

(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2012-128807
(P2012-128807A)

(43) 公開日 平成24年7月5日(2012.7.5)

(51) Int.Cl.

G06F 9/46 (2006.01)

F I

G06F 9/46 350

テーマコード (参考)

審査請求 未請求 請求項の数 8 O L (全 66 頁)

(21) 出願番号 特願2010-282158 (P2010-282158)
(22) 出願日 平成22年12月17日 (2010.12.17)

(71) 出願人 000005223
富士通株式会社
神奈川県川崎市中原区上小田中4丁目1番1号
(74) 代理人 100074099
弁理士 大菅 義之
(74) 代理人 100133570
弁理士 ▲徳▼永 民雄
(72) 発明者 三吉 貴史
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

(54) 【発明の名称】 情報処理装置

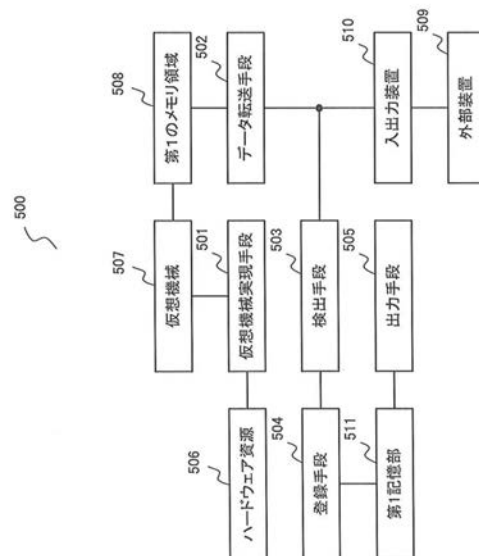
(57) 【要約】 (修正有)

【課題】DMAライトによるゲスト物理メモリの変更内容の保持に必要なハードウェア資源を少なくする。

【解決手段】1以上の仮想機械507を実現する仮想機械実現手段501と、前記仮想機械に割り当てられる第1のメモリ領域508のアドレスと、前記第1のメモリ領域の実メモリである第2のメモリ領域のアドレスと、を相互に変換することにより、入出力装置510から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送を行なうデータ転送手段502と、前記入出力装置から、前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送されるデータを検出する検出手段503と、検出したデータが一定の条件を満たす場合に、該検出したデータにより変更される前記第1のメモリ領域に関する更新情報を第1記憶部511に記憶する登録手段504と、前記第1記憶部に記憶される更新情報を出力する出力手段505と、を備える。

【選択図】 図5

情報処理装置を説明する図



【特許請求の範囲】**【請求項 1】**

ハードウェア資源を管理することにより、1以上の仮想機械を実現する仮想機械実現手段と、

前記仮想機械に割り当てられる第1のメモリ領域のアドレスと、前記第1のメモリ領域の実メモリである第2のメモリ領域のアドレスと、を相互に変換することにより、外部装置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送を行なうデータ転送手段と、

前記入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送されるデータを検出する検出手段と、

検出したデータが一定の条件を満たす場合に、該検出したデータにより変更される前記第1のメモリ領域に関する更新情報を生成し、前記更新情報を第1記憶部に記憶する登録手段と、

前記第1記憶部に記憶される更新情報を出力する出力手段と、

を備えることを特徴とする情報処理装置。

【請求項 2】

検出したデータにより変更される前記第1のメモリ領域と、最後に変更された前記第1のメモリ領域を記憶する第2記憶部に記憶した前記第1のメモリ領域と、が異なる場合、前記更新情報を生成して該更新情報を前記第1記憶部に記憶するとともに、前記変更される前記第1のメモリ領域を前記第2記憶部に記憶する、

ことを特徴とする請求項1に記載の情報処理装置。

【請求項 3】

前記登録手段は、前記検出手段により検出したデータを出力した前記外部装置と、前記外部装置を識別する識別情報を記憶する外部装置記憶手段に記憶された識別情報が示す前記外部装置と、が同じ場合、前記更新情報を生成して該更新情報を前記第1記憶部に記憶する、

ことを特徴とする請求項1に記載の情報処理装置。

【請求項 4】

前記更新情報は、前記検出手段により検出したデータにより変更された前記第1のメモリ領域のページ数を計数した計数情報を含み、

検出したデータにより変更される前記第1のメモリ領域を含むページと、前記第2記憶部に記憶した前記第1のメモリ領域を含むページの次のページと、が同一の場合、前記更新情報に含まれる計数情報が示す計数値に1を加算するとともに、前記変更される前記第1のメモリ領域を前記第2記憶部に記憶する、

ことを特徴とする請求項2に記載の情報処理装置。

【請求項 5】

前記登録手段は、過去に変更された複数の前記第1のメモリ領域を前記第2記憶部に記憶し、前記検出手段により検出したデータにより変更される前記第1のメモリ領域が、前記第2記憶部に記憶した前記第1のメモリ領域のいずれかとも異なる場合に、前記更新情報を生成して該更新情報を前記第1記憶部に記憶するとともに、前記変更される前記第1のメモリ領域を前記第2記憶部に記憶する、

ことを特徴とする請求項2に記載の情報処理装置。

【請求項 6】

前記更新情報は、変更された前記第1のメモリ領域を含むページのサイズを拡張する拡張情報を含み、

検出したデータにより変更される前記第1のメモリ領域と一定範囲内の前記第1のメモリ領域が前記第2記憶部に含まれる場合、該第2記憶部に含まれる前記第1のメモリ領域についての前記更新情報の拡張情報を、前記検出手段により検出したデータにより変更される前記第1のメモリ領域を含むページのサイズに更新する、

ことを特徴とする請求項5に記載の情報処理装置。

10

20

30

40

50

【請求項 7】

ハードウェア資源を管理することにより、1以上の仮想機械を実現するステップと、
前記仮想機械に割り当てられる第1のメモリ領域のアドレスと、前記第1のメモリ領域
の実メモリである第2のメモリ領域のアドレスと、を相互に変換することにより、外部装
置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第1
のメモリ領域に直接的にデータ転送を行なうステップと、

前記入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデ
ータ転送されるデータを検出するステップと、

検出したデータが一定の条件を満たす場合に、前記検出したデータにより変更される前
記第1のメモリ領域に関する更新情報を生成し、前記更新情報を第1記憶部に記憶するス
テップと、

前記第1記憶部に記憶される更新情報を出力するステップと、
を備えることを特徴とする情報処理装置の仮想化方法。

【請求項 8】

ハードウェア資源を管理することにより、1以上の仮想機械を実現する仮想機械実現手
段と、

前記仮想機械に割り当てられる第1のメモリ領域のアドレスと、前記第1のメモリ領域
の実メモリである第2のメモリ領域のアドレスと、を相互に変換することにより、外部装
置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第1
のメモリ領域に直接的にデータ転送を行なうデータ転送手段と、

移動対象の前記仮想機械に割り当てられる第1のメモリ領域に格納されるデータを取得
して移動先に転送するメモリ領域転送手段と、

前記入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデ
ータ転送されるデータを検出する検出手段と、

検出したデータが一定の条件を満たす場合に、該検出したデータにより変更される前記
第1のメモリ領域に関する更新情報を生成し、前記更新情報を第1記憶部に記憶する登録
手段と、

前記第1記憶部から更新情報を取得する更新情報取得手段と、

前記更新情報取得手段により取得した更新情報に基づいて、移動対象の前記仮想機械に
割り当てられる第1のメモリ領域において前記検出手段により検出したデータにより変更
された更新データを、前記移動先に転送する更新データ転送手段と、

を備えることを特徴とする情報処理装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、1以上の仮想機械を実現する情報処理装置に関する。

【背景技術】

【0002】

従来、情報処理装置としてのサーバ上で複数のVM (Virtual Machine
: 仮想機械) といわれる仮想的なサーバを動作させるサーバ仮想化技術が知られている。

図1は、サーバ仮想化技術の概要を示す図である。

【0003】

サーバ101では、VMM (Virtual Machine Monitor) 102と呼ばれるソフトウェアを動作させている。このVMM 102は、サーバ101のメモ
リやIO (Input/Output: 入出力) などに関するハードウェア資源を管理す
る。また、VMM 102は、各VM #0 ~ #2に対して各VMが必要とするハードウェア
資源をエミュレート (模倣) して提供する。

【0004】

複数のVMを動作させる場合、複数のVMがサーバ101に備わる物理メモリの同一ア
ドレスの記憶領域を重複して使用しないように、メモリの排他制御が行なわれる。この排

10

20

30

40

50

他制御は、VMM102のメモリ管理によって行なわれる。

【0005】

図2は、VMM102によるメモリ管理の概要を示す図である。

VMが認識するメモリは、ゲスト物理メモリと呼ばれる。このゲスト物理メモリは、VMからは連続するメモリ領域として把握される。一方、VMM102から見るとゲスト物理メモリはVM毎に存在する別々のメモリ空間として把握される。

【0006】

図2は、VM#0とVM#1の2つのVMがVMM102上で動作している場合のゲスト物理メモリを示している。この場合、VM#0のゲスト物理メモリ#0とVM#1のゲスト物理メモリ#1の2つのゲスト物理メモリが存在する。

10

【0007】

例えば、アドレスXは、ゲスト物理メモリ#0および#1にそれぞれ1つずつ存在する。このアドレスXをそのままサーバ101に備わる実メモリ、すなわちホスト物理メモリに割り当てると競合が発生する。

【0008】

そこで、VMM102は、ゲスト物理メモリ#0のアドレスXをホスト物理メモリのアドレスZに割り当て、ゲスト物理メモリ#1のアドレスXをホスト物理メモリのアドレスYに割り当てる。このように、VMM102は、ホスト物理メモリを別々のホスト物理メモリに割り当てることにより、メモリの競合を回避する。

