSに割り当てられる、IP（Internet Protocol）アドレスのようなネットワークアドレスである。

【0035】

マスターゲストOS情報は、マスタとして動作している（つまりマスター権を有している）ゲストOS（マスターゲストOS）を示す。ここでは、マスターゲストOS情報として、マスターゲストOSのゲストOS識別子が用いられる。

【0036】

なお、図2に示されるゲストOS情報領域323には、便宜的に、ゲストOS識別子が「1」のゲストOS（つまりゲストOS31-1）のゲストOS情報のみが格納されている状態が示されている。しかし、ゲストOS情報領域323には、ネットワークデバイス13を共有する他のゲストOS31-2、31-3のゲストOS情報も格納されているものとする。

【0037】

DMAバッファ領域324は、ネットワークデバイス13が送受信データのDMA（Direct Memory Access）転送に使用するデータバッファ（以下、DMAバッファと称する）324aのための領域である。

【0038】

送信ディスクリプタチェーン領域325は、送信ディスクリプタチェーン325aを格納するのに用いられる。送信ディスクリプタチェーン325aは、DMAバッファ324a内の領域（へのリンク）及びHW仕様に基づいた送信ディスクリプタを管理するためのディスクリプタ管理構造を有する。送信ディスクリプタチェーン325aに含まれる各送信ディスクリプタは、チェーンでリンクされている順にサイクリックに利用される。送信ディスクリプタチェーン325aを一巡するサイクルをディスクリプタチェーンサイクルと呼ぶ。

【0039】

図3は、送信ディスクリプタチェーン325aの一例を示す。図3の例では、送信ディスクリプタチェーン325aは、送信ディスクリプタ325-1～325-10を含む。本実施形態において送信ディスクリプタ325-1～325-10はリング状のチェーンによりリンクされている。送信ディスクリプタ325-j（j = 1～10）は、ステータス部3251、ポインタ部3252及び利用ゲストOS情報部（図示せず）を含む。

【0040】

ポインタ部3252は、送信されるべきデータが設定（格納）されているDMAバッファ324a内の領域を指し示すポインタを格納する。利用ゲストOS情報部は、当該利用ゲストOS情報部を含む送信ディスクリプタ325-jが、いずれのゲストOSによって参照（使用）中であることを示す情報（例えばゲストOS識別子）を格納する。

【0041】

ステータス部3251は、送信ディスクリプタ325-jのステータスを示す。送信ディスクリプタ325-jの取り得るステータスは、例えば、「empty」、「not ready」、「ready」及び「dma」の4種類である。

【0042】

10

20

30

40

50

「empty」は、送信ディスクリプタ325-jが未使用の状態にあることを示す。

「not ready」は、送信ディスクリプタ325-jがデータ送信のために使用予定である（使用が予約されている）が、まだ送信されるべきデータの設定が終わっていない状態にあることを示す。

【0043】

「ready」は、送信ディスクリプタ325-jがデータ送信のために使用予定であり、且つ送信されるべきデータの設定も終了している状態にあることを示す。

【0044】

「dma」は送信ディスクリプタ325-jのポインタ部3252で指定されるデータの送信のためのDMA転送が実行中であることを示す。

【0045】

送信ディスクリプタチェーン325aは、送信中先頭ポインタ、送信前先頭ポインタ及び空き先頭ポインタの3種類のポインタで管理される。これらのポインタは、任意の時点において、

送信中先頭ポインタ 送信前先頭ポインタ 空き先頭ポインタ
のような不等号で示される順番の位置関係（大きいほうが進行方向）にある送信ディスクリプタを指し示す。

【0046】

送信中先頭ポインタの指し示す位置から送信前先頭ポインタの指し示す位置の直前までには、ステータスが「dma」の送信ディスクリプタのみが存在する。図3の例では、送信ディスクリプタ325-3、325-4が、これに該当する。

【0047】

送信前先頭ポインタの指し示す位置から空き先頭ポインタの指し示す位置の直前までには、ステータスが「ready」の送信ディスクリプタとステータスが「not ready」の送信ディスクリプタとが混在し得る。図3の例では、送信ディスクリプタ325-5～325-9が、これに該当する。

【0048】

空き先頭ポインタから先には、ステータスが「empty」の送信ディスクリプタのみが存在する。図3の例では、送信ディスクリプタ325-10、325-1、325-2が、これに該当する。

これら3種類のポインタは、ネットワークデバイス13によって操作される。

【0049】

再び図2を参照すると、受信ディスクリプタチェーン領域326は、受信ディスクリプタチェーンを格納するのに用いられる。受信ディスクリプタチェーンは、DMAバッファ324a内の領域（へのリンク）及びHW仕様に基づいた受信ディスクリプタを管理するためのディスクリプタ管理構造を有する。受信ディスクリプタチェーンに含まれる各受信ディスクリプタは、送信ディスクリプタチェーン325aに含まれる上記各送信ディスクリプタと同様に、チェーンでリンクされている順にサイクリックに利用される。この受信ディスクリプタチェーンを一巡するサイクルもディスクリプタチェーンサイクルと呼ばれる。

【0050】

図4は、受信ディスクリプタチェーンの一例を示す。図4の例では、受信ディスクリプタチェーンは、受信ディスクリプタ326-1～326-10を含む。本実施形態において受信ディスクリプタ326-1～326-10はリング状のチェーンによりリンクされている。受信ディスクリプタ326-j（j=1～10）は、送信ディスクリプタ325-jと同様に、ステータス部3261、ポインタ部3262及び利用ゲストOS情報部（図示せず）を含む。

【0051】

ポインタ部3262は、受信されるべきデータが設定されているDMAバッファ324a内の領域を指し示すポインタを格納する。利用ゲストOS情報部は、当該利用ゲストO

10

20

30

40

50

S 情報部を含む受信ディスクリプタ 3 2 6 -j が、いずれのゲスト OS によって参照（使用）中であるかを示す情報（ゲスト OS 識別子）を格納する。

【 0 0 5 2 】

ステータス部 3 2 6 1 は、受信ディスクリプタ 3 2 6 -j のステータスを示す。受信ディスクリプタ 3 2 6 -j の取り得るステータスは、例えば、「 e m p t y 」、「 r e c e i v e d 」及び「 d m a 」の 3 種類である。

【 0 0 5 3 】

「 e m p t y 」は、受信ディスクリプタ 3 2 6 -j が未使用の状態にあることを示す。

「 r e c e i v e d 」は、受信ディスクリプタ 3 2 6 -j のポインタ部 3 2 6 2 で指定される DMA バッファ 3 2 4 a 内の領域にデータを受信済みであるが、その旨をゲスト OS が認識していない状態にあることを示す。

【 0 0 5 4 】

「 d m a 」は、受信ディスクリプタ 3 2 6 -j のポインタ部 3 2 6 2 で指定される DMA バッファ 3 2 4 a 内の領域に受信データを一時格納するための DMA 転送が実行中であることを示す。

【 0 0 5 5 】

受信ディスクリプタチェーンは、受信中先頭ポインタ、DMA 中先頭ポインタ及び空き先頭ポインタの 3 種類のポインタで管理される。これらのポインタは、任意の時点において、

受信中先頭ポインタ DMA 中先頭ポインタ 空き先頭ポインタ

のような不等号で示される順番の位置関係（大きいほうが進行方向）にある受信ディスクリプタを指し示す。受信中先頭ポインタの指し示す位置から DMA 中先頭ポインタの指し示す位置の直前までには、ステータスが「 r e c e i v e d 」の受信ディスクリプタと「 e m p t y 」の受信ディスクリプタとが混在する。図 4 の例では、受信ディスクリプタ 3 2 6 -3, 3 2 6 -4 が、これに該当する。

【 0 0 5 6 】

DMA 中先頭ポインタの指し示す位置から空き先頭ポインタの指し示す位置の直前までには、ステータスが「 d m a 」の受信ディスクリプタのみが存在する。図 4 の例では、受信ディスクリプタ 3 2 6 -5, 3 2 6 -6 が、これに該当する。

【 0 0 5 7 】

空き先頭ポインタから先には、ステータスが「 e m p t y 」の受信ディスクリプタのみが存在する。図 4 の例では、受信ディスクリプタ 3 2 6 -7 ~ 3 2 6 -10, 3 2 6 -1, 3 2 6 -2 が、これに該当する。

これら 3 種類のポインタは、ネットワークデバイス 1 3 によって操作される。

【 0 0 5 8 】

他共有情報領域 3 2 7 は、各ゲスト OS 3 1 -i (i = 1 , 2 , 3) のドライバ 3 1 0 -i 間で共有すべき情報、例えば特殊処理で使用する必要がある情報を格納するのに用いられる。

【 0 0 5 9 】

次に、上述の構成の仮想計算機システムの特徴について説明する。

【 0 0 6 0 】

まず、図 1 の仮想計算機システムでは、ネットワークデバイス 1 3 を共有するゲスト OS 3 1 -1 ~ 3 1 -3 上でそれぞれドライバ 3 1 0 -1 ~ 3 1 0 -3 が動作する。ドライバ 3 1 0 -1 ~ 3 1 0 -3 の 1 つはマスターとして動作して、DMA 処理や、DMA 処理のためのネットワークデバイス 1 3 からの割り込みに対する処理など、当該ネットワークデバイス 1 3 の管理の主要な部分を取り扱う。

【 0 0 6 1 】

図 1 の仮想計算機システムの特徴は、ゲスト OS 3 1 -1 ~ 3 1 -3 上でそれぞれ動作するドライバ 3 1 0 -1 ~ 3 1 0 -3 に共通して必要となる情報、及びゲスト OS 3 1 -1 ~ 3 1 -3 に固有のゲスト OS 情報のうち必要な部分を相互参照できるようにした点にある。この相

10

20

30

40

50

互参照を可能とするために、VMM 20によってゲストOS 31-1~31-3に提供される共有メモリ32が使用される。ゲストOS 31-1~31-3のそれぞれドライバ310-1~310-3は共有メモリ32にアクセスする。

【0062】

共有メモリ32は、HW 10に含まれているメモリ12に確保された共有メモリ領域120を用いて構築される仮想化されたメモリである。このため、共有メモリ32へのアクセスは、物理的にはメモリ12内の共有メモリ領域120へのアクセスによって実現される。

【0063】

本実施形態ではまた、DMAバッファ324aも共有メモリ32に配置される。これにより、DMAバッファ324aも各ゲストOS 31-1~31-3から直接参照（アクセス）することが可能となり、従来は大きな問題であったデータコピーに伴うオーバーヘッドを後述するように削減することができる。

【0064】

また本実施形態では、マスターとなるゲストOSを変更可能とする仕組み、つまりマスターゲストOSが機能しなくなっても他のゲストOSがネットワークデバイス13を使用し続けることを可能とする仕組みが提供される。本実施形態では更に、効率的なネットワークデバイス13の使用を支援する仕組みも提供される。

【0065】

以下、上述の特徴を有する仮想計算機システムにおけるゲストOSのドライバを中心とする動作について、「open（オープン）処理」、「送信処理」、「受信処理」及び「エラー割り込み処理」を例に順次説明する。