20

【0009】

VMがDMA(Direct Memory Access)を起動する場合にも同様のメモリ管理が必要となる。

例えば、VM#0が、DMA転送により、ゲスト物理メモリ#0のアドレスXにデータを書き込む場合を考える。なお、DMA転送により所定のメモリ領域にデータを書き込むことを「DMAライト」という。また、DMA転送によるデータの書き込み先を示すメモリアドレスを「DMAアドレス」という。

【0010】

この場合、VM#0がIOアダプタにDMAアドレスとしてアドレスXを設定すると、DMAライトはホスト物理メモリのアドレスXに対して行われる。ホスト物理メモリのアドレスXはゲスト物理メモリ#0のアドレスXとは無関係のメモリ領域であるため、メモリ破壊が引き起こされる。メモリ破壊が引き起こされた場合には、システムにパニックが発生し、その後システムが停止する。

30

【0011】

したがって、VM#0によるDMAライトの実行は、以下のような手順で行なわれる。

(1) VM#0は、IOアダプタに対して、ゲスト物理メモリ#0のアドレスXへのDMAライトを要求する。

(2) VMM102は、VM#0のDMAライトの要求をトラップし、ゲスト物理メモリ#0のアドレスXをホスト物理メモリのバッファ領域のアドレスWに変換する。そして、VMM102は、変換したアドレスWをIOアダプタのDMAアドレス設定レジスタに設定する。

40

(3) VM#0は、IOアダプタに対してDMA開始を指示する。

(4) IOアダプタは、アドレスWに対してDMAライトを実行する。

(5) DMAライトが完了すると、IOアダプタは、DMA完了割り込みによりDMAライトの完了を通知する。

(6) VMM102は、ホスト物理メモリのアドレスWに格納されたデータをゲスト物理メモリ#0のアドレスXにコピーする。

(7) VMM102は、DMA完了割り込みによりVM#0にDMAライトの完了を通知する。

(8) VM#0は、ゲスト物理メモリ#0のアドレスXからデータを取り出す。

【0012】

50

サーバ仮想化に関する機能の一つに、ライブ・マイグレーション(Live Migration)という技術がある。このライブ・マイグレーションとは、あるサーバで稼働中のVMを、別のサーバに動作を停止することなく当該動作を移行させる技術である。

【0013】

図3は、ライブ・マイグレーションの概要を示す図である。

ライブ・マイグレーションは、図3に示すように、2台のサーバ#0および#1と、サーバ#0および#1が共有するストレージ320と、サーバ#0および#1が接続するネットワーク330と、を備える環境で行うことが可能である。ライブ・マイグレーションは、以下のような手順で行なわれる。

(1) 移動先サーバ#1で動作するVMM311は、VM312を用意する。

10

(2) 移動元VM302が業務を継続している状態で、移動元VMM301は、移動元VM302が使用しているメモリ内容を移動先VM310に転送する。この処理を「プレ・コピー」という。

(3) 移動元VM302の業務を一時的に停止し、移動元VMM301は、移動元VM302が使用しているメモリ内容を移動先VM310に転送する。この処理を「ストップ・アンド・コピー」という。

(4) 移動先VMM311からの指示に応じて、移動先VM312は業務を再開する。

【0014】

上述のように、プレ・コピーやストップ・アンド・コピーなどのメモリ・コピーはすべてVMMが行なう。VMMがメモリ・コピーを実行できるのは、VMMが宿主物理メモリを管理しているからである。

20

【0015】

VMMは、CPU(Central Processing Unit)から宿主物理メモリへのライト処理だけでなくDMAライトについても介在し、移動元VMのメモリ・データの変更分を検出する。

【0016】

IO仮想化技術としてIOMMU(Input/Output Memory Management Unit)という技術がある。このIOMMUは、上述のようなVMM介在によるメモリアクセスのオーバーヘッドを軽減することで高速化を図る技術である。

【0017】

IOMMUは、IOアダプタと宿主物理メモリとに接続し、ゲスト物理メモリアドレスと宿主物理メモリアドレスの変換を実施するメモリ管理ユニットである。

30

IOMMUを実装すると、図4に示すように、IO装置からVMのゲスト物理メモリに直接DMAを実行することが可能となる。この技術を使用すると、DMAによるゲスト物理メモリへのアクセスの際に、VMMが介在する必要がなくなる。また、VMMが宿主物理メモリ-ゲスト物理メモリ間でのデータコピーを行なう必要もなくなる。

【0018】

しかし、上述のように、IOMMUを使用すると、DMAアクセス時にVMMは介在しないことになる。この場合、VMMは、DMAライトによるゲスト物理メモリの更新を把握できないので、VMMは、メモリ・コピーを正しく実行することができない。そのため、例えば、図4に示したIOMMUを実装したサーバに、DMAライトによるゲスト物理メモリの変更内容を保持させるなどして対応していた。

40

【0019】

上記技術に関連して、仮想マシン間で、DMAなど、物理アドレスに対していくらかのアクセスができるようにしながら、ゲストから宿主への物理アドレスの参照に使用される参照テーブルからの物理アドレス空間のページなどを可能にする方法およびシステムが知られている。

【先行技術文献】

【特許文献】

【0020】

50

【特許文献1】特開2006-252554号公報

【発明の概要】

【発明が解決しようとする課題】

【0021】

しかし、図4に示したIOMMUを実装したサーバに、DMAライトによるゲスト物理メモリの変更内容を保持させる場合、DMAライトによるゲスト物理メモリの変更内容を保持させるために大容量の記憶装置が必要となる。すなわち、IOMMUを実装したサーバでライブ・マイグレーションを実行するには、大量のハードウェア資源が必要となる。

【0022】

1つの側面では、本発明は、少ないハードウェア資源でライブ・マイグレーションを実行できる情報処理装置を提供することを目的とする。

10

【課題を解決するための手段】

【0023】

1つの態様によれば、本情報処理装置は、以下の手段を備える。

仮想機械実現手段は、ハードウェア資源を管理することにより、1以上の仮想機械を実現する。

【0024】

データ転送手段は、前記仮想機械に割り当てられる第1のメモリ領域のアドレスと、前記第1のメモリ領域の実メモリである第2のメモリ領域のアドレスと、を相互に変換する。そして、データ転送手段は、外部装置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送を行なう。

20

【0025】

検出手段は、前記入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送されるデータを検出する。

登録手段は、検出したデータが一定の条件を満たす場合に、該検出したデータにより変更される前記第1のメモリ領域に関する更新情報を生成し、前記更新情報を第1記憶部に記憶する。

【0026】

出力手段は、前記第1記憶部に記憶される更新情報を出力する。

30

【発明の効果】

【0027】

1つの態様では、少ないハードウェア資源でライブ・マイグレーションを実行することが可能となる。

【図面の簡単な説明】

【0028】

【図1】サーバ仮想化技術の概要を示す図である。

【図2】VMMによるメモリ管理の概要を示す図である。

【図3】ライブ・マイグレーションの概要を示す図である。

【図4】IOMMUについて説明する図である。

【図5】情報処理装置を説明する図である。

40

【図6】情報処理装置の構成例を示す図である。

【図7】ノース・ブリッジにおけるDMA処理の動作を説明する図である。

【図8】PCIeスイッチにおけるDMA処理の動作を説明する図である。

【図9】パケット検出部の構成例を示す図である。

【図10】FIFO#0の構成例を示す図である。

【図11】DMAパケットのヘッダの構成例を示す図である。

【図12】ダーティ・ページ管理ユニットによるダーティ・ページ情報の登録処理を示すフローチャートである。

【図13】ダーティ・ページ管理ユニットによるダーティ・ページ情報の出力処理を示すフローチャートである。

50

【図 1 4】本実施例に係るライブ・マイグレーションの概要を説明する図である。

【図 1 5】本実施例に係るライブ・マイグレーションの概要を示すフローチャートである。

【図 1 6】プレ・コピーの具体的な処理を示すフローチャートである。

【図 1 7】ストップ・アンド・コピーの具体的な処理を示すフローチャートである。

【図 1 8】条件 1 を用いたパケット検出部の具体例を示す図である。

【図 1 9】図 1 8 に示したパケット検出部の処理を示すフローチャートである。

【図 2 0】条件 2 を用いたパケット検出部の具体例を示す図である。

【図 2 1】図 2 0 に示したパケット検出部の処理を示すフローチャートである。

【図 2 2】ライブ・マイグレーション時における S I D 記憶部への S I D 設定処理を示すフローチャートである。 10

【図 2 3】条件 3 を用いたパケット検出部の具体例を示す図である。

【図 2 4】図 2 3 に示したパケット検出部の処理を示すフローチャートである。

【図 2 5】条件 4 を用いたパケット検出部の具体例を示す図である。

【図 2 6】図 2 5 に示したパケット検出部の処理を示すフローチャートである。

【図 2 7】条件 5 を用いたパケット検出部の具体例を示す図である。

【図 2 8】図 2 7 に示したパケット検出部の処理を示すフローチャートである。

【図 2 9】条件 1 および 2 を用いたパケット検出部の具体例を示す図である。

【図 3 0】図 2 9 に示したパケット検出部の処理を示すフローチャートである。

【図 3 1】条件 2 および 3 を用いたパケット検出部の具体例を示す図である。 20

【図 3 2】図 3 1 に示したパケット検出部の処理を示すフローチャートである。

【図 3 3】条件 2 および 4 を用いたパケット検出部の具体例を示す図である。

【図 3 4】図 3 3 に示したパケット検出部の処理を示すフローチャートである。

【図 3 5】条件 2 および 5 を用いたパケット検出部の具体例を示す図である。

【図 3 6】図 3 5 に示したパケット検出部の処理を示すフローチャートである。

【図 3 7】条件 3 を用いたパケット検出部を実現する場合に使用する F I F O の構成例を示す図である。

【図 3 8】条件 5 を用いたパケット検出部を実現する場合に使用する F I F O の構成例を示す図である。

【発明を実施するための形態】 30

【0029】

以下、本実施形態の一例について、図 5 ~ 図 3 8 に基づいて説明する。なお、以下に説明する実施形態はあくまでも例示であり、以下に明示しない種々の変形や技術の適用を排除する意図ではない。すなわち、本実施形態は、その趣旨を逸脱しない範囲で、各実施例を組み合わせるなど種々変形して実施することができる。