【0066】

[open処理]

まず、ゲストOSのドライバ（のopen処理部313）がネットワークデバイス13の使用を開始する際に実行する「open処理（ドライバのopen処理）」の手順について、当該ドライバがゲストOS 31-1のドライバ310-1である場合を例に、図5のフローチャートを参照して説明する。

【0067】

今、ゲストOS 31-1から当該ゲストOS 31-1のドライバ310-1に対して、「open要求」が与えられたものとする。するとドライバ310-1は、共有メモリ32の領域を探して、当該共有メモリ32の領域が存在するかを判定する（ステップS1）。

【0068】

もし、ステップS1の判定がYesであるならば、ドライバ310-1は、マスターとなっている別のゲストOSのドライバからの要求によって共有メモリ32の領域が確保されているとしてステップS8に進む。

【0069】

これに対し、ステップS1の判定がNoであるならば、ドライバ310-1はVMM 20に依頼して、共有メモリ32の領域を確保する（ステップS2）。つまりドライバ310-1は、VMM 20により、HW 10に含まれているメモリ12の記憶領域の一部を共有メモリ領域120として確保させることにより、共有メモリ32（仮想化された共有メモリメモリ）を構築させる。

【0070】

次にドライバ310-1は、共有メモリ32を初期化するための初期化処理を行う（ステップS3）。この初期化処理は、共有メモリ32における領域321~327の割り当て、物理レジスタマップ領域321に物理レジスタの群（ネットワークデバイス13に所属する物理レジスタの群）をマップする処理を含む。また、上記初期化処理は、特殊処理ステータス領域322に当該初期化処理が進行中であることを示すステータスを記録する処理を含む。ここでは、初期化処理が進行中であることを示すビットがオンされる。以降、ドライバ310-1は、初期化処理の進捗状況を、例えば予め定められたタイミング毎に記録

10

20

30

40

50

する。

【0071】

またドライバ310-1は、ネットワークデバイス13を初期化する(ステップS4)。次にドライバ310-1は、ネットワークデバイス13からの割り込み(HW割り込み)の配信先を自身に設定する(ステップS5)。

【0072】

次にドライバ310-1は、共有メモリ32のゲストOS情報領域323に、自身を含む(自身が動作している)ゲストOS31-1のゲストOS情報を書き込む(ステップS6)。即ちドライバ310-1は、ゲストOS31-1のゲストOS識別子、SW割り込みハンドラ312が使用するSW割り込み番号等を含む、ゲストOS31-1のゲストOS情報を、ゲストOS情報領域323に書き込む。また、ステップS6においてドライバ310-1は、マスターゲストOSがゲストOS31-1であることを示すマスターゲストOS情報をゲストOS情報領域323に登録する。そしてドライバ310-1は、共有メモリ32の特殊処理ステータス領域322に初期化処理が終了したことを記録して(ステップS7)、「open処理」を終了する。

【0073】

一方、ステップS1の判定がYesの場合、ドライバ310-1は、マスターとなっている別のゲストOSのドライバからの要求によって構築された共有メモリ32の特殊処理ステータス領域322を参照して当該別のゲストOSのドライバによる初期化処理(ステップS2~S7)の進捗状況を確認し、当該初期化処理が継続していれば、当該初期化処理が終了するまで待つ(ステップS8)。そして初期化処理が終了すると、ドライバ310-1は、共有メモリ32のゲストOS情報領域323にゲストOS31-1に固有のゲストOS情報(自身のゲストOS情報)を書き込んで(ステップS9)、「open処理」を終了する。

【0074】

[送信処理]

次に、ゲストOSのドライバ(の送信処理部315)がネットワークデバイス13を通じてパケットをネットワーク2に送信する際に実行される「送信処理」の手順について、当該ドライバがゲストOS31-1のドライバ310-1である場合を例に、図6のフローチャートを参照して説明する。この「送信処理」において実行されるネットワークデバイス13へのパケットデータの転送には、DMAが用いられる。

【0075】

今、ゲストOS31-1から当該ゲストOS31-1のドライバ310-1に対して、「送信要求」が与えられたものとする。するとドライバ310-1は、共有メモリ32の送信ディスクリプタチェーン領域325における送信ディスクリプタチェーン325aから、空き先頭ポインタの指し示す送信ディスクリプタを先頭とする、送信に必要な数の送信ディスクリプタ325-j(図3の例では、送信ディスクリプタ325-jは送信ディスクリプタ325-10を含む)を確保する(ステップS11)。

【0076】

次にドライバ310-1は、確保した送信ディスクリプタ325-jのステータス(ステータス部3251に設定されているステータス情報の示すステータス)を「empty」から「not ready」に変更すると共に、確保した送信ディスクリプタ325-jの数だけ空き先頭ポインタを更新する(ステップS12)。つまりドライバ310-1は、確保した送信ディスクリプタ325-jの数だけ空き先頭ポインタの指し示す位置を移動する。

【0077】

次にドライバ310-1は、確保した送信ディスクリプタ325-jを(当該ドライバ310-1を含む)ゲストOS31-1が使用していることを示す情報を、当該ディスクリプタ325-jに追加する。上記ステップS11~S13は、各ゲストOSのドライバ間で排他的に実行される必要がある。この制御には、共有メモリを使用する際に従来から用いられる一般的な排他制御手法が適用可能である。この排他制御については、以降の手順でも同様

10

20

30

40

50

の条件が必要となる箇所があるが、説明を省略する。

【0078】

次にドライバ310-1は、ゲストOS31-1から要求されたパケットの送信のために、次のように送信ディスクリプタ325-jの準備を行う(ステップS14)。まずドライバ310-1は、送信ディスクリプタ325-jのポインタ部3252に格納されているポインタに基づいて共有メモリ32上のDMAバッファ324a内の領域に直接アクセスすることにより、送信されるべきパケットのデータを書き込む。ドライバ310-1は、このデータの書き込み(つまり送信の準備)が完了した送信ディスクリプタ325-jのステータスを「not ready」から「ready」に変更する。

【0079】

本実施形態において、DMA転送はマスターゲストOSが管理する。そこでドライバ310-1は、当該ドライバ310-1を含む(当該ドライバ310-1動作している)ゲストOS31-1がマスターゲストOSであるかを判定する(ステップS15)。

【0080】

もし、ゲストOS31-1がマスターゲストOSでないならば(ステップS15がNo)、ドライバ310-1はマスターゲストOS(のドライバ)にSW割り込みによる通知を行って、送信のためのDMA転送を要求する(ステップS16)。即ちドライバ310-1はマスターゲストOS(のドライバ)に、送信のためのDMA要求の割り込み(送信DMA要求割り込み)を行う。マスターゲストOS(のドライバ)への割り込み配信先は、共有メモリ32のゲストOS情報領域323を参照することによって知ることができる。ドライバ310-1は、ステップS16を実行すると、送信処理を終了する。

【0081】

上述のように、本実施形態では、ゲストOS間の通知にSW割り込みが使用される。しかし、VMM20が提供するゲストOS間の通知手段であれば、SW割り込み以外の手段を用いても構わない。本実施形態では、複数の用途でゲストOS間の通知にSW割り込みが使用される。そこで本実施形態では、割り込み配信先が用途別に設定され、また、割り込み要因を示す仮想の割り込み要因レジスタが共有メモリ32内に配置される構成とすることによって、割り込み要因の種別を対象のゲストOSが認識できるようになっているものとする。

【0082】

一方、ステップS15で、ドライバ310-1を含むゲストOS31-1がマスターゲストOSであると判定された場合、つまり送信要求元のゲストOS31-1がマスターゲストOSである場合(ステップS15がYes)、ドライバ310-1はステップS17に進む。また、ゲストOS31-1がマスターゲストOSであり、且つ当該ゲストOS31-1に、他のゲストOSからの送信要求に従って当該他のゲストOSのドライバから送信DMA要求割り込み(送信要求割り込み)が入った場合にも、ドライバ310-1はステップS17に進む。

【0083】

ステップS17においてドライバ310-1は、ネットワークデバイス13に対するDMA処理を開始する。具体的には、ドライバ310-1は、送信前先頭ポインタで指し示される送信ディスクリプタチェーン325a内の位置からステータスが「ready」である連続する送信ディスクリプタの分だけ、ネットワークデバイス13に対してDMA要求を発行する。ここで、複数のゲストOSが送信ディスクリプタを確保して送信のための準備中である可能性もあるため、送信前先頭ポインタの指し示す位置から空き先頭ポインタの指し示す位置の直前までには、ステータスが「not ready」の送信ディスクリプタとステータスが「ready」の送信ディスクリプタとが混在している可能性がある。ステップS17においてドライバ310-1は、DMA処理が開始された送信ディスクリプタのステータスを「dma」に変更し、送信前先頭ポインタを当該DMA処理が開始された送信ディスクリプタの数の分だけ進める。

【0084】

ドライバ 3 1 0 -1からのDMA要求に応じてネットワークデバイス 1 3によって実行されるDMA転送が終了すると、当該ネットワークデバイス 1 3からドライバ 3 1 0 -1(マスターゲストOS 3 1 -1のドライバ 3 1 0 -1)に、DMA転送終了(DMA終了)がHW割り込みによって通知される。ドライバ 3 1 0 -1は、このDMA終了のHW割り込みに応じて、送信中先頭ポインタの指し示す位置から始まる送信終了した送信ディスクリプタのステータスを「empty」に変更し、変更の対象となった送信ディスクリプタの数の分だけ、当該送信中先頭ポインタを更新する(ステップS 1 8)。このステップS 1 8においてドライバ 3 1 0 -1は、変更の対象となった各送信ディスクリプタの利用ゲストOS情報部に格納されている、当該ディスクリプタを参照(使用)中のゲストOSを示す情報をクリアする。

10

【0085】

ドライバ 3 1 0 -1(マスターゲストOS 3 1 -1のドライバ 3 1 0 -1)は、ステップS 1 8を実行すると、空き先頭ポインタと送信前先頭ポインタとが一致しているかを判定する(ステップS 1 9)。もし、空き先頭ポインタと送信前先頭ポインタとが一致していなければ(ステップS 1 9がNo)、ドライバ 3 1 0 -1は、DMA要求が未発行の送信ディスクリプタがあるはずであると判断する。この場合、ドライバ 3 1 0 -1は、ネットワークデバイス 1 3に対して再びDMA要求を発行するために、ステップS 1 7に戻る。

【0086】

これに対し、空き先頭ポインタと送信前先頭ポインタとが一致しているならば(ステップS 1 9がYes)、ドライバ 3 1 0 -1は、DMA要求が未発行の送信ディスクリプタがないとして送信処理を終了する。

20

【0087】

[受信処理]

次に、ゲストOSのドライバ(の受信処理部 3 1 6)がネットワークデバイス 1 3を通じてパケットを受信する際に実行される「受信処理」の手順について、当該ドライバがゲストOS 3 1 -1のドライバ 3 1 0 -1である場合を例に、図7のフローチャートを参照して説明する。

【0088】

ネットワークデバイス 1 3は、ネットワーク 2から自身が受信すべきパケットを検知すると、空きの受信ディスクリプタ 3 2 6 -jのポインタ部 3 2 5 2で指定されるDMAバッファ 3 2 4 a内の領域へ検知されたパケットをDMA転送する動作(受信DMA動作)を開始する(ステップS 2 1)。DMAバッファ 3 2 4 a内の領域の使用順序は例えばHWを用いた設定によって予め指定することができる。本実施形態では、DMAバッファ 3 2 4 a内の領域の使用順序は、受信ディスクリプタチェーンの順番通りになるよう設定されている。

30

【0089】

さて、ネットワークデバイス 1 3による上述の受信DMA動作(ステップS 2 1)が完了したものとする。即ち、ネットワークデバイス 1 3がネットワーク 2からパケットを受信し、その受信されたパケットを、受信ディスクリプタチェーンに従って空きの受信ディスクリプタ 3 2 6 -jのポインタ部 3 2 5 2で指定されるDMAバッファ 3 2 4 a内の領域へDMA転送する受信DMA動作(ステップS 2 1)が完了したものとする。この場合、ネットワークデバイス 1 3はマスターゲストOSにHW割り込みを発行する(ステップS 2 2)。ここでは、マスターゲストOSはゲストOS 3 1 -1であるものとする。

40

【0090】

マスターゲストOS 3 1 -1のドライバ 3 1 0 -1は、ネットワークデバイス 1 3からのHW割り込みを受け付けると、DMA中先頭ポインタの指し示す受信ディスクリプタチェーンの位置から、受信DMA(受信データを一時格納するためのDMA転送)が完了した受信ディスクリプタ 3 2 6 -jを調べる(ステップS 2 3)。このステップS 2 3においてドライバ 3 1 0 -1は、受信DMAが完了した受信ディスクリプタ 3 2 6 -jのステータスを「empty」から「received」に更新すると共に、DMA中先頭ポインタを更新

50

する。また、ステップ S 2 3 においてドライバ 3 1 0 -1 は、受信 D M A が完了した受信ディスクリプタ 3 2 6 -j の内容（受信ディスクリプタ 3 2 6 -j で示される受信データ）を必要としているゲスト O S を特定する。即ちドライバ 3 1 0 -1 は、ゲスト O S 情報領域 3 2 3 に格納されている各ゲスト O S のゲスト O S 情報の中から、受信パケットに含まれている宛先ネットワークアドレスに一致するネットワークアドレスを含むゲスト O S 情報を選択し、当該選択されたゲスト O S 情報に含まれているゲスト O S 識別子から上述の受信ディスクリプタ 3 2 6 -j の内容を必要とするゲスト O S を特定する。

【 0 0 9 1 】

ドライバ 3 1 0 -1 は、受信 D M A が完了した受信ディスクリプタ 3 2 6 -j をいずれのゲスト O S が参照すべきかを示す情報を、当該受信ディスクリプタ 3 2 6 -j の利用ゲスト O S 情報部に設定する。もし、ネットワークデバイス 1 3 によってネットワーク 2 から受信されたパケットがブロードキャストパケットのような、複数のゲスト O S によって参照されるべきパケットである場合、当該複数のゲスト O S の情報が、該当する受信ディスクリプタ 3 2 6 -j の利用ゲスト O S 情報部に設定される。

【 0 0 9 2 】

ゲスト O S 3 1 -1 のドライバ 3 1 0 -1 は、受信ディスクリプタ 3 2 6 -j の情報を必要とするゲスト O S に対して、受信 D M A の完了を S W 割り込みによって通知する（ステップ S 2 4 ）。ここで、ドライバ 3 1 0 -1 からの S W 割り込みを受けたゲスト O S がゲスト O S 3 1 -2 及び 3 1 -3 であるものとする。

【 0 0 9 3 】

ゲスト O S 3 1 -1 からの S W 割り込みを受けたゲスト O S 3 1 -2 及び 3 1 -3 のそれぞれドライバ 3 1 0 -2 及び 3 1 0 -3 は、自身に割り当てられた受信ディスクリプタ 3 2 6 -j を参照して受信に必要な処理を行う（ステップ S 2 5 ）。この受信に必要な処理は、受信ディスクリプタ 3 2 6 -j のポインタ部 3 2 6 2 によって指し示される D M A バッファ 3 2 4 a 内の領域のデータを、ゲスト O S 3 1 -2 及び 3 1 -3 に固有の記憶領域に転送する処理を含む。ドライバ 3 1 0 -2 及び 3 1 0 -3 は、受信に必要な処理を完了して、受信ディスクリプタ 3 2 6 -j を参照する必要がなくなると、当該ディスクリプタ 3 2 6 -j からそれぞれゲスト O S 3 1 -2 及び 3 1 -3 に関する情報を削除する。

【 0 0 9 4 】

マスターゲスト O S 3 1 -1 のドライバ 3 1 0 -1 は、受信ディスクリプタ 3 2 6 -j から当該ディスクリプタ 3 2 6 -j を参照すべきゲスト O S 3 1 -2 及び 3 1 -3 に関する情報が全て削除されると、当該ディスクリプタ 3 2 6 -j のステータスを「 r e c e i v e d 」から「 e m p t y 」に更新する（ステップ S 2 6 ）。つまりドライバ 3 1 0 -1 は、受信ディスクリプタ 3 2 6 -j を参照する必要のあった全てのゲスト O S にとって、もはや当該ディスクリプタ 3 2 6 -j を参照する必要がなくなると、当該ディスクリプタ 3 2 6 -j のステータスを「 e m p t y 」に更新する。ステップ S 2 6 においてドライバ 3 1 0 -1 は、受信先頭ポインタの指し示す受信ディスクリプタ 3 2 6 -j のステータスが「 e m p t y 」に更新されたならば、当該受信先頭ポインタを 1 受信ディスクリプタ分だけ進める。

【 0 0 9 5 】

[エラー割り込み処理]

次に、ネットワークデバイス 1 3 での動作でエラーが発生した場合に実行される「エラー割り込み処理」の手順について、ゲスト O S 3 1 -1 がマスターゲスト O S である場合を例に、図 8 のフローチャートを参照して説明する。

【 0 0 9 6 】

今、ネットワークデバイス 1 3 での動作でエラーが発生したものとする。この場合、ネットワークデバイス 1 3 はマスターゲスト O S 3 1 -1 のドライバ 3 1 0 -1 にエラー発生を通知するための H W 割り込みを発行する（ステップ S 3 0 ）。

【 0 0 9 7 】

ドライバ 3 1 0 -1 （の H W 割り込みハンドラ 3 1 1 ）は、ネットワークデバイス 1 3 から H W 割り込み（ステップ S 3 0 ）を受けると、共有メモリ 3 2 の物理レジスタマップ領

10

20

30

40

50

域 3 2 1 にマップされた、当該 H W 割り込みの要因を示す割り込み要因レジスタ（例えば当該ネットワークデバイス 1 3 のエラーステータスを示すエラーステータスレジスタ）を参照することによって当該 H W 割り込みの要因となったエラーの種類を識別する（ステップ S 3 1）。このステップ S 3 1 においてドライバ 3 1 0 -1（のエラー割り込み処理部 3 1 7）は、識別されたエラーに対応する処理（エラー処理）が必要ならば、共有メモリ 3 2 の特殊処理ステータス領域 3 2 2 に格納されているビットマップ中の当該エラー処理に対応するビットを、エラー処理が開始されていることを示す状態に設定する。つまりドライバ 3 1 0 -1（マスターゲスト O S 3 1 -1 のドライバ 3 1 0 -1）は特殊処理ステータス領域 3 2 2 に、識別されたエラーに対応する処理（エラー処理）が開始されていることを示すステータス（特殊処理ステータス）を設定する（ステップ S 3 1）。

10

【 0 0 9 8 】

次にマスターゲスト O S 3 1 -1 のドライバ 3 1 0 -1（のエラー割り込み処理部 3 1 7）は、他のゲスト O S 3 1 -2 及び 3 1 -3 のそれぞれドライバ 3 1 0 -2 及び 3 1 0 -3 に対し、エラー処理を開始したことを、S W 割り込みを使用して通知する（ステップ S 3 2）。その後、マスターゲスト O S 3 1 -1 のドライバ 3 1 0 -1 は、エラー処理が終了すると、共有メモリ 3 2 の特殊処理ステータス領域 3 2 2 に格納されているビットマップ中の当該処理に対応するビットの状態を、処理が行われていないことを示す状態に切り替える。つまりドライバ 3 1 0 -1 は、特殊処理ステータス領域 3 2 2 に設定されている、エラー処理が開始されていることを示すステータスを解除する（ステップ S 3 3）。

20

【 0 0 9 9 】

一方、マスターゲスト O S 3 1 -1 以外のゲスト O S 3 1 -2 及び 3 1 -3 のそれぞれドライバ 3 1 0 -2 及び 3 1 0 -3（のエラー割り込み処理部 3 1 7）は、ドライバ 3 1 0 -1 からエラー処理の開始を通知する S W 割り込み（ステップ S 3 2）を受けると、現在実行中の処理を中断または中止するために必要な処理を、当該エラーの種類に応じて行う（ステップ S 4 1）。

【 0 1 0 0 】

次にドライバ 3 1 0 -2 及び 3 1 0 -3 は、それぞれ共有メモリ 3 2 の特殊処理ステータス領域 3 2 2 のステータスを監視し、エラー処理が開始されていることを示すステータスが解除されていることを確認したならば、つまりエラー処理が終了したことを確認したならば、通常の処理に復帰する（ステップ S 4 2）。

30

【 0 1 0 1 】

以上、本実施形態において適用される、ゲスト O S 3 1 -i（ $i = 1, 2, 3$ ）のドライバ 3 1 0 -i（に設けられた o p e n 処理部 3 1 3、送信処理部 3 1 5、受信処理部 3 1 6 及びエラー割り込み処理部 3 1 7）が有する 4 つの基本的な機能、つまり「o p e n 処理機能」、「送信処理機能」、「受信処理機能」及び「エラー割り込み処理機能」について詳述した。これらの機能、特に「送信処理機能」及び「受信処理機能」により、D M A バッファのデータのコピーによるオーバーヘッドを除くことができる。

【 0 1 0 2 】

さて本実施形態では、ドライバ 3 1 0 -i に、「マスター権の委譲処理機能」及び「マスター権の取得処理機能」（をそれぞれ有するマスター権委譲処理部 3 1 8 a 及びマスター権取得処理部 3 1 8 b）が追加されている。以下、「マスター権の委譲処理」及び「マスター権の取得処理」について順次説明する。

40

【 0 1 0 3 】

[マスター権の委譲処理]

「マスター権の委譲処理」とは、マスターゲスト O S が自律的に、マスター権を他のゲスト O S へと移すための処理である。以下、「マスター権の委譲処理」について、ゲスト O S 3 1 -1 がマスターゲスト O S である場合を例に、図 9 のフローチャートを参照して説明する。