【実施例】

【0030】

図 5 は、本実施例に係る情報処理装置 5 0 0 を説明する図である。

情報処理装置 5 0 0 は、仮想機械実現手段 5 0 1 と、データ転送手段 5 0 2 と、検出手段 5 0 3 と、登録手段 5 0 4 と、出力手段 5 0 5 と、を備える。 40

【0031】

仮想機械実現手段 5 0 1 は、ハードウェア資源 5 0 6 を管理することにより、1 または 2 以上の仮想機械 5 0 7 を実現する。仮想機械実現手段 5 0 1 は、情報処理装置 5 0 0 に含まれる CPU に所定のプログラム命令を実行させることによって実現することができる。ハードウェア資源 5 0 6 には、情報処理装置 5 0 0 に備わるメモリや、I / O などに関するハードウェアを含むことができる。なお、図 5 では、第 1 のメモリ領域 5 0 8 とハードウェア資源 5 0 6 とを別々に記載しているが、第 1 のメモリ領域 5 0 8 はハードウェア資源 5 0 6 に含むこともできる。

【0032】

データ転送手段 5 0 2 は、仮想機械 5 0 7 に割り当てられる第 1 のメモリ領域 5 0 8 の 50

アドレスと、第1のメモリ領域508の実メモリである第2のメモリ領域のアドレスと、を相互に変換する。これにより、データ転送手段502は、入出力装置510から仮想機械507に割り当てられる第1のメモリ領域508に直接的にデータ転送を行なう。このようなデータ転送は、例えば、IOMMUなどを利用して実現することができる。なお、入出力装置510は、情報処理装置500と通信可能に接続する外部装置509とのデータの入出力を制御する装置である。

【0033】

検出手段503は、入出力装置510から仮想機械507に割り当てられる第1のメモリ領域508に直接的にデータ転送されるデータを検出する。

登録手段504は、検出手段503により検出したデータが一定の条件を満たす場合に、検出手段503が検出したデータにより変更される第1のメモリ領域508に関する更新情報を生成し、更新情報を第1記憶部511に記憶する。更新情報には、検出手段503が検出したデータにより変更される第1のメモリ領域508のアドレスと、データの出力元の出力もとのユニットを識別する識別情報と、を含むことができる。

【0034】

出力手段505は、第1記憶部511に記憶される更新情報を出力する。

以上のように、登録手段504は、検出手段503が検出したデータにより変更される第1のメモリ領域508に関する更新情報を第1記憶部511に記憶する。そして、出力手段505は、第1記憶部511に記憶される更新情報を出力する。

【0035】

したがって、情報処理装置500で動作する仮想機械507は、データ転送手段502によって、入出力装置510から仮想機械507に割り当てられる第1のメモリ領域508に直接的にデータ転送されたメモリ領域の更新情報を得ることが可能となる。

【0036】

その結果、入出力装置510から、仮想機械507に割り当てられる第1のメモリ領域508に直接的にデータ転送するデータ転送手段502を備える情報処理装置500であっても、仮想機械507のライブ・マイグレーションを実行することが可能となる。

【0037】

また、登録手段504は、検出手段503により検出したデータが一定の条件を満たす場合にだけ更新情報を生成し、更新情報を第1記憶部511に記憶する。その結果、第1記憶部511を実現するためのハードウェア資源を少なくすることができる。

【0038】

以上のように、入出力装置510から第1のメモリ領域508に直接的にデータ転送するデータ転送手段502を備える情報処理装置500であっても、少ないハードウェア資源で仮想機械507のライブ・マイグレーションを実行することが可能となる。

【その他の実施例】

【0039】

図6は、情報処理装置600の構成例を示す図である。

情報処理装置600は、演算処理装置としてのCPU610および611と、主記憶装置としてのメモリ620と、メモリ制御装置としてのノース・ブリッジ630と、を備える。又、情報処理装置600は、入出力制御装置としてのIOアダプタ640、641および642と、PCIeスイッチ(PCI Express Switch)650と、を備える。

【0040】

CPU610および611は、メモリ620に展開されるプログラムを実行する演算処理装置である。CPU610および611は、所定のプログラムを実行することにより、サーバの仮想化を実現する。また、CPU610および611は、ストレージ#0や#1などから所定のプログラムを読み出して実行することにより、本実施例に係るライブ・マイグレーションを実現する。

【0041】

10

20

30

40

50

なお、サーバの仮想化に必要なハードウェア資源やメモリ等管理などの技術については従来技術を利用することができる。

メモリ620は、CPU610および611が実行するプログラムやデータを記憶する揮発性メモリである。例えば、RAM(Random Access Memory)などがメモリ620として使用される。なお、必要に応じて不揮発性メモリを使用してもよい。

【0042】

ノース・ブリッジ630は、CPU610および611と、メモリ620およびPCIeスイッチ650と接続している。ノース・ブリッジ630は、CPU610および611と、メモリ620やPCIeスイッチ650と、が相互に通信可能となるようにデータ転送路を制御する。

10

【0043】

また、ノース・ブリッジ630は、IOMMU631と、IOTテーブル記憶部632と、を備える。

IOMMU631は、ゲスト物理メモリのアドレスとホスト物理メモリのアドレスとの変換を実施するなどのメモリ管理を行なうメモリ管理ユニットである。以下、ゲスト物理メモリのアドレスを「ゲスト物理アドレス」といい、「GPA(Guest Physical Address)」と略記する。また、ホスト物理メモリのアドレスを「ホスト物理アドレス」といい、「HPA(Host Physical Address)」と略記する。

20

【0044】

IOTテーブル記憶部632は、IOMMU631がGPAとHPAの変換を行なうときに使用するIOTテーブルを記憶する記憶装置である。例えば、キャッシュメモリなどがIOTテーブル記憶部632として使用される。IOTテーブルについては、図7で説明する。

【0045】

IOアダプタ640、641および642は、情報処理装置600と接続するIO装置とのインタフェースである。DMAを行なうDMA回路は、各IOアダプタ640、641および642に備わる。IO装置とは、例えば、図6に示すデータを記憶するストレージ#0、#1やネットワークと接続するネットワーク機器などの入出力装置である。

【0046】

PCIeスイッチ650は、ノース・ブリッジ630、IOアダプタ640、641および642と接続している。また、PCIeスイッチ650は、ノース・ブリッジ630とIOアダプタ640、641および642との間を結ぶデータ転送路の切り替え制御などを行なう。情報処理装置600は、データ転送路にシリアル転送インタフェースの規格である「PCI Express 2.0」を使用している。

30

【0047】

また、PCIeスイッチ650は、ダーティ・ページ管理ユニット651と、ダーティ・ページ記憶部652と、を備える。

ダーティ・ページ管理ユニット651は、IOアダプタ640、641または642からメモリ620に、DMAライトによって転送されるデータとしてのパケットを検出する。そして、ダーティ・ページ管理ユニット651は、DMAライトによって追加、変更または更新等されたデータに関する情報をページ単位でダーティ・ページ記憶部652に記憶する。

40

【0048】

以下、DMAライトによって追加、変更または更新等されたページ単位のデータを、「ダーティ・ページ」という。また、ダーティ・ページに関する情報を、「ダーティ・ページ情報」という。

【0049】

ダーティ・ページ記憶部652は、ダーティ・ページ情報を記憶する記憶装置である。例えば、スタティックRAMなどがダーティ・ページ記憶部652として使用される。

50

情報処理装置 600 は、ライブ・マイグレーションを実施するために情報処理装置 601 とネットワークを介して接続している。また、情報処理装置 600 と情報処理装置 601 は、ストレージ #0 および #1 を共有する。

【0050】

情報処理装置 601 は、CPU 660 および 661 と、メモリ 670 と、ノース・ブリッジ 680 と、I/O アダプタ 690、691 および 692 と、PCIe スイッチ 700 と、を備える。

【0051】

CPU 660、661、メモリ 670 および ノース・ブリッジ 680 は、CPU 610、611、メモリ 620 および ノース・ブリッジ 630 と、同様の機能を有する。また、I/O アダプタ 690、691、692 および PCIe スイッチ 700 は、I/O アダプタ 640、641、642 および PCIe スイッチ 650 と、同様の機能を有する。

10

【0052】

ただし、ノース・ブリッジ 680 には、IOMMU 631 および I/O テーブル記憶部 632 を備えなくてもよい。同様に、PCIe スイッチ 700 には、ダーティ・ページ管理ユニット 651 および ダーティ・ページ記憶部 652 を備えなくてもよい。

【0053】

情報処理装置 600 と接続する情報処理装置 601 は、図 3 で示した一般的なライブ・マイグレーションが実行可能であればよい。

【0054】

20

図 6 に示した構成は、本実施例に係る情報処理装置 600 の一例である。したがって、図 6 に示す構成に情報処理装置 600 の構成を限定する趣旨ではない。例えば、CPU の数や I/O アダプタの数、情報処理装置 600 を構成する各ユニットの配置などを図 6 に示すものに限定する趣旨ではない。

【0055】

また、情報処理装置 600 は、CD や DVD などの可搬記憶媒体を駆動する媒体駆動装置 643 を備えてもよい。本実施例に係るライブ・マイグレーションを CD や DVD などの可搬記憶媒体に記憶した場合、CPU 610 および 611 は、媒体駆動装置 643 を介して可搬記憶媒体から所定のプログラムを読み出して実行することにより、本実施例に係るライブ・マイグレーションを実現する。

30

【0056】

図 7 は、ノース・ブリッジ 630 における DMA 処理の動作を説明する図である。

DMA パケット 710 および 711 は、DMA 処理により I/O アダプタ 640、641 または 642 から所定の転送先に転送されるデータとしてのパケットである。