【 0 1 0 4 】

マスターゲスト O S 3 1 -1 のドライバ 3 1 0 -1（のマスター権委譲処理部 3 1 8 a）は

50

、マスター権の委譲が必要となった場合、現在ネットワークデバイス 13 を共有している他のゲスト OS の中から、マスター権が委譲されるべきゲスト OS を 1 つ選択する（ステップ S 5 1）。ここでは、ゲスト OS 3 1-2 が選択されたものとする。ステップ S 5 1 において、マスターゲスト OS 3 1-1 のドライバ 3 1 0-1 は、SW 割り込みを使用して、マスターゲスト OS 3 1-1 以外の全てのゲスト OS（ここでは、ゲスト OS 3 1-2 及び 3 2-3）に、ゲスト OS 3 1-1 からゲスト OS 3 1-2 へのマスターゲスト OS 変更を通知する。

【0105】

SW 割り込みによってマスターゲスト OS 変更の通知を受けたゲスト OS のうち、マスター権が委譲されるゲスト OS 3 1-2 のドライバ 3 1 0-2 は、共有メモリ 3 2 のゲスト OS 情報領域 3 2 3 のマスターゲスト OS に関連する情報を当該ゲスト OS 3 1-2 を示すように更新する（ステップ S 5 2）。

10

【0106】

マスターゲスト OS 変更の通知を受けたゲスト OS のうち、マスター権が委譲されるゲスト OS（新マスターゲスト OS）3 1-2 のドライバ 3 1 0-2 は、ネットワークデバイス 13 からの HW 割り込み先がゲスト OS 3 1-2（ドライバ 3 1 0-2）自身になるように設定を変更する（ステップ S 6 1）。もし、マスターゲスト OS 変更を通知したゲスト OS（旧マスターゲスト OS）3 1-1 が特殊処理を実行中であったならば、新マスターゲスト OS 3 1-2 は、共有メモリ 3 2 の特殊処理ステータス領域 3 2 2 を参照して当該特殊処理を途中から引き継ぐ（ステップ S 6 2）。

20

【0107】

[マスター権の取得処理]

「マスター権の取得処理」とは、マスター以外のゲスト OS が自身へマスター権を移すための処理、つまりマスター以外のゲスト OS がマスター権を取得するための処理である。以下、「マスター権の取得処理」について、ゲスト OS 3 1-1 がマスターゲスト OS であり、ゲスト OS 3 1-2 がマスター権を取得しようとするゲスト OS である場合を例に、図 10 のフローチャートを参照して説明する。

【0108】

ゲスト OS 3 1-2 のドライバ 3 1 0-2（のマスター権取得処理部 3 1 8 b）は、当該ゲスト OS 3 1-2 がマスター権を取得しようとする場合、その旨を当該ゲスト OS 3 1-2 以外の、マスターゲスト OS 3 1-1 を含む全てのゲスト OS（ここでは、ゲスト OS 3 1-1 及び 3 2-3）に対して SW 割り込みを使用して通知する（ステップ S 7 1）。この通知により、ゲスト OS 3 1-2 以外のゲスト OS でのマスターゲスト OS に関する処理が抑止され、ゲスト OS 3 1-2 はマスター権を排他的に取得することが可能となる。

30

【0109】

ゲスト OS 3 1-2 のドライバ 3 1 0-2 はステップ S 7 1 を実行すると、共有メモリ 3 2 のゲスト OS 情報領域 3 2 3 のマスターゲスト OS 情報を当該ゲスト OS 3 1-2 を示すように更新する（ステップ S 7 2）。これによりゲスト OS 3 1-2 はマスター権を取得して、新たにマスターゲスト OS（新マスターゲスト OS）となる。

40

【0110】

ゲスト OS 3 1-2 のドライバ 3 1 0-2 は、ネットワークデバイス 13 からの HW 割り込み先がゲスト OS 3 1-2 自身になるように設定を変更する（ステップ S 7 3）。

【0111】

もし、ゲスト OS（新マスターゲスト OS）3 1-2 がマスター権を取得するまでマスターであったゲスト OS（旧マスターゲスト OS）3 1-1 が特殊処理を実行中であったならば、新マスターゲスト OS 3 1-2 は、共有メモリ 3 2 の特殊処理ステータス領域 3 2 2 を参照して当該特殊処理を途中から引き継ぐ（ステップ S 7 4）。

【0112】

[c l o s e 処理]

本実施形態では、ゲスト OS のドライバ（の c l o s e 処理部 3 1 4）がネットワーク

50

デバイス 13 の使用を終了する際に実行される「close 処理 (ドライバの close 処理)」で、上述の「マスター権の委譲処理」が利用される。以下、この「close (クローズ) 処理」の手順について、ゲスト OS 31-1 のドライバ 310-1 がネットワークデバイス 13 の使用を終了する場合を例に、図 11 のフローチャートを参照して説明する。

【0113】

今、ゲスト OS 31-1 から当該ゲスト OS 31-1 のドライバ 310-1 に対して、「close 要求」が与えられたものとする。するとドライバ 310-1 (の close 処理部 314) は、当該ドライバ 310-1 を含むゲスト OS 31-1 がマスターゲスト OS であるかを判定する (ステップ S81)。もし、ゲスト OS 31-1 がマスターゲスト OS でないならば (ステップ S81 が No)、ドライバ 310-1 は後述するステップ S84 に進む。

10

【0114】

一方、ゲスト OS 31-1 がマスターゲスト OS であるならば (ステップ S81 が Yes)、ドライバ 310-1 は、他にネットワークデバイス 13 を共有しているゲスト OS (他のゲスト OS) が存在するかを判定する (ステップ S82)。

【0115】

もし、ステップ S82 の判定が Yes であるならば、即ちゲスト OS 31-1 がマスターゲスト OS で、且つ他にネットワークデバイス 13 を共有しているゲスト OS があるならば、当該ゲスト OS 31-1 のドライバ 310-1 (の close 処理部 314) は (マスター権委譲処理部 318a を用いて) 上述の「マスター権の委譲処理」を行う (ステップ S83)。そしてドライバ 310-1 はステップ S84 に進む。

20

【0116】

ステップ S84 においてドライバ 310-1 (旧マスターゲスト OS 31-1 のドライバ 310-1) は、共有メモリ 32 から当該ドライバ 310-1 を含むゲスト OS (旧マスターゲスト OS) 31-1 に関連する情報を削除する。このステップ S84 での処理は、ドライバ 310-1 が確保・参照していた各ディスクリプタの解放を含む。ドライバ 310-1 はステップ S84 を実行すると、「close 処理」を終了する。

【0117】

一方、ステップ S82 の判定が No であるならば、即ちゲスト OS 31-1 がマスターゲスト OS で、且つ他にネットワークデバイス 13 を共有しているゲスト OS がないならば、当該ゲスト OS 31-1 のドライバ 310-1 は、ネットワークデバイス 13 を停止させるための処理 (ネットワークデバイス終了処理) を行う (ステップ S85)。そしてドライバ 310-1 は、VMM 20 により共有メモリ 32 (共有メモリ領域 120) を解放させて (ステップ S86)、「close 処理」を終了する。

30

【0118】

このような「close 処理」により、たとえマスターゲスト OS のドライバが close されても、つまりマスターゲスト OS がネットワークデバイス 13 の使用を終了しても、他のゲスト OS で運用を続けることが可能となる。

【0119】

< 相互監視に基づいたゲスト OS の排除 >

40

本実施形態のように複数のゲスト OS のドライバ (ネットワークデバイスドライバ) が協調して動作する場合、次のような問題が生じる可能性がある。例えば、あるゲスト OS が「panic」或いは「hang」と呼ばれるような動作不能な状態に陥った場合に、ネットワークデバイス 13 の HW 仕様によってはそのゲスト OS が使用しているディスクリプタの処理が進まなくなると、結果として送受信の処理を進めなくなる可能性がある。また、マスターゲスト OS が動作不能になった場合は、通信自体を行うことができなくなってしまう。

【0120】

そこで本実施形態では、ゲスト OS 31-1 ~ 31-3 が正常に動作しているかを、当該ゲスト OS 31-1 ~ 31-3 のドライバ 310-1 ~ 310-3 が相互に監視する仕組みと、正常

50

に動作していないと判断されたゲストOSをネットワークデバイス13の共有状態から強制的に排除する仕組みとが用意される。この2つの仕組みにより、ゲストOS（例えばマスターゲストOS）が動作不能になっても、残りのゲストOSによって運用が続けられる。

【0121】

正常に動作していないゲストOSを検出する手法は種々考えられる。ここでは以下に示す簡単な検出手法が適用されるものとする。

まず、共有メモリ32（共有メモリ領域120）に、ゲストOS31-1～31-3に共通のカウンタ（カウンタ領域）が用意される。

【0122】

ゲストOS31-1～31-3のドライバ310-1～310-3は、それぞれ特定のタイミングで上記共通のカウンタをカウントアップし、そのカウントアップ後の当該カウンタの値（カウンタ値）を、共有メモリ32のゲストOS情報領域323に格納される自身のゲストOS情報内に保持する。なお、上記特定のタイミングを、ゲストOS内のタイマで管理しても良いし、通信時における特定の処理の機会としても良い。

【0123】

ゲストOS31-1～31-3のドライバ310-1～310-3は、それぞれ他のゲストOSによって当該他のゲストOSのゲストOS情報に保持されているカウンタ値（つまり他のゲストOSのカウンタ値）を監視する。そしてドライバ310-1～310-3は、他のゲストOSのカウンタ値と共通のカウンタの値との差が閾値を超えていたなら、当該他のゲストOSにSW割り込みを配信する。なお、上記の監視のタイミングを、ゲストOS内のタイマで管理しても良いし、通信時における特定の処理の機会としても良い。

【0124】

ドライバ310-1～310-3は、SW割り込みの配信先のゲストOSが当該SW割り込みに反応（例えば上記共通のカウンタのカウントアップ）できなければ、そのゲストOSは正常に動作していない、異常ゲストOSであると判断する。

【0125】

[ゲストOS排除処理]

以下、正常に動作していないと判断された異常ゲストOSを（ネットワークデバイス13を共有する状態から）排除するための「ゲストOS排除処理」の手順について、図12のフローチャートを参照して説明する。

【0126】

今、ゲストOS31-2のドライバ310-2（のゲストOS排除処理部319）が、例えばゲストOS31-1を、正常に動作していない異常ゲストOS（つまり排除されるべきゲストOS）として検出したものとする（ステップS91）。異常ゲストOS31-1を検出したドライバ310-2は、当該異常ゲストOS31-1がマスターゲストOSであるかを当該異常ゲストOS31-1のゲストOS情報に基づいて判定する（ステップS92）。