【0057】

I/O テーブル 720 は、I/O テーブル記憶部 632 に記憶される I/O テーブルを示している。I/O テーブル 720 は、GPA と HPA のアドレス変換テーブルである。I/O テーブル 720 は、GPA とソース ID とを含むアドレスに対応する HPA を定義する。ソース ID とは、パケットの転送元のユニットを識別する ID (Identification Data) である。以下では、ソース ID を「SID (Source ID)」と略記する。

40

【0058】

I/O テーブル 720 は、後述するページテーブル 730 のエントリのうち頻繁に使用する一部のエントリをキャッシュ (一時保持) するものである。

ページテーブル 730 は、GPA と HPA のアドレス変換テーブルである。ページテーブル 730 は、GPA に対応する HPA を定義している。

【0059】

IOMMU 631 は、DMA パケット 710 を検出すると、検出した DMA パケット 710 のヘッダから、GPA と SID とを取得する。

すると、IOMMU 631 は、I/O テーブル 720 を参照する。そして、IOMMU 6

50

31は、DMAパケット710から取得したGPAとSIDとを含むアドレスが、IOTテーブル720に登録されているか否かを判別する。

【0060】

DMAパケット710から取得したGPAとSIDとを含むアドレスが登録されている場合、IOMMU631は、そのGPAとSIDとを含むアドレスに対応するHPAを取得する。

【0061】

そして、IOMMU631は、DMAパケット710のヘッダに設定されている転送先を示すGPAを、IOTテーブル720から取得したHPAに変更する。

一方、DMAパケット710から取得したGPAとSIDとを含むアドレスがIOTテーブル720に登録されていない場合、IOMMU631は、メモリ620に格納されているページテーブル730を参照する。

【0062】

そして、IOMMU631は、ページテーブル730のアドレスであって、DMAパケット710から取得したGPAと同一のアドレスに登録されているHPAを取得する。

そして、IOMMU631は、DMAパケット710のヘッダに設定されている転送先を示すGPAを、ページテーブル730から取得したHPAに変更する。

【0063】

以上のようにして転送先のアドレスがGPAからHPAに変換されたDMAパケット711は、メモリ620に出力される。

【0064】

図8は、PCIeスイッチ650におけるDMA処理の動作を説明する図である。

ダーティ・ページ管理ユニット651は、制御I/F(Interface)部810と、パケット検出部820と、ライトポインタ830と、リードポインタ840と、を備える。

【0065】

なお、ダーティ・ページ管理ユニット651の動作の理解を容易にするために、図8では、ダーティ・ページ管理ユニット651の中にダーティ・ページ記憶部652を記載している。

【0066】

ダーティ・ページ記憶部652は、FIFO(First In First Out:先入れ先出しメモリ)を実現するメモリ#0および#1を備える。以下、FIFOを実現するメモリ#0および#1をそれぞれ「FIFO#0」、「FIFO#1」という。

【0067】

制御I/F部810は、ダーティ・ページ管理ユニット651のインタフェースである。VMM102は、制御I/F部810を介してFIFO#0または#1に記憶されているダーティ・ページ情報を取得する。

【0068】

制御I/F部810は、ステータスレジスタ811と、データレジスタ812と、コントロールレジスタ813と、を備える。

ステータスレジスタ811は、VMM102などのソフトウェアからリード/ライトが可能なレジスタである。ステータスレジスタ811は、以下のような状態表示ビットを備える。

【0069】

オーバーフロー情報ビット#0:FIFO#0がオーバーフローを起こしたことを示す。

オーバーフロー情報ビット#1:FIFO#1がオーバーフローを起こしたことを示す。

【0070】

有効データカウント#0:FIFO#0において、ライトポインタ830が示すアドレ

10

20

30

40

50

スとリードポインタ 8 4 0 が示すアドレスとの差をダーティ・ページ情報のサイズで除した値を示す。まだ読み出していないダーティ・ページ情報が F I F O # 0 にどれだけ残っているかを示す。

【 0 0 7 1 】

有効データカウンタ # 1 : F I F O # 1 において、ライトポインタが示すアドレスとリードポインタが示すアドレスとの差をダーティ・ページ情報のサイズで除した値を示す。まだ読み出していないダーティ・ページ情報が F I F O # 1 にどれだけ残っているかを示す。

【 0 0 7 2 】

データレジスタ 8 1 2 は、V M M 1 0 2 などのソフトウェアからリードが可能なレジスタである。F I F O # 0 または # 1 に記憶されているデータを V M M 1 0 2 から読み出すために使用するレジスタである。

10

【 0 0 7 3 】

データレジスタ 8 1 2 は、コントロールレジスタ 8 1 3 のリードセレクトビットで指定した F I F O # 0 または # 1 から読み出されたダーティ・ページ情報がセットされる。ダーティ・ページ管理ユニット 6 5 1 は、データレジスタ 8 1 2 に格納されているダーティ・ページ情報がリードされると同時に、リードポインタ 8 4 0 に、次のダーティ・ページ情報が記憶されるアドレスに設定する。

【 0 0 7 4 】

コントロールレジスタ 8 1 3 は、V M M 1 0 2 などのソフトウェアからリード/ライトが可能なレジスタである。コントロールレジスタ 8 1 3 は、以下のような制御ビットを備える。

20

【 0 0 7 5 】

スタートビット : F I F O # 0 または # 1 へのダーティ・ページ情報の登録を開始する。

ストップビット : F I F O # 0 または # 1 へのダーティ・ページ情報の登録を停止する。

【 0 0 7 6 】

クリアビット : F I F O # 0 または # 1 のポインタを 0 に戻し、F I F O # 0 または # 1 を空の状態にする。空の状態とは、F I F O に登録されているダーティ・ページ情報がない状態である。

30

【 0 0 7 7 】

ライトセレクトビット : 2 つの F I F O # 0 または # 1 のどちらにダーティ・ページ情報を登録するかを選択する。

リードセレクトビット : 2 つの F I F O # 0 または # 1 のどちらからダーティ・ページ情報を読み出すかを選択する。

【 0 0 7 8 】

パケット検出部 8 2 0 は、コントロールレジスタ 8 1 3 のスタートビットが「1」に設定されている間、P C I e スイッチ 6 5 0 に入力される D M A ライトパケット 8 5 0 を検出する。そして、パケット検出部 8 2 0 は、検出した D M A ライトパケット 8 5 0 のヘッダから、G P A と S I D を取得する。さらに、パケット検出部 8 2 0 は、取得した G P A からページ部を取得する。以下、G P A から取得したページ部のデータを「G P A ページ」という。G P A ページについては、図 1 0 で後述する。さらに、パケット検出部 8 2 0 は、G P A ページと S I D とを含むデータをダーティ・ページ情報として、選択中の F I F O # 0 または # 1 に書き込む。同時に、パケット検出部 8 2 0 は、ライトポインタ 8 3 0 に、次にダーティ・ページ情報を記憶する領域のアドレスを設定する。

40

【 0 0 7 9 】

ここで、F I F O # 0 および # 1 は、オーバーフローを避けるために、サイズを大きくする必要がある。例えば、D M A 性能が 1 0 G b p s 、F I F O # 0 と # 1 との切り替え時間が 1 m s の場合、約 1 万エントリのダーティ・ページ情報を記憶できる巨大な F I F

50

が必要となる。大量のハードウェア資源が必要となる。

【0080】

さらに、FIFOのサイズを大きくした場合でも、DMAライトが大量に発生すると、容易にFIFOがオーバーフローしてしまう可能性が増大する。CPUがFIFOに記憶されているダーティ・ページ情報を読み出す速度よりも、DMA転送の速度、すなわち、パケット検出部820がDMAライトを検出してダーティ・ページ情報をFIFOに書き込む速度の方が高速だからである。

【0081】

そこで、本実施例に係るパケット検出部820は、取得したGPAページとSIDの一方または両方が、特定の条件を満たすか否かを判別する。そして、特定の条件を満たす場合、パケット検出部820は、GPAページとSIDとを含むデータをダーティ・ページ情報として、選択中のFIFO#0または#1に書き込む。同時に、パケット検出部820は、ライトポインタ830に、次にダーティ・ページ情報を記憶する領域のアドレスを設定する。

10

【0082】

ライトポインタ830は、FIFO#0または#1へ最後に書き込まれたダーティ・ページ情報のアドレスを示す。本実施例では、図8に示すように、FIFO#0用とFIFO#1用とにそれぞれライトポインタ830が備えられている。

【0083】

リードポインタ840は、FIFO#0または#1から最後に読み出されたダーティ・ページ情報のアドレスを示す。本実施例では、図8に示すように、FIFO#0用とFIFO#1用とにそれぞれリードポインタ840を備える。

20

【0084】

なお、FIFO#0が選択されている場合には、FIFO#0についてのライトポインタおよびリードポインタがライトポインタ830およびリードポインタ840に格納される。

【0085】

同様に、FIFO#1が選択されている場合には、FIFO#1についてのライトポインタおよびリードポインタがライトポインタ830およびリードポインタ840に格納される。

30

【0086】

図9は、パケット検出部820の構成例を示す図である。なお、図9に記載のFIFO#Nは、選択中のFIFO#0または#1を示す。同様に、図18、図20、図23、図25、図27、図29、図31、図33および図35に記載のFIFO#Nも、選択中のFIFO#0または#1を示す。

【0087】

パケット検出部820は、FIFO制御部901と、判定部902と、を備える。

FIFO制御部901は、制御I/F部810に含まれるコントロールレジスタ813のスタートビットが「1」に設定されている間、DMAライトパケット850を検出する。そして、FIFO制御部901は、検出したDMAライトパケット850のヘッダから、GPAページとSIDを取得する。

40

【0088】

FIFO制御部901は、ダーティ・ページ情報を、選択中のFIFO#0または#1に書き込む必要がある、との判定結果を判定部902から得る。この場合、FIFO制御部901は、DMAライトパケット850から取得したGPAページとSIDとを含むデータを、ダーティ・ページ情報として、選択中のFIFO#0または#1に書き込む。同時に、FIFO制御部901は、ライトポインタ830に、次にダーティ・ページ情報を記憶する領域のアドレスを設定する。

【0089】

また、FIFO制御部901は、ダーティ・ページ情報を、選択中のFIFO#0また

50

は # 1 に書き込む必要がない、との判定結果を判定部 9 0 2 から得る。この場合、F I F O 制御部 9 0 1 は、ダーティ・ページ情報の、選択中の F I F O # 0 または # 1 への書き込みは行わない。

【 0 0 9 0 】

判定部 9 0 2 は、D M A ライトパケット 8 5 0 から取得した G P A ページと S I D の一方または両方が、特定の条件を満たすか否かを判別する。特定の条件を満たす場合、判定部 9 0 2 は、ダーティ・ページ情報を、選択中の F I F O # 0 または # 1 に書き込む必要があると判断し、その旨を F I F O 制御部 9 0 1 に通知する。

【 0 0 9 1 】

また、特定の条件を満たさない場合、判定部 9 0 2 は、ダーティ・ページ情報を、選択中の F I F O # 0 または # 1 に書き込む必要がないと判断し、その旨を F I F O 制御部 9 0 1 に通知する。

10

【 0 0 9 2 】

判定部 9 0 2 が使用する「特定の条件」として、下記「条件 1」～「条件 5」や、下記「条件 1」～「条件 5」の組み合わせなど、様々な条件を用いることができる。

条件 1：D M A ライトパケット 8 5 0 から取得した D M A アドレスが、選択中の F I F O # 0 または # 1 に最後に書き込まれた D M A アドレスと異なる場合だけ、ダーティ・ページ情報を、選択中の F I F O # 0 または # 1 に書き込む。条件 1 を用いた場合のパケット検出部 8 2 0 の具体例については、図 1 8 および図 1 9 で後述する。

【 0 0 9 3 】

条件 2：D M A ライトパケット 8 5 0 の送信元のユニットが、ライブ・マイグレーションの対象である場合だけ、ダーティ・ページ情報を、選択中の F I F O # 0 または # 1 に書き込む。条件 2 を用いた場合のパケット検出部 8 2 0 の具体例については、図 2 0 ~ 図 2 2 で後述する。

20

【 0 0 9 4 】

条件 3：D M A ライトが、選択中の F I F O # 0 または # 1 に最後に書き込まれた D M A アドレスを先頭アドレスとしたアクセスであって、複数ページをまたぐ連続アクセスでない場合だけ、ダーティ・ページ情報を、選択中の F I F O # 0 または # 1 に書き込む。条件 3 を用いた場合のパケット検出部 8 2 0 の具体例については、図 2 3 および図 2 4 で後述する。

30

【 0 0 9 5 】

条件 4：D M A ライトパケット 8 5 0 から取得した D M A アドレスが、選択中の F I F O # 0 または # 1 に書き込まれた複数の D M A アドレスのいずれとも異なる場合だけ、ダーティ・ページ情報を、選択中の F I F O # 0 または # 1 に書き込む。条件 4 を用いた場合のパケット検出部 8 2 0 の具体例については、図 2 5 および図 2 6 で後述する。

【 0 0 9 6 】

条件 5：D M A ライトパケット 8 5 0 から取得した D M A アドレスが、選択中の F I F O # 0 または # 1 に書き込まれた複数の D M A アドレスのいずれとも一定範囲内でない場合だけ、ダーティ・ページ情報を、選択中の F I F O # 0 または # 1 に書き込む。条件 5 を用いた場合のパケット検出部 8 2 0 の具体例については、図 2 7 および図 2 8 で後述する。

40

【 0 0 9 7 】

また、上記「条件 1」～「条件 5」を組み合わせ使用した場合のパケット検出部 8 2 0 の具体例については、図 2 9 ~ 図 3 8 で後述する。

【 0 0 9 8 】

図 1 0 は、F I F O # 0 の構成例を示す図である。なお、図 1 0 には、F I F O # 0 についてのみ示すが、F I F O # 1 についても同様の構成である。

【 0 0 9 9 】

F I F O # 0 は、D M A パケットの送信元の I O アダプタを示す S I D と、その D M A パケットの転送先を示す G P A の G P A ページと、を備える。P C I E x p r e s s の

50

場合、S I Dは、D M Aパケットの送信元のバス番号と、そのD M Aパケットの送信元のファンクション番号と、を含む。

【0100】

本実施例では、1ページのサイズを例えば4KBとする。この場合、64ビット幅のアドレスデータの下位12ビットは同一ページ内のアドレスを表す。したがって、本実施例では、G P Aのうち、ビット63～12までのデータを含むG P A [63 : 12]、すなわちG P AページをF I F O # 0に登録する。

【0101】

G P Aに32ビット長のデータを使用する場合、G P Aのビット63～32に「0」を補完した64ビット長のアドレスデータについてのG P AページをF I F O # 0に登録する。

10

上述のS I D、G P Aページは、D M Aパケットのヘッダから取得することができる。

【0102】

図11は、D M Aパケットのヘッダの構成例を示している。本実施例では、P C I eスイッチ650を使用しているので、D M Aパケットのヘッダは規格「P C I E x p r e s s 2 . 0」に準じる。

【0103】

図11の(1)に示すD M Aパケットのヘッダ1101は、G P Aに32ビット幅のアドレスデータを使用する場合のD M Aパケットのヘッダの構成例を示す。また、図11の(2)に示すD M Aパケットのヘッダ1102は、G P Aに64ビット幅のアドレスデータを使用する場合のD M Aパケットのヘッダの構成例を示す。

20

【0104】

「R」は、リザーブ領域である。常に「0」が設定される。

「F m t (F o r m a t)」は、ペイロードに格納されたデータの有無、ヘッダ長を示す2ビット幅のデータである。

【0105】

「T y p e」は、パケットのタイプを示す5ビット幅のデータである。

「T C (T r a n s a c t i o n C l a s s)」は、パケットの優先度を示す3ビット幅のデータである。

【0106】

「T D (T L P (T r a n s a c t i o n L a y e r P a c k e t) D i g e s t)」は、エラーチェック符号であるE C R C (E x t e n d e d C y c l i c a l R e d u n d a n c y C h e c k)の有無を示す1ビット幅のデータである。

30

【0107】

「E P (E r r o r P o i s o n e d)」は、ペイロードに格納されたデータが壊れている可能性を示す1ビット幅のデータである。

「A t t r (A t t r i b u t e s)」は、パケットの順序やプロトコルに関する補足情報を示す2ビット幅のデータである。

【0108】

「A T (A d d r e s s T r a n s l a t i o n)」は、アドレス変換に関する補足情報を示す2ビット幅のデータである。

40

「L e n g t h」は、ペイロードに格納されているデータのデータ長を示す10ビット幅のデータである。

【0109】

「バス番号」は、D M Aパケットの送信元のバス番号を示す8ビット幅のデータである。

「ファンクション番号」は、D M Aパケットの送信元のファンクション番号を示す8ビット幅のデータである。

【0110】

「T a g」は、D M Aパケットの管理番号を示す8ビット幅のデータである。

50

「Last DW BE」は、Last DW Byteの有効・無効を示す4ビット幅のデータである。

【0111】

「1st DW BE」は、1st DW Byteの有効・無効を示す4ビット幅のデータである。

「アドレス」は、30ビット幅または62ビット幅のアドレスデータである。

【0112】

図12は、ダーティ・ページ管理ユニット651によるダーティ・ページ情報の登録処理を示すフローチャートである。なお、図12中に記載のFIFOは、FIFO#0または#1のうち、コントロールレジスタ813のライトセレクトビットに設定されたFIFOのことを示す。

10

【0113】

VMM102などのソフトウェアが、ダーティ・ページ管理ユニット651に備わるコントロールレジスタ813のスタートビットを「1」に設定すると、ダーティ・ページ管理ユニット651は、ダーティ・ページ情報の登録処理を開始する(ステップS1200)。

【0114】

ステップS1201において、パケット検出部820は、IOアダプタ640、641または642などのIOアダプタからPCIeスイッチ650に入力したDMAパケットのヘッダを参照する。

20

【0115】

そして、パケット検出部820は、DMAパケットのヘッダから「Fmt」と「Type」を取得する。

パケット検出部820は、IOアダプタからメモリ620に転送され、Fmt = 10 (2進数) またはFmt = 11 (2進数)、かつType = 00000 (2進数) であるDMAライトパケットを検出すると、ステップS1202に移行する(ステップS1201 YES)。

【0116】

なお、Fmt = 10 (2進数) の場合、GPAに32ビット幅のアドレスデータを使用するDMAパケットであることを示している。また、Fmt = 11 (2進数) の場合、GPAに64ビット幅のアドレスデータを使用するDMAパケットであることを示している。

30

また、Type = 00000 (2進数) は、DMAライトパケットであることを示している。

【0117】

パケット検出部820は、PCIeスイッチ650に入力したDMAパケットが、DMAライトパケットでないと判断すると、ステップS1201を再度実行する(ステップS1201 NO)。

【0118】

ステップS1202において、パケット検出部820は、ライトポインタ830が指すアドレスとリードポインタ840が指すアドレスとの差が、FIFOのサイズにダーティ・ページ情報のサイズを除いた値に等しいか否かにより、FIFOがFullか否かを判別する。

40

【0119】

ライトポインタ830が指すアドレスとリードポインタ840が指すアドレスとの差が0の場合、パケット検出部820は、FIFOがFullと判断する(ステップS1202 YES)。この場合、パケット検出部820は、処理をステップS1203に移行する。

【0120】

ステップS1203において、パケット検出部820は、ステータスレジスタ811の

50

オーバーフロー情報ビットに「1」を設定すると、処理をステップS1209に移行する。そして、ダーティ・ページ管理ユニット651は、ダーティ・ページ情報の登録処理を終了する。