【0127】

もし、異常ゲストOS31-1がマスターゲストOSでないならば（ステップS92がNo）、当該異常ゲストOS31-1を検出したドライバ310-2は後述するステップS94に進む。これに対して異常ゲストOS31-1がマスターゲストOSであるならば（ステップS92がYes）、当該異常ゲストOS31-1を検出したドライバ310-2が（マスター権取得処理部318bを用いて）上述の「マスター権の取得処理」を行う（ステップS93）。ドライバ310-2はステップS93を実行するとステップS94に進む。

【0128】

ステップS94においてドライバ310-2は、排除されるべき異常ゲストOS31-1に、その旨をSW割り込みを使用して通知する（ステップS94）。この通知は、異常ゲストOS31-1（正常に動作していないと判断された異常ゲストOS）が動作不能状態から復帰した場合に、ネットワークデバイス13を共有する状態から排除されているにもかかわらず自身が当該ネットワークデバイス13を使用しているものとして誤動作しないよ

10

20

30

40

50

うに知らせるために行われる。

【 0 1 2 9 】

次にドライバ 3 1 0 -2は、共有メモリ 3 2 から異常ゲスト OS 3 1 -1に関連する情報を削除する（ステップ S 9 5）。このステップ S 9 5での処理は、異常ゲスト OS 3 1 -1のドライバ 3 1 0 -1が確保・参照していた各ディスクリプタの解放を含む。ドライバ 3 1 0 -2は、ステップ S 9 5の実行により、異常ゲスト OS 3 1 -1をネットワークデバイス 1 3を共有する状態から強制的に排除すると、「ゲスト OS 排除処理」を終了する。

【 0 1 3 0 】

さて、排除されたゲスト OS 3 1 -1のドライバ 3 1 0 -1が復帰したものとする。そして、復帰したドライバ 3 1 0 -1が、ドライバ 3 1 0 -2からの SW 割り込み（ステップ S 9 4）によって当該ドライバ 3 1 0 -1を含むゲスト OS 3 1 -1が排除されたことが認識できたものとする。この場合、ドライバ 3 1 0 -1は、ゲスト OS 3 1 -1の排除に対応した後処理、例えばゲスト OS 3 1 -1をエラー扱いとする処理を行う（ステップ S 1 0 0）。

【 0 1 3 1 】

< 動的なマスター権の移行 >

上述の説明では、マスター権の移行は、マスターゲスト OS の終了や、マスターゲスト OS に異常が発生した際に特別に行われる。しかし、マスターゲスト OS が正常稼働中でも、動作効率の向上のために、動的にマスター権を移行することも可能である。このような動作効率向上のためのマスター権移行の例を以下に列挙する。

【 0 1 3 2 】

（ 1 ）最も通信回数の多いゲスト OS がマスターとなる。この場合、通信の際のマスターゲスト OS とその他のゲスト OS との間の通知やディスパッチによるオーバーヘッドを軽減することができる。通信回数に関しては、共有メモリ 3 2 -1～ 3 2 1-2のゲスト OS 情報領域 3 2 3 に格納されるゲスト OS 情報に含まれている通信回数統計情報を利用すれば良い。

【 0 1 3 3 】

（ 2 ）ディスパッチ機会（時間）の最も多いゲスト OS （つまりディスパッチ機会の最も多いゲスト OS ）がマスターとなる。各ゲスト OS のディスパッチ機会（時間）の情報を取得するには、当該各ゲスト OS のディスパッチポリシ（スケジュール）を当該ゲスト OS 自身が把握すれば良い。本実施形態において VMM 2 0 は、ゲスト OS 3 1 -1～ 3 1 -3のディスパッチポリシを提供する手段を有している。ゲスト OS 3 1 -1～ 3 1 -3のそれぞれドライバ 3 1 0 -1～ 3 1 0 -3（のマスター権移行処理部 3 1 8）は、この手段を利用して、ゲスト OS 3 1 -1～ 3 1 -3のディスパッチポリシを把握する。ディスパッチポリシを提供する手段は、例えばゲスト OS 3 1 -1～ 3 1 -3に CPU 1 1（CPU 時間）を割り当てるための、VMM 2 0 に設けられたスケジューラまたはディスパッチャに含まれているものとする。ここでは、ディスパッチポリシを提供する手段は、ゲスト OS 3 1 -1～ 3 1 -3のディスパッチ時間の比率を指定するディスパッチ指定手段である。

【 0 1 3 4 】

（ 3 ）各ゲスト OS に予め優先度（ゲスト OS 優先度）を付与し、稼働中のゲスト OS の中で最も優先度が高いゲスト OS がマスターとなる。各ゲスト OS の優先度は、例えば、システムの構築者のようなユーザの指定によって、ネットワーク 2 に接続されたクライアント端末から付与されるものとする。各ゲスト OS の優先度は、前述のように、共有メモリ 3 2 -1～ 3 2 1-2のゲスト OS 情報領域 3 2 3 に格納されるゲスト OS 情報に含まれている。

【 0 1 3 5 】

[通信回数の多いゲスト OS へのマスター権の移行処理]

次に上記（ 1 ）に対応した「通信回数の多いゲスト OS へのマスター権の移行処理」について説明する。

【 0 1 3 6 】

本実施形態では、マスターゲスト OS 以外のゲスト OS が通信を行う度に、マスターゲ

10

20

30

40

50

ストOS自身が通信を行う場合と比較して、マスターゲストOSへの割り込み通知とディスパッチによるオーバーヘッドが発生する。したがって、より通信回数の多いゲストOSがマスターゲストOSとなる方が、システム全体のオーバーヘッドが少なくなり、通信のスループットが良くなることが期待される。

【0137】

本実施形態では、通信回数の判断に、各ゲストOSが使用したディスクリプタの数を使用される。ここでは通信回数の判断に、送信ディスクリプタチェーン325a（における送信ディスクリプタ325-j）が使用されるものとする。しかし通信回数の判断に、受信ディスクリプタチェーン（における受信ディスクリプタ326-j）を用いても良い。また、通信回数をより高精度に判断するために、両ディスクリプタチェーン（におけるディスクリプタ325-j及び326-j）を用いても構わない。

10

【0138】

マスター権を移行するかを判定するのに、「最も多くディスクリプタを使用したゲストOSをマスターゲストOSとする」ような簡単な手法が考えられる。この手法を適用した場合、マスター権の移行処理が頻発する可能性がある。また、マスター権の移行処理自体がオーバーヘッドとなる可能性もある。

【0139】

そこで本実施形態では、マスター権の移行処理が頻発するのを防ぐために、以下に挙げる2つの条件に基づきマスター権を移行するかが判定される。

1a) マスターゲストOS以外のある1つのゲストOSの通信回数が著しく多い場合：ディスクリプタチェーンに含まれるディスクリプタ数の一定割合以上（例えば1/2以上）を、ある1つのゲストOSが使用していた場合が、これに該当する。

20

【0140】

1b) マスターゲストOSの通信回数が著しく少ない場合：

N個のゲストOSがネットワークデバイス13を共有しているものとする、ディスクリプタチェーンに含まれるディスクリプタの数の一定割合未満、例えば1/(2N)個未満しか、マスターゲストOSが使用していなかった場合が、これに該当する。

【0141】

以下、上述の2つの条件に基づいて行われる「通信回数の多いゲストOSへのマスター権の移行処理」の手順について、ゲストOS31-1がマスターゲストOSである場合を例に、図13のフローチャートを参照して説明する。

30

【0142】

マスターゲストOS31-1のドライバ310-1（のマスター権移行処理部318）は、当該マスターゲストOS31-1からの送信要求、または他のゲストOS（ゲストOS31-2または31-3）からのDMA要求割り込みに応じて送信処理が行われる際に、共有メモリ32内にある該当するゲストOS（送信要求元またはDMA要求割り込み元のゲストOS）の通信回数の統計情報を更新する（ステップS111）。

【0143】

次にマスターゲストOS31-1のドライバ310-1は、通信回数の統計情報によって示されるゲストOS31-1～31-3の通信回数の合計zが、予め定められた一定個数、例えば1200個に達したかを判定する（ステップS112）。もし、上記zが1200個に達したならば（ステップS112がYes）、ドライバ310-1は通信回数の調査を次のように行う。なお、調査を行うタイミングが、例えば一定時間毎のように、zの数以外の別のタイミングであっても構わない。

40

【0144】

まずマスターゲストOS31-1のドライバ310-1は、マスターゲストOS31-1以外のゲストOS、つまりゲストOS31-2及び31-3の通信回数の統計情報を調査し、当該ゲストOS31-2及び31-3の中に、ディスクリプタ使用数（通信回数）bが上記一定個数（1200個）の一定割合以上、例えば1/2以上（つまり600個以上）のゲストOSが存在するかを判定する（ステップS113）。もし、bが600個以上のゲストOS

50

が見つければ（ステップS 1 1 3がY e s）、ドライバ3 1 0 -1（のマスター権移行処理部3 1 8に含まれているマスター権委譲処理部3 1 8 a）は、そのゲストOSにマスター権を委譲するための前述の「マスター権委譲処理」を行う（ステップS 1 1 4）。そしてドライバ3 1 0 -1は、後述するステップS 1 1 6に進む。

【0 1 4 5】

一方、bが6 0 0個以上のゲストOSが見つからないならば（ステップS 1 1 3がN o）、マスターゲストOS 3 1 -1のドライバ3 1 0 -1は、当該マスターゲストOS 3 1 -1の通信回数の統計情報を調査し、当該マスターゲストOS 3 1 -1のディスクリプタ使用数（通信回数）aが上記一定個数（1 2 0 0個）の一定割合未満、例えば1 / (2 N) 未満（Nが3の場合、2 0 0個未満）であるかを判定する（ステップS 1 1 5）。

10

【0 1 4 6】

もし、マスターゲストOS 3 1 -1のディスクリプタ使用数が1 2 0 0の1 / (2 N) 未満であれば（ステップS 1 1 5がY e s）、当該マスターゲストOS 3 1 -1のドライバ3 1 0 -1はステップS 1 1 4に進む。このステップS 1 1 4においてマスターゲストOS 3 1 -1のドライバ3 1 0 -1（のマスター権移行処理部3 1 8に含まれているマスター権委譲処理部3 1 8 a）は、当該ゲストOS 3 1 -1を除くゲストOS 3 1 -2及び3 1 -3のうち、最もディスクリプタ使用数の多いゲストOSにマスター権を委譲するための前述の「マスター権委譲処理」を行う（ステップS 1 1 4）。そしてドライバ3 1 0 -1は、ステップS 1 1 6に進む。

20

【0 1 4 7】

ステップS 1 1 6においてドライバ3 1 0 -1は、共有メモリ3 2 内にあるゲストOS 3 1 -1～3 1 -3の通信回数（ディスクリプタ使用数）の統計情報をクリアする。

【0 1 4 8】

次に、上記「通信回数の多いゲストOSへのマスター権の移行処理」の具体例（例1乃至3）について説明する。ここでは、ゲストOS 3 1 -1、ゲストOS 3 1 -2及びゲストOS 3 1 -3を、それぞれゲストOS # 1、ゲストOS # 2 及びゲストOS # 3と表現する。また、ステップS 1 1 2が実行される時点において、ゲストOS # 1がマスターゲストOSであり、送信ディスクリプタチェーン3 2 5 aのディスクリプタ数の合計が1 2 0 0であるものとする。

30

【0 1 4 9】

（例1）

ゲストOS # 1のディスクリプタ使用数：4 0 0

ゲストOS # 2のディスクリプタ使用数：7 0 0

ゲストOS # 3のディスクリプタ使用数：1 0 0

例1では、ステップS 1 1 3の判定条件にゲストOS # 2が合致する。この場合、ゲストOS # 1からゲストOS # 2にマスター権が移行される。

【0 1 5 0】

（例2）

ゲストOS # 1のディスクリプタ使用数：1 0 0

ゲストOS # 2のディスクリプタ使用数：5 0 0

ゲストOS # 3のディスクリプタ使用数：6 0 0

例2では、ステップS 1 1 3の判定条件に合致するゲストOSは存在しない。ここでは、マスターゲストOS # 1のディスクリプタ使用数は、1 2 0 0 / (3 × 2)、つまり2 0 0よりも少ない。したがって、マスターゲストOS # 1のディスクリプタ使用数（通信数）は、ステップS 1 1 5の判定条件に合致する。この場合、残りのゲストOSのうち、最もディスクリプタ使用数の多いゲストOS # 3にマスター権が移行される。

40

【0 1 5 1】

（例3）

ゲストOS # 1のディスクリプタ使用数：3 0 0

ゲストOS # 2のディスクリプタ使用数：5 0 0

50

ゲストOS # 3 のディスクリプタ使用数 : 4 0 0

例 3 では、ステップ S 1 1 3 の判定条件に合致するゲストOS は存在しない。また、マスターゲストOS # 1 のディスクリプタ使用数は、ステップ S 1 1 5 の判定条件に合致しない。この場合、ゲストOS # 2 のディスクリプタ使用数は最も多いものの、当該ゲストOS # 2 へのマスター権の移行は行われない。これにより、マスター権の移行処理が頻発するのが防止される。

【 0 1 5 2 】

[ディスパッチ機会の多いゲストOS へマスター権の移行処理]

次に、上記 (2) に対応した「ディスパッチ機会の多いゲストOS へのマスター権の移行処理」について説明する。

10

【 0 1 5 3 】

ディスパッチされる機会 (または時間) が最も多いゲストOS がマスター権を持つと、その分だけネットワークデバイス 1 3 へのアクセスや特殊処理を行う機会が増える。このため、結果としてスループットが向上する可能性がある。

【 0 1 5 4 】

そこで本実施形態では、ゲストOS 3 1 -1 ~ 3 1 -3 のそれぞれドライバ 3 1 0 -1 ~ 3 1 0 -3 は、VMM 2 0 が有するディスパッチ指定手段を利用して、ゲストOS 3 1 -1 ~ 3 1 -3 のディスパッチ時間の比率を把握する。

【 0 1 5 5 】

ゲストOS 3 1 -1 ~ 3 1 -3 のドライバ 3 1 0 -1 ~ 3 1 0 -3 は、当該ゲストOS 3 1 -1 ~ 3 1 -3 のディスパッチ時間の比率に基づき、当該ドライバの o p e n 時における「マスター権の取得処理」及び当該ドライバの c l o s e 時における「マスター権の委譲処理」を次のように実行する。なお、マスターゲストOS はゲストOS 3 1 -1 であるものとする。またマスターゲストOS 以外のゲストOS 3 1 -2 またはゲストOS 3 1 -3 をゲストOS 3 1 -k (k は 2 または 3) と表現するものとする。

20

【 0 1 5 6 】

(ドライバの o p e n 時)

ゲストOS 3 1 -k のドライバ 3 1 0 -k の o p e n 時において、マスターゲストOS 3 1 -1 よりも当該ゲストOS 3 1 -k のディスパッチ時間の比率の設定が大きいものとする。この場合、ゲストOS 3 1 -k のドライバ 3 1 0 -k は、マスター権の取得処理を行う。

30

【 0 1 5 7 】

(ドライバの c l o s e 時)

マスターゲストOS 3 1 -1 のドライバ 3 1 0 -1 の c l o s e 時において、他のゲストOS の中でゲストOS 3 1 -k のディスパッチ時間の比率の設定が最も大きいものとする。この場合、マスターゲストOS 3 1 -1 のドライバ 3 1 0 -1 は、自身の o p e n 時に、ゲストOS 3 1 -k へマスター権を委譲する。

なお、ディスパッチ時間に代えて、ディスパッチ機会 (回数) を用いても良い。

【 0 1 5 8 】

[優先度の高いゲストOS へのマスター権の移行処理]

次に上記 (3) に対応した「優先度の高いゲストOS へのマスター権の移行処理」について説明する。

40

【 0 1 5 9 】

前述のように、ゲストOS 3 1 -1 ~ 3 1 -3 の優先度 (ゲストOS 優先度) は、VMM 2 0 が、ディスパッチされるべきゲストOS を選択する際に選定の基準とする値である。このため優先度が高いゲストOS ほど、当該ゲストOS へのディスパッチの機会が多く、CPU 割当時間が多くなると想定される。したがって、ゲストOS 3 1 -1 ~ 3 1 -3 の優先度は、当該ゲストOS 3 1 -1 ~ 3 1 -3 のディスパッチのされやすさの指標 (つまりディスパッチされる時間もしくは回数の指標) であるといえる。このため、ゲストOS 3 1 -1 ~ 3 1 -3 のディスパッチ時間の比率に代えて、当該ゲストOS 3 1 -1 ~ 3 1 -3 の優先度を用いることも可能である。

50

【0160】

そこでゲストOS 31-1~31-3のドライバ310-1~310-3は、共有メモリ32のゲストOS情報領域323に格納されているゲストOS情報に含まれているゲストOS優先度に基づいて、当該ドライバのopen時における「マスター権の取得処理」及び当該ドライバのclose時における「マスター権の委譲処理」を次のように実行する。なお、マスターゲストOSはゲストOS 31-1であるものとする。またマスターゲストOS以外のゲストOS 31-2またはゲストOS 31-3をゲストOS 31-k (kは2または3)と表現するものとする。

【0161】

(ドライバのopen時)

ゲストOS 31-kのドライバ310-kのopen時において、マスターゲストOS 31-1よりも当該ゲストOS 31-kの優先度が高いものとする。この場合、ゲストOS 31-kのドライバ310-kは、マスター権の取得処理を行う。

【0162】

(ドライバのclose時)

マスターゲストOS 31-1のドライバ310-1のclose時において、他のゲストOSの中でゲストOS 31-kの優先度が最も高いものとする。この場合、マスターゲストOS 31-1のドライバ310-1は、自身のopen時に、ゲストOS 31-kへマスター権を委譲する。

【0163】

[ゲストOSの優先度の変更時におけるマスター権の移行処理]

ゲストOS 31-1~31-3の優先度は、例えばVMM 20によって変更され、その変更がVMM 20からゲストOS 31-1~31-3に通知されるものとする。またゲストOS 31-1~31-3の優先度は、当該ゲストOS 31-1~31-3からVMM 20への要求によっても変更されるものとする。

【0164】

以下、「ゲストOSの優先度の変更時における(優先度の高いゲストOSへの)マスター権の移行処理」について、ゲストOS 31-1の優先度に変更された場合を例に、図14のフローチャートを参照して説明する。

【0165】

今、ゲストOS 31-1のドライバ310-1(のマスター権移行処理部318)が、当該ゲストOS 31-1の優先度(ゲストOS優先度)がVMM 20によって変更されたことを、当該VMM 20からの通知によって検出したものとする。するとゲストOS 31-1のドライバ310-1は、共有メモリ32のゲストOS情報領域323に格納されているゲストOS情報に含まれている当該ゲストOS 31-1のゲストOS優先度を当該変更後の優先度に更新する(ステップS121)。なお、ゲストOS 31-1のドライバ310-1が上記ゲストOS情報領域323に格納されているゲストOS情報を例えば定期的に参照することにより、当該ゲストOS 31-1のゲストOS優先度に変更されたことを検出しても良い。

【0166】

次にゲストOS 31-1のドライバ310-1は、当該ゲストOS 31-1がマスターゲストOSであるかを判定する(ステップS122)。もし、ゲストOS 31-1がマスターゲストOSであるならば(ステップS122がYes)、当該ゲストOS 31-1のドライバ310-1は、当該ゲストOS 31-1の優先度が下がったかを判定する(ステップS123)。

【0167】

もし、ゲストOS 31-1の優先度が下がってないならば(ステップS123がNo)、当該ゲストOS 31-1のドライバ310-1は何もせずにマスター権の移行処理を終了する。

【0168】

これに対し、ゲストOS 31-1の優先度が下がったならば(ステップS123がYes

10

20

30

40

50

）、当該ゲストOS 31-1のドライバ310-1は、共有メモリ32のゲストOS情報領域323に格納されているゲストOS情報に基づき、当該ゲストOS 31-1（つまりマスターゲストOS）よりも優先度の高いゲストOSが存在するかを判定する（ステップS124）

もし、ゲストOS 31-1（マスターゲストOS）よりも優先度の高いゲストOSが存在しないならば（ステップS124がNo）、当該ゲストOS 31-1のドライバ310-1は何もせずにマスター権の移行処理を終了する。

【0169】

これに対し、ゲストOS 31-1（マスターゲストOS）よりも優先度の高いゲストOSが存在するならば（ステップS124がYes）、当該ゲストOS 31-1のドライバ310-1（のマスター権移行処理部318に含まれているマスター権委譲処理部318a）は、当該ゲストOS 31-1（マスターゲストOS）よりも優先度の高いゲストOSのうち最も優先度の高いゲストOSにマスター権を委譲して（ステップS125）、マスター権の移行処理を終了する。即ちゲストOS 31-1のドライバ310-1は、当該ゲストOS 31-1がマスターゲストOSであり、且つ当該ゲストOS 31-1の優先度が下がって他のゲストOSの優先度のほうが高くなっていれば、「マスター権委譲処理」を行う。

【0170】

一方、ゲストOS 31-1がマスターゲストOSでない場合（ステップS122がNo）にも、ドライバ310-1は、当該ゲストOS 31-1の優先度が下がったかを判定する（ステップS126）。

【0171】

もし、ゲストOS 31-1の優先度が下がったならば（ステップS126がYes）、当該ゲストOS 31-1のドライバ310-1は何もせずにマスター権の移行処理を終了する。

【0172】

これに対し、ゲストOS 31-1の優先度が下がっていないならば（ステップS126がNo）、当該ゲストOS 31-1のドライバ310-1は、当該ゲストOS 31-1の優先度が上がってマスターゲストOSの優先度よりも高くなったかを判定する（ステップS127）。

【0173】

もし、ゲストOS 31-1の優先度がマスターゲストOSの優先度よりも依然として低いならば（ステップS127のNo）、当該ゲストOS 31-1のドライバ310-1は何もせずにマスター権の移行処理を終了する。