【0121】

ステップS1202において、ライトポインタ830が指すアドレスとリードポインタ840が指すアドレスとの差がFIFOのサイズにダーティ・ページ情報のサイズを除いた値に等しくない場合、パケット検出部820は、FIFOがFullでないと判断する(ステップS1202 NO)。この場合、パケット検出部820は、処理をステップS1204に移行する。

【0122】

ステップS1204において、パケット検出部820は、DMAライトパケットのヘッダからSIDとGPAを取得する。そして、パケット検出部820は、GPAからGPAページを取得する。

【0123】

ステップS1205において、パケット検出部820は、ステップS1204で取得したSIDおよびGPAページのいずれか一方または両方に基づいて、FIFOへのダーティ・ページ情報の書き込みが必要か否かを判定する。FIFOへのダーティ・ページ情報の書き込みは不要と判定した場合(ステップS1205 NO)、パケット検出部820は、ダーティ・ページ情報の登録処理を終了する(ステップS1209)。

【0124】

また、FIFOへのダーティ・ページ情報の書き込みが必要と判定した場合(ステップS1205 YES)、パケット検出部820は、処理をステップS1206に移行する。この場合、パケット検出部820は、ライトポインタ830が示すFIFOのアドレスに、ステップS1204で取得したSIDとGPAページとを含むデータをダーティ・ページ情報として登録する。

【0125】

ステップS1207において、パケット検出部820は、ダーティ・ページ情報を格納する領域サイズだけ、ライトポインタ830に格納されているアドレスをインクリメントする。

【0126】

ステップS1208において、パケット検出部820は、制御I/F部810に備わるコントロールレジスタ813を参照する。そして、コントロールレジスタ813のストップビットに「1」が設定されている場合(ステップS1208 YES)、パケット検出部820は、処理をステップS1209に移行する。そして、ダーティ・ページ管理ユニット651は、ダーティ・ページ情報の登録処理を終了する。

【0127】

ステップS1208において、コントロールレジスタ813のストップビットが「1」でない場合、パケット検出部820は、処理をステップS1201に移行する(ステップS1208 NO)。

【0128】

図13は、ダーティ・ページ管理ユニット651によるダーティ・ページ情報の出力処理を示すフローチャートである。なお、図13中に記載のFIFOは、FIFO#0または#1のうち、コントロールレジスタ813のリードセレクトビットに設定されたFIFOのことを示す。

【0129】

ステップS1301において、VMM102などのソフトウェアが、ダーティ・ページ管理ユニット651に備わるデータレジスタ812からダーティ・ページ情報を読み出すと、制御I/F部810は、処理をステップS1302に移行する。

【0130】

ステップS1302において、制御I/F部810は、ダーティ・ページ情報を格納す

10

20

30

40

50

る領域サイズだけ、リードポインタ 8 4 0 に格納されているアドレスをインクリメントする。

【 0 1 3 1 】

ステップ S 1 3 0 3 において、制御 I / F 部 8 1 0 は、リードポインタ 8 4 0 が示す F I F O のアドレスからダーティ・ページ情報を取得し、データレジスタ 8 1 2 にラッチする。

【 0 1 3 2 】

以上の処理が終了すると、制御 I / F 部 8 1 0 は、処理をステップ S 1 3 0 4 に移行する。そして、制御 I / F 部 8 1 0 は、ダーティ・ページ情報の出力処理を終了する。

【 0 1 3 3 】

図 1 4 は、本実施例に係るライブ・マイグレーションの概要を説明する図である。

図 1 4 に示す情報処理システム 1 4 0 0 は、サーバ # 0 と、サーバ # 0 とネットワークを介して互いに通信可能に接続するサーバ # 1 と、サーバ # 0 とサーバ # 1 とが共有するストレージ # 0 および # 1 と、を備える。

【 0 1 3 4 】

サーバ # 0 は、図 6 に示した情報処理装置 6 0 0 である。また、サーバ # 0 は、図 6 に示した情報処理装置 6 0 1 である。

サーバ # 0 では V M M # 0 が動作する。V M M # 0 は、V M # 0 および V M # 1 を実現する。一方、サーバ # 1 では V M M # 1 が動作する。また、V M # 0 はストレージ # 0 を専有し、V M # 1 はストレージ # 1 を専有する。

【 0 1 3 5 】

以上の状態において、サーバ # 0 で動作する V M M # 0 が実現する V M # 0 を、サーバ # 1 で動作する V M M # 1 上に移動するライブ・マイグレーションを行なう場合について以下に説明する。

【 0 1 3 6 】

図 1 5 は、本実施例に係るライブ・マイグレーションの概要を示すフローチャートである。

ステップ S 1 5 0 1 において、V M M # 0 は、V M # 0 の移動先である V M M # 1 に対して、新たな V M に割り当てるメモリ領域の確保を要求する。以下、V M に割り当てるメモリ領域を「V M 領域」という。

【 0 1 3 7 】

一方、V M M # 1 から V M 領域確保の要求を受けると、V M M # 1 は、サーバ # 1 に備わるメモリの所定の領域に V M 領域を確保する。そして、サーバ # 1 は、サーバ # 0 に対して V M 領域の確保が完了したことを通知する。

【 0 1 3 8 】

サーバ # 1 から V M 領域の確保が完了した旨の通知を受けると、サーバ # 0 は、処理をステップ S 1 5 0 2 に移行する。

ステップ S 1 5 0 2 において、V M M # 0 は、ある時点において、移動の対象である V M # 0 に割り当てられているホスト物理メモリの領域に格納されているデータを取得する。この、ある時点（チェックポイント時点）において V M に割り当てられているホスト物理メモリの領域に格納されているデータを「スナップショット」という。

【 0 1 3 9 】

ステップ S 1 5 0 3 において、V M M # 0 は、ダーティ・ページ管理ユニット 6 5 1 のコントロールレジスタ 8 1 3 のスタートビットを「1」に設定し、ダーティ・ページ情報の記録を開始する。

【 0 1 4 0 】

ステップ S 1 5 0 4 において、V M M # 0 は、ステップ S 1 5 0 2 で取得したスナップショットを、V M # 0 の移動先のサーバ # 1 に備わるメモリの所定の領域であって、ステップ S 1 5 0 1 で確保した V M 領域に、メモリ・コピーする。

【 0 1 4 1 】

10

20

30

40

50

ステップ S 1 5 0 5 において、VMM # 0 は、ダーティ・ページ管理ユニット 6 5 1 から制御 I / F 部 8 1 0 のステータスレジスタ 8 1 1 の有効データカウンタを介して、残りのダーティ・ページ情報数を取得する。そして、VMM # 0 は、残りのダーティ・ページ情報数とあらかじめ決められた閾値とを比較する。

【 0 1 4 2 】

残りのダーティ・ページ情報数が閾値より大きい場合 (ステップ S 1 5 0 5 YES)、VMM # 0 は、処理をステップ S 1 5 0 6 に移行する。また、残りのダーティ・ページ情報数が閾値以下の場合 (ステップ S 1 5 0 5 NO)、VMM # 0 は、処理をステップ S 1 5 0 7 に移行する。

【 0 1 4 3 】

ステップ S 1 5 0 6 において、VMM # 0 は、ダーティ・ページ管理ユニット 6 5 1 から制御 I / F 部 8 1 0 のデータレジスタ 8 1 2 を介して、ダーティ・ページ記憶部 6 5 2 に現在までに登録されている全てのダーティ・ページ情報を取得する。

【 0 1 4 4 】

そして、VMM # 0 は、取得したダーティ・ページ情報に含まれる GPA ページを HPA に変換する。

そして、VMM # 0 は、サーバ # 0 に備わるメモリから、変換した HPA に格納されているデータを取得する。VMM # 0 は、取得したデータを、VM # 0 の移動先のサーバ # 1 に備わるメモリの領域であって、ステップ S 1 5 0 1 で確保した VM 領域に、メモリ・コピーする。

【 0 1 4 5 】

メモリ・コピーを実行すると、VMM # 0 は、処理をステップ S 1 5 0 5 に移行する。

以上に説明したステップ S 1 5 0 2 ~ S 1 5 0 6 の処理が「プレ・コピー」である。プレ・コピーは、VM # 0 が業務を継続中に行なわれる。したがって、プレ・コピー中も VMM # 0 は、IOMMU 6 3 1 を介して、ストレージ # 0 に対するリード/ライト処理を行なう。そのため、以下のストップ・アンド・コピーが必要となる。

【 0 1 4 6 】

ステップ S 1 5 0 7 において、VMM # 0 は、VM # 0 の業務を停止する。

ステップ S 1 5 0 8 において、VMM # 0 は、ステップ S 1 5 0 6 と同様の処理を実行する。

【 0 1 4 7 】

以上に説明したステップ S 1 5 0 7 ~ S 1 5 0 8 の処理が「ストップ・アンド・コピー」である。

ステップ S 1 5 0 9 において、VMM # 0 は、VM # 0 の移動先のサーバ # 1 で動作している VMM # 1 に、メモリ・コピーの完了を通知する。

【 0 1 4 8 】

ステップ S 1 5 1 0 において、VMM # 0 からメモリ・コピーの完了通知を受けると、VMM # 1 は、VMM # 1 上で VM # 0 の業務を開始させる。VMM # 1 上で VM # 0 が業務を再開すると、VMM # 1 は、処理をステップ S 1 5 1 1 に移行する。そして、VMM # 0 および # 1 は、ライブ・マイグレーションを終了する。

【 0 1 4 9 】

図 1 6 は、プレ・コピーの具体的な処理を示すフローチャートである。

ステップ S 1 6 0 1 において、VMM # 0 は、制御 I / F 部 8 1 0 に備わるコントロールレジスタ 8 1 3 のクリアビットに「1」を設定することにより、FIFO # 0 および # 1 を初期化する。