【0174】

これに対し、ゲストOS 31-1の優先度がマスターゲストOSの優先度よりも高くなったならば（ステップS127のYes）、当該ゲストOS 31-1のドライバ310-1（のマスター権移行処理部318に含まれているマスター権取得処理部318b）はマスター権を取得して（ステップS128）、マスター権の移行処理を終了する。即ちゲストOS 31-1のドライバ310-1は、当該ゲストOS 31-1がマスターゲストOSでなく、且つ当該ゲストOS 31-1の優先度が上がってマスターゲストOSの優先度よりも高くなっていれば、「マスター権取得処理」を行う。

【0175】

なお、本発明は、上記実施形態そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化できる。また、上記実施形態に開示されている複数の構成要素の適宜な組み合わせにより種々の発明を形成できる。例えば、実施形態に示される全構成要素から幾つかの構成要素を削除してもよい。

【図面の簡単な説明】

【0176】

【図1】本発明の一実施形態に係る仮想計算機システムの構成を示すブロック図。

【図2】図1に示される共有メモリにおける領域の割り当て例を示す図。

【図3】図2に示される送信ディスクリプタチェーン領域に格納される送信ディスクリプ

10

20

30

40

50

タチェーンの一例を示す図。

【図４】図２に示される受信ディスクリプタチェーン領域に格納される受信ディスクリプタチェーンの一例を示す図。

【図５】同実施形態における「open処理」の手順を示すフローチャート。

【図６】同実施形態における「送信処理」の手順を示すフローチャート。

【図７】同実施形態における「受信処理」の手順を示すフローチャート。

【図８】同実施形態における「エラー割り込み処理」の手順を示すフローチャート。

【図９】同実施形態における「マスター権の委譲処理」の手順を示すフローチャート。

【図１０】同実施形態における「マスター権の取得処理」の手順を示すフローチャート。

【図１１】同実施形態における「close処理」の手順を示すフローチャート。

【図１２】同実施形態における「ゲストOS排除処理」の手順を示すフローチャート。

【図１３】同実施形態における「通信回数の多いゲストOSへのマスター権の移行処理」の手順を示すフローチャート。

【図１４】同実施形態における「ゲストOSの優先度の変更時におけるマスター権の移行処理」の手順を示すフローチャート。

【符号の説明】

【０１７７】

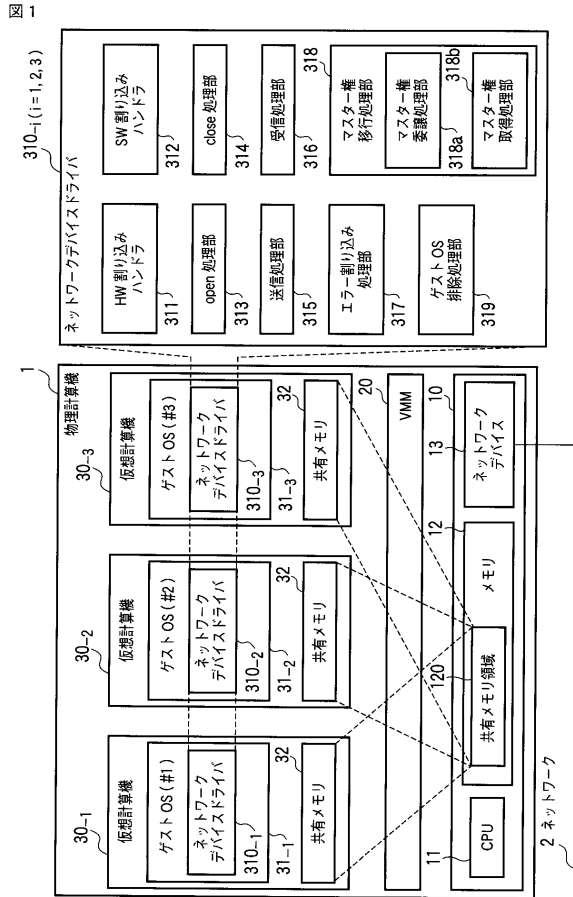
１…物理計算機、２…ネットワーク、１０…HW（ハードウェア）、１１…CPU、１２…メモリ、１３…ネットワークデバイス、２０…VMM（仮想計算機マネージャ）、３０-１～３０-３…仮想計算機、３１-１～３１-３…ゲストOS、３２…共有メモリ、１２０…共有メモリ領域、３１０-１～３１０-３…ネットワークデバイスドライバ、３１１…HW割り込みハンドラ、３１２…SW割り込みハンドラ、３１３…open処理部、３１４…close処理部、３１５…送信処理部、３１６…受信処理部、３１７…エラー割り込み処理部、３１８…マスター権移行処理部、３１８a…マスター権委譲処理部、３１８b…マスター権取得処理部、３１９…ゲストOS排除処理部、３２１…物理レジスタマップ領域、３２２…特殊処理ステータス領域、３２３…ゲストOS情報領域、３２４…DMAバッファ領域、３２４a…DMAバッファ、３２５…送信ディスクリプタチェーン領域、３２５a…送信ディスクリプタチェーン、３２５-１～３２５-１０…送信ディスクリプタ、３２６…受信ディスクリプタチェーン領域、３２６a…受信ディスクリプタチェーン、３２６-１～３２６-１０…受信ディスクリプタ、３２７…他共有情報領域。