【 0 1 5 0 】

また、VMM # 0 は、コントロールレジスタ 8 1 3 のライトセレクトビットを「0」に設定することにより、ダーティ・ページ情報を記録する FIFO を FIFO # 0 に設定する。同様に、VMM # 0 は、コントロールレジスタ 8 1 3 のリードセレクトビットを「0」に設定することにより、ダーティ・ページ情報を読み出す FIFO を FIFO # 0 に設

10

20

30

40

50

定する。

【0151】

ステップS1602において、VMM#0は、ダーティ・ページ管理ユニット651に備わるコントロールレジスタ813のスタートビットを「1」に設定することにより、ダーティ・ページ情報の記録を開始する。

【0152】

ステップS1603において、VMM#0は、メモリ・コピーを実行する。このステップS1603では、図15に示したステップS1502～S1504で説明したメモリ・コピー処理が行なわれる。

【0153】

ステップS1604において、VMM#0は、一定時間ウェイト状態に移行する。一定時間経過すると、VMM#0は、処理をステップS1605に移行する。

ステップS1605において、VMM#0は、制御I/F部810に備わるコントロールレジスタ813のライトセレクトビットに「0」または「1」を設定することにより、ダーティ・ページ情報を記録するFIFOを変更する。

【0154】

例えば、現在選択中のFIFOがFIFO#0の場合、VMM#0は、コントロールレジスタ813のライトセレクトビットに「1」を設定することにより、ダーティ・ページ情報を記録するFIFOをFIFO#1に変更する。

【0155】

同様に、現在選択中のFIFOがFIFO#1の場合、VMM#0は、コントロールレジスタ813のライトセレクトビットに「0」を設定することにより、ダーティ・ページ情報を記録するFIFOをFIFO#0に変更する。

【0156】

ステップS1606において、VMM#0は、制御I/F部810に備わるステータスレジスタ811を読み出す。

ステップS1607において、VMM#0は、現在選択中のFIFOのオーバーフロー情報ビットが「1」に設定されているか否かを判別する。

【0157】

現在選択中のFIFOのオーバーフロー情報ビットが「1」に設定されている場合（ステップS1607 YES）、VMM#0は、処理をステップS1608に移行する。この場合、VMM#0は、実行中のライブ・マイグレーションを中断、または最初からやり直す。

【0158】

ステップS1607において、現在選択中のFIFOのオーバーフロー情報ビットが「0」に設定されている場合（ステップS1607 NO）、VMM#0は、処理をステップS1609に移行する。

【0159】

ステップS1609において、VMM#0は、ステータスレジスタ811の有効データカウントに設定されている数だけデータレジスタ812を読み出し、ダーティ・ページ情報を取得する。

【0160】

ステップS1610において、VMM#0は、ステップS1609で読み出したダーティ・ページ情報のうち、SIDがストレージ#0を示すダーティ・ページ情報のみを抽出する。ステップS1610または後述するS1507の処理で抽出するダーティ・ページ情報を「コピー対象ダーティ・ページ情報」という。

【0161】

ステップS1611において、VMM#0は、コピー対象ダーティ・ページ情報に含まれるGPAページを、ページテーブル730にしたがって、HPAに変換する。この変換されたHPAもページを示す。ページの単位は、GPAページと同じである。

10

20

30

40

50

【 0 1 6 2 】

なお、図 2 3 で後述するように、条件 3 を用いたパケット検出部 8 2 0 を実現する場合、ダーティ・ページ情報には、カウント値が含まれる。この場合、コピー対象ダーティ・ページ情報に含まれる G P A ページを X とすると、V M M # 0 は、X から X + カウント値までの G P A ページを H P A に変換する。例えば、コピー対象ダーティ・ページ情報に含まれる G P A ページが 0 x 1 0 0、カウント値が 2 の場合、V M M # 0 は、0 x 1 0 0、0 x 1 0 1 および 0 x 1 0 2 の G P A ページを H P A に変換する。

【 0 1 6 3 】

また、図 2 7 で後述するように、条件 5 を用いたパケット検出部 8 2 0 を実現する場合、ダーティ・ページ情報には、拡張サイズが含まれる。この場合、V M M # 0 は、コピー対象ダーティ・ページ情報に含まれる G P A ページを、拡張サイズが示すページ単位に変換する。そして、V M M # 0 は、変換された値を H P A に変換する。例えば、コピー対象ダーティ・ページ情報に含まれる G P A ページの拡張サイズが 8 K B の場合、V M M # 0 は、G P A ページの下位 0 ~ 1 2 ビットをマスクした値を H P A に変換する。

10

【 0 1 6 4 】

ステップ S 1 6 1 2 において、V M M # 0 は、ステップ S 1 6 1 1 で変換された H P A が示すページの範囲に格納されているデータを、サーバ # 0 に備わるメモリ 6 2 0 から取得する。そして、V M M # 0 は、取得したデータを、V M # 0 の移動先である V M M # 1 に転送する。

【 0 1 6 5 】

一方、V M M # 1 は、V M M # 0 からデータを受信すると、受信したデータを図 1 5 に示したステップ S 1 5 0 1 の処理で確保した V M 領域に格納する。

20

ステップ S 1 6 1 3 において、V M M # 0 は、制御 I / F 部 8 1 0 に備わるコントロールレジスタ 8 1 3 のリードセレクトビットに「 0 」または「 1 」を設定することにより、ダーティ・ページ情報を読み出す F I F O を変更する。

【 0 1 6 6 】

例えば、現在選択中の F I F O が F I F O # 0 の場合、V M M # 0 は、コントロールレジスタ 8 1 3 のリードセレクトビットに「 1 」を設定することにより、ダーティ・ページ情報を読み出す F I F O を F I F O # 1 に変更する。

【 0 1 6 7 】

同様に、現在選択中の F I F O が F I F O # 1 の場合、V M M # 0 は、コントロールレジスタ 8 1 3 のリードセレクトビットに「 0 」を設定することにより、ダーティ・ページ情報を読み出す F I F O を F I F O # 0 に変更する。

30

【 0 1 6 8 】

ステップ S 1 6 1 4 において、V M M # 0 は、プレ・コピーの完了条件を満たしているか否かを判別する。このステップ S 1 6 1 4 の処理では、図 1 5 に示したステップ S 1 5 0 5 の処理が行なわれる。

【 0 1 6 9 】

プレ・コピーの完了条件が満たされていると判断した場合（ステップ S 1 6 1 4 Y E S）、V M M # 0 は、処理をステップ S 1 6 1 5 に移行する。この場合、V M M # 0 は、プレ・コピーを終了する。

40

【 0 1 7 0 】

また、プレ・コピーの完了条件を満たされていないと判断した場合（ステップ S 1 6 1 4 N O）、V M M # 0 は、処理をステップ S 1 6 0 4 に移行する。

【 0 1 7 1 】

図 1 7 は、ストップ・アンド・コピーの具体的な処理を示すフローチャートである。

図 1 6 に示したプレ・コピーが完了すると、V M M # 0 は、ストップ・アンド・コピー処理を開始する。V M M # 0 は、V M # 0 の業務を停止すると、処理をステップ S 1 7 0 1 に移行する（ステップ S 1 7 0 0）。

【 0 1 7 2 】

50

ステップ S 1 7 0 1 において、VMM # 0 は、一定時間ウェイト状態に移行する。DMA ライト処理が完全に停止するまでの待つためである。一定時間経過すると、VMM # 0 は、処理をステップ S 1 7 0 2 に移行する。

【 0 1 7 3 】

ステップ S 1 7 0 2 において、VMM # 0 は、制御 I / F 部 8 1 0 に備わるコントロールレジスタ 8 1 3 のライトセレクトビットに「 0 」または「 1 」を設定することにより、ダーティ・ページ情報を記録する FIFO を変更する。

【 0 1 7 4 】

例えば、現在選択中の FIFO が FIFO # 0 の場合、VMM # 0 は、コントロールレジスタ 8 1 3 のライトセレクトビットに「 1 」を設定することにより、ダーティ・ページ情報を記録する FIFO を FIFO # 1 に変更する。

10

【 0 1 7 5 】

同様に、現在選択中の FIFO が FIFO # 1 の場合、VMM # 0 は、コントロールレジスタ 8 1 3 のライトセレクトビットに「 0 」を設定することにより、ダーティ・ページ情報を記録する FIFO を FIFO # 0 に変更する。

【 0 1 7 6 】

ステップ S 1 7 0 3 において、VMM # 0 は、制御 I / F 部 8 1 0 に備わるステータスレジスタ 8 1 1 を読み出す。そして、VMM # 0 は、現在選択中の FIFO のオーバーフロー情報ビットが「 1 」に設定されているか否かを判別する。

【 0 1 7 7 】

20

ステップ S 1 7 0 4 において、VMM # 0 は、現在選択中の FIFO のオーバーフロー情報ビットが「 1 」に設定されている場合（ステップ S 1 7 0 4 YES）、VMM # 0 は、処理をステップ S 1 7 0 5 に移行する。この場合、VMM # 0 は、ライブ・マイグレーションを中断、または最初からやり直す。

【 0 1 7 8 】

ステップ S 1 7 0 4 において、現在選択中の FIFO のオーバーフロー情報ビットが「 0 」に設定されている場合（ステップ S 1 7 0 4 NO）、VMM # 0 は、処理をステップ S 1 7 0 6 に移行する。

【 0 1 7 9 】

ステップ S 1 7 0 6 において、VMM # 0 は、ステータスレジスタ 8 1 1 の有効データカウントに設定されている数だけデータレジスタ 8 1 2 を読み出し、ダーティ・ページ情報を取得する。

30

【 0 1 8 0 】

ステップ S 1 7 0 7 において、VMM # 0 は、ステップ S 1 7 0 6 で読み出したダーティ・ページ情報のうち、ダーティ・ページ情報に含まれる SID がストレージ # 0 を示すコピー対象ダーティ・ページ情報を抽出する。