10

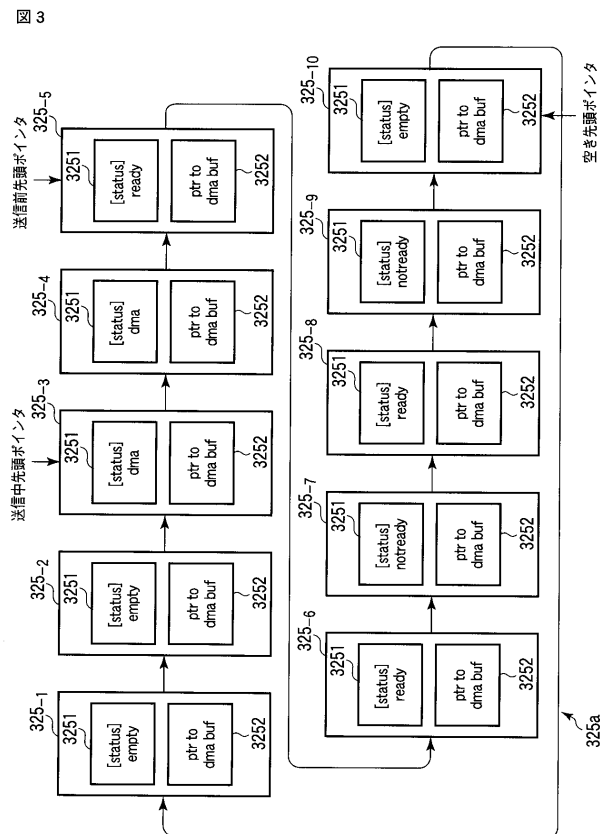
20

30

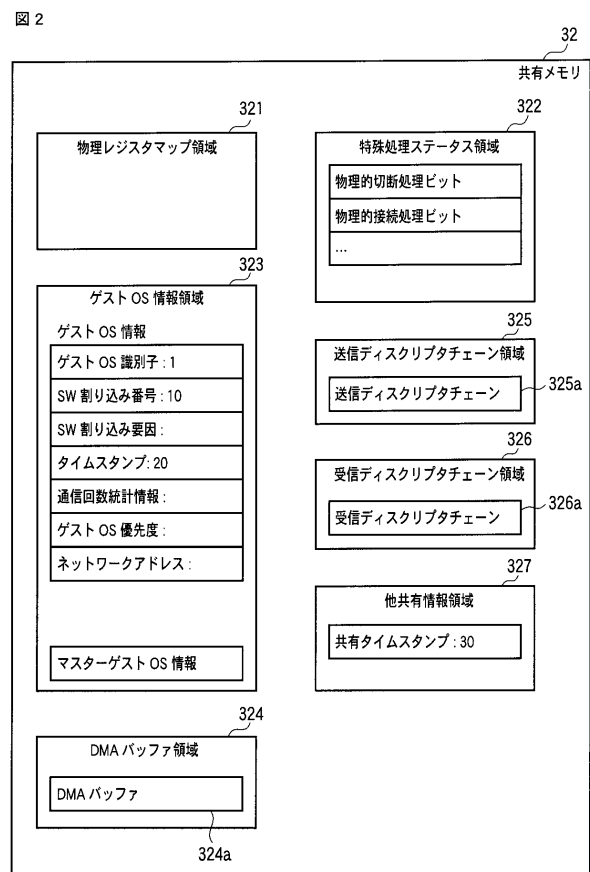
【 図 1 】



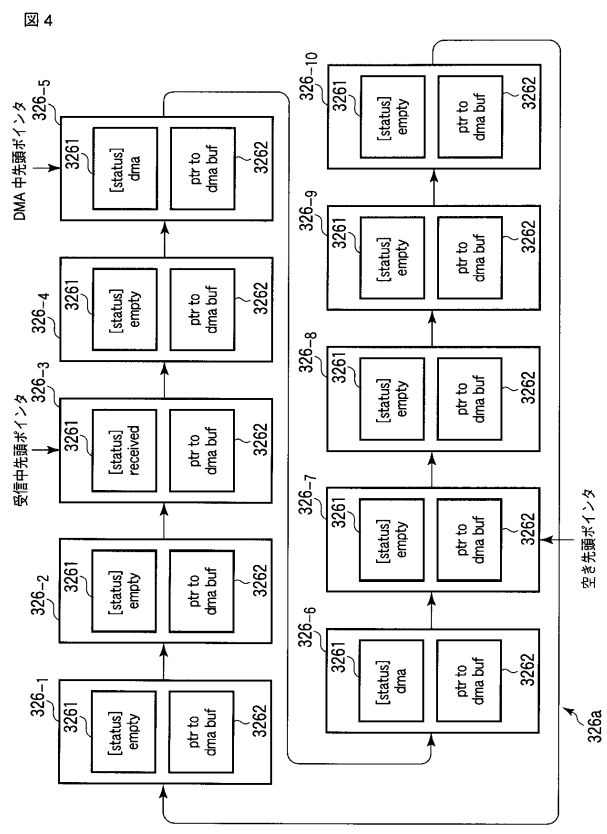
【 図 3 】



【 図 2 】

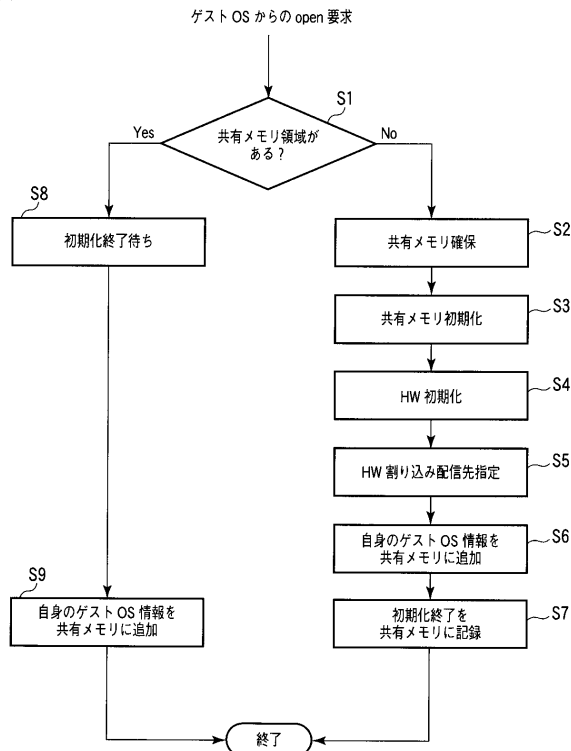


【圖 4】



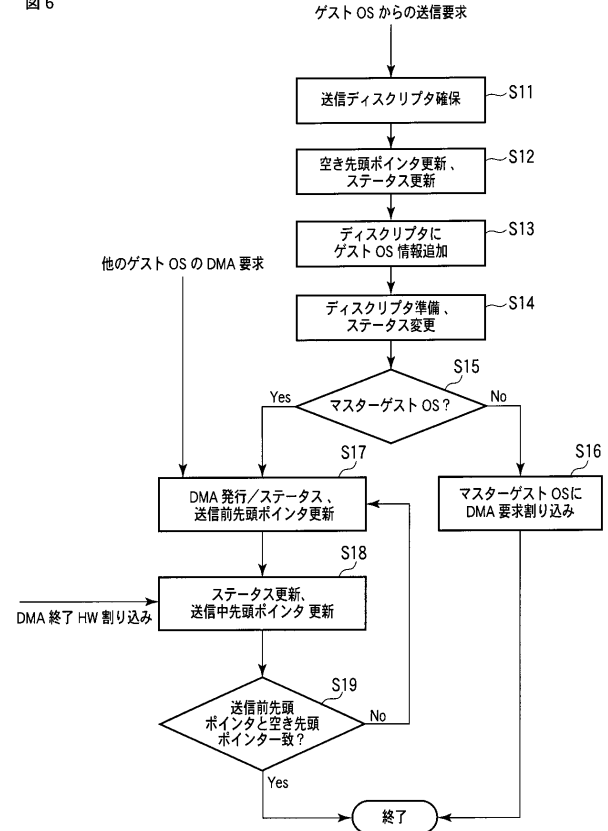
【図 5】

図 5



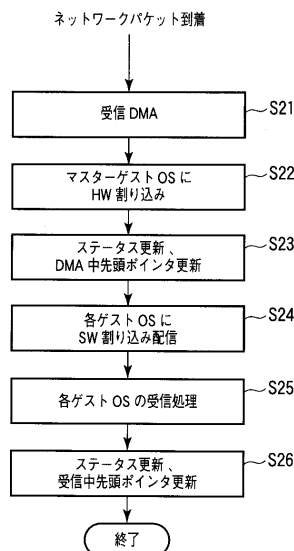
【図 6】

図 6



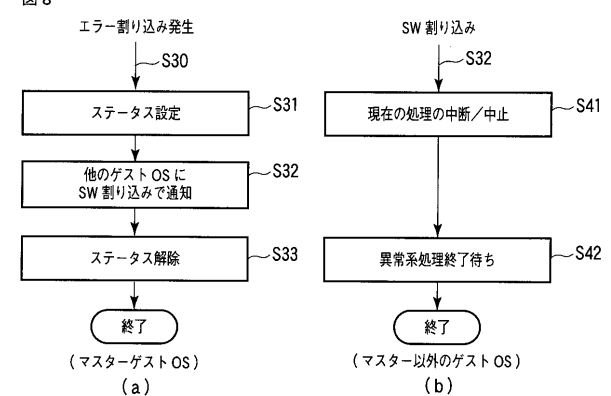
【図 7】

図 7



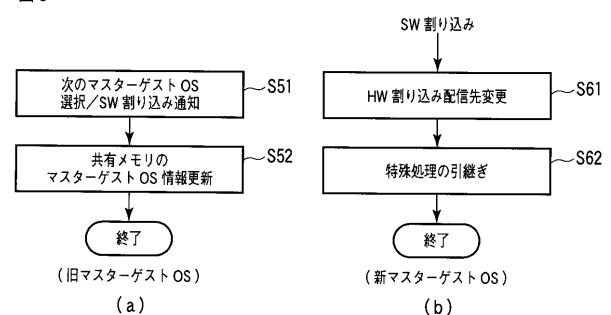
【図 8】

図 8



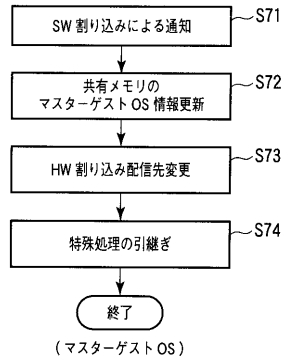
【図 9】

図 9



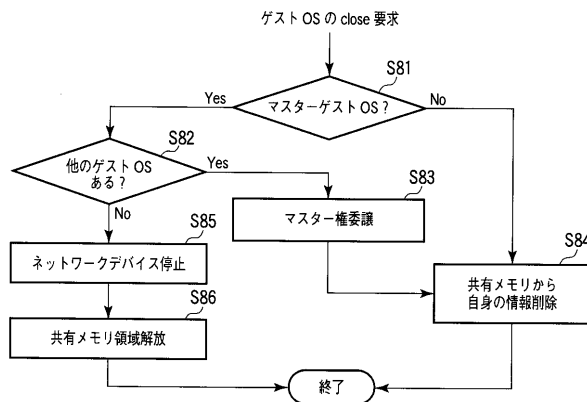
【図 10】

図 10



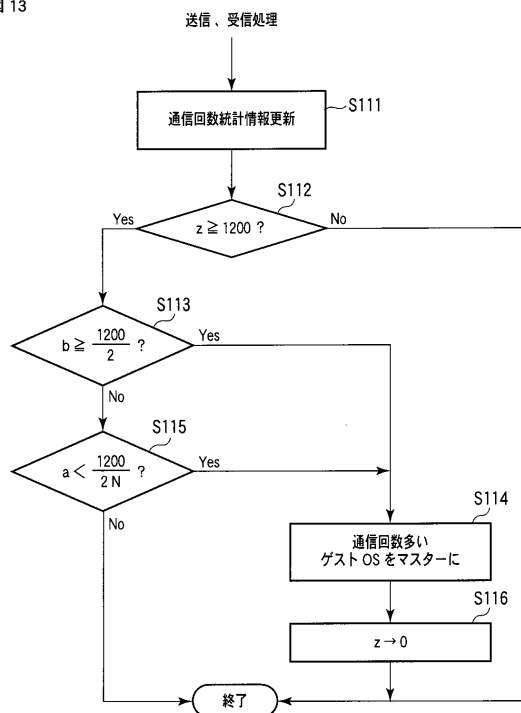
【図 11】

図 11



【図 13】

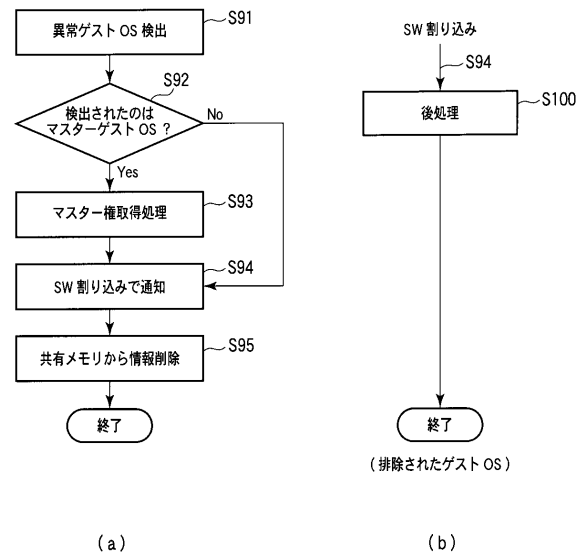
図 13



z : 各ゲスト OS の通信回数の合計
 N : ネットワークデバイスを共有するゲスト OS の数
 a : マスターゲスト OS の通信回数
 b : その他ゲスト OS の通信回数のうちの最大値

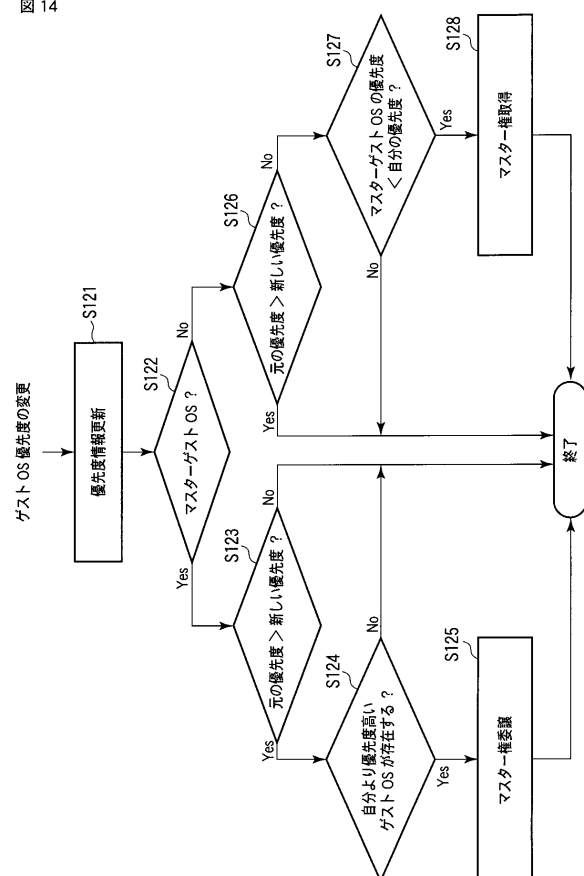
【図 12】

図 12



【図 14】

図 14



フロントページの続き

(74)代理人 100109830
弁理士 福原 淑弘

(74)代理人 100075672
弁理士 峰 隆司

(74)代理人 100095441
弁理士 白根 俊郎

(74)代理人 100084618
弁理士 村松 貞男

(74)代理人 100103034
弁理士 野河 信久

(74)代理人 100119976
弁理士 幸長 保次郎

(74)代理人 100153051
弁理士 河野 直樹

(74)代理人 100140176
弁理士 砂川 克

(74)代理人 100100952
弁理士 風間 鉄也

(74)代理人 100101812
弁理士 勝村 紘

(74)代理人 100070437
弁理士 河井 将次

(74)代理人 100124394
弁理士 佐藤 立志

(74)代理人 100112807
弁理士 岡田 貴志

(74)代理人 100111073
弁理士 堀内 美保子

(74)代理人 100134290
弁理士 竹内 将訓

(74)代理人 100127144
弁理士 市原 卓三

(74)代理人 100141933
弁理士 山下 元

(72)発明者 高山 雅陽
東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内