【 0 1 8 1 】

ステップ S 1 7 0 8 において、VMM # 0 は、ステップ S 1 6 1 1 と同様に、コピー対象ダーティ・ページ情報に含まれる GPA ページを、ページテーブル 7 3 0 にしたがって、HPA に変換する。この変換された HPA もページを示す。ページの単位は、GPA ページと同じである。

40

【 0 1 8 2 】

ステップ S 1 7 0 9 において、VMM # 0 は、制御 I / F 部 8 1 0 に備わるステータスレジスタ 8 1 1 の有効データカウントビットから、FIFO にデータページ情報があるか否かを判別する。

【 0 1 8 3 】

ステータスレジスタ 8 1 1 の有効データカウントが「 0 」の場合（ステップ S 1 7 0 9 YES）、VMM # 0 は、FIFO が empty と判断し、処理をステップ S 1 7 1 0 に移行する。

【 0 1 8 4 】

50

ステップ S 1 7 1 0 において、VMM # 0 は、ステップ S 1 7 0 8 で変換された H P A が示すページの範囲に格納されているデータを、サーバ # 0 に備わるメモリ 6 2 0 から取得する。そして、VMM # 0 は、取得したデータを、V M # 0 の移動先である V M M # 1 に転送する。

【 0 1 8 5 】

一方、VMM # 1 は、VMM # 0 からデータを受信すると、受信したデータを図 1 5 に示したステップ S 1 5 0 1 の処理で確保した V M 領域に格納する。

VMM # 1 へのデータ転送が完了すると、VMM # 0 は、処理をステップ S 1 7 1 1 に移行する。そして、VMM # 0 は、ストップ・アンド・コピーを終了する（ステップ S 1 7 1 1 ）。

10

【 0 1 8 6 】

一方、ステップ S 1 7 0 9 において、ステータスレジスタ 8 1 1 の有効データカウンタが「 0 」でない場合（ステップ S 1 7 0 9 N O ）、VMM # 0 は、F I F O が e m p t y でないと判断し、処理をステップ S 1 7 1 2 に移行する。

【 0 1 8 7 】

ステップ S 1 7 1 2 において、VMM # 0 は、制御 I / F 部 8 1 0 に備わるコントロールレジスタ 8 1 3 のリードセレクトビットに「 0 」または「 1 」を設定することにより、ダーティ・ページ情報を読み出す F I F O を変更する。

【 0 1 8 8 】

例えば、現在選択中の F I F O が F I F O # 0 の場合、VMM # 0 は、コントロールレジスタ 8 1 3 のリードセレクトビットに「 1 」を設定することにより、ダーティ・ページ情報を読み出す F I F O を F I F O # 1 に変更する。

20

【 0 1 8 9 】

同様に、現在選択中の F I F O が F I F O # 1 の場合、VMM # 0 は、コントロールレジスタ 8 1 3 のリードセレクトビットに「 0 」を設定することにより、ダーティ・ページ情報を読み出す F I F O を F I F O # 0 に変更する。

ダーティ・ページ情報を読み出す F I F O を変更すると、VMM # 0 は、処理をステップ S 1 7 0 1 に移行する。

【 0 1 9 0 】

（パケット検出部 8 2 0 の具体例）

30

図 1 8 は、条件 1 を用いたパケット検出部 8 2 0 の具体例を示す図である。以下、条件 1 を用いたパケット検出部 8 2 0 を「パケット検出部 1 8 0 0」という。

【 0 1 9 1 】

パケット検出部 1 8 0 0 は、F I F O 制御部 9 0 1 と、判定部 1 8 1 0 と、を備える。

判定部 1 8 1 0 は、ページ記憶部 1 8 1 1 と、比較部 1 8 1 2 と、を備える。

ページ記憶部 1 8 1 1 は、DMA ライトパケット 8 5 0 のヘッダから取得した G P A ページを記憶するレジスタである。比較部 1 8 1 2 は、DMA ライトパケット 8 5 0 のヘッダから取得した G P A ページと、ページ記憶部 1 8 1 1 に記憶されている G P A ページと、を比較する。

【 0 1 9 2 】

40

DMA ライトパケット 8 5 0 のヘッダから取得した G P A ページと、ページ記憶部 1 8 1 1 に記憶されている G P A ページと、が一致しない場合、比較部 1 8 1 2 は、F I F O にダーティ・ページ情報を登録する必要があると判断する。この場合、比較部 1 8 1 2 は、F I F O 制御部 9 0 1 に対して、F I F O にダーティ・ページ情報を登録する必要がある旨を通知する。そして、比較部 1 8 1 2 は、DMA ライトパケット 8 5 0 のヘッダから取得した G P A ページをページ記憶部 1 8 1 1 に記憶する。

【 0 1 9 3 】

DMA ライトパケット 8 5 0 のヘッダから取得した G P A ページと、ページ記憶部 1 8 1 1 に記憶されている G P A ページと、が一致する場合、比較部 1 8 1 2 は、F I F O にダーティ・ページ情報を登録する必要があるないと判断する。この場合、比較部 1 8 1 2 は、

50

F I F O制御部 9 0 1 に対して、F I F Oにダーティ・ページ情報を登録する必要がない旨を通知する。

【 0 1 9 4 】

図 1 9 は、パケット検出部 1 8 0 0 の処理を示すフローチャートである。図 1 9 に示す処理は、図 1 2 に示したステップ S 1 2 0 4 ~ S 1 2 0 6 に対応する処理である。

図 1 2 に示したステップ S 1 2 0 4 に処理が移行すると、パケット検出部 1 8 0 0 は、以下の処理を開始する（ステップ S 1 9 0 0 ）。

【 0 1 9 5 】

ステップ S 1 9 0 1 において、パケット検出部 1 8 0 0 は、DMAライトパケットのヘッダから S I D と G P A を取得する。そして、パケット検出部 1 8 0 0 は、G P A から G P A ページを取得する。

10

【 0 1 9 6 】

ステップ S 1 9 0 2 において、パケット検出部 1 8 0 0 は、ステップ S 1 9 0 1 で取得した G P A ページと、ページ記憶部 1 8 1 1 に記憶されている G P A ページと、を比較する。

【 0 1 9 7 】

ステップ S 1 9 0 1 で取得した G P A ページと、ページ記憶部 1 8 1 1 に記憶されている G P A ページと、が一致する場合（ステップ S 1 9 0 2 Y E S ）、パケット検出部 1 8 0 0 は、ダーティ・ページ情報を F I F O に登録することなく、処理をステップ S 1 9 0 5 に移行する。

20

【 0 1 9 8 】

また、ステップ S 1 9 0 1 で取得した G P A ページと、ページ記憶部 1 8 1 1 に記憶されている G P A ページと、が一致しない場合（ステップ S 1 9 0 2 N O ）、パケット検出部 1 8 0 0 は、処理をステップ S 1 9 0 3 に移行する。この場合、パケット検出部 1 8 0 0 は、ステップ S 1 9 0 1 で取得した S I D と G P A ページを含むダーティ・ページ情報を F I F O に登録する（ステップ S 1 9 0 3 ）。

【 0 1 9 9 】

ステップ S 1 9 0 4 において、パケット検出部 1 8 0 0 は、ステップ S 1 9 0 1 で取得した G P A ページをページ記憶部 1 8 1 1 に記憶する。そして、パケット検出部 1 8 0 0 は、処理をステップ S 1 9 0 5 に移行する。

30

以上の処理が終了すると、パケット検出部 1 8 0 0 は、図 1 2 に示したステップ S 1 2 0 7 から処理を開始する。

【 0 2 0 0 】

（パケット検出部 8 2 0 のその他の具体例）

図 2 0 は、条件 2 を用いたパケット検出部 8 2 0 の具体例を示す図である。以下、条件 2 を用いたパケット検出部 8 2 0 を「パケット検出部 2 0 0 0」という。

【 0 2 0 1 】

パケット検出部 2 0 0 0 は、F I F O制御部 9 0 1 と、判定部 2 0 1 0 と、を備える。

判定部 2 0 1 0 は、S I D記憶部 2 0 1 1 と、比較部 2 0 1 2 と、を備える。

S I D記憶部 2 0 1 1 は、ライブ・マイグレーションの対象となるユニットの S I D を記憶するレジスタである。

40

【 0 2 0 2 】

制御 I / F 部 8 1 0 のコントロールレジスタ 8 1 3 に所定の値が入力されると、F I F O制御部 9 0 1 は、データレジスタ 8 1 2 に記憶されている S I D を、S I D記憶部 2 0 1 1 に記憶する。このようにして、S I D は、あらかじめ S I D記憶部 2 0 1 1 に記憶しておくことができる。

【 0 2 0 3 】

比較部 2 0 1 2 は、DMAライトパケット 8 5 0 のヘッダから取得した S I D と、S I D記憶部 2 0 1 1 に記憶されている S I D と、を比較する。

そして、DMAライトパケット 8 5 0 のヘッダから取得した S I D と、S I D記憶部 2

50

011に記憶されているSIDと、が一致しない場合、比較部2012は、FIFOにダーティ・ページ情報を登録する必要があると判断する。この場合、比較部2012は、FIFO制御部901に対して、FIFOにダーティ・ページ情報を登録する必要がある旨を通知する。そして、比較部2012は、DMAライトパケット850のヘッダから取得したSIDをSID記憶部2011に記憶する。

【0204】

また、DMAライトパケット850のヘッダから取得したSIDと、SID記憶部2011に記憶されているSIDと、が一致する場合、比較部2012は、FIFOにダーティ・ページ情報を登録する必要があると判断する。この場合、比較部2012は、FIFO制御部901に対して、FIFOにダーティ・ページ情報を登録する必要がある旨を通知する。

10

【0205】

図21は、パケット検出部2000の処理を示すフローチャートである。図21に示す処理は、図12に示したステップS1204~S1206に対応する処理である。

図12に示したステップS1204に処理が移行すると、パケット検出部2000は、以下の処理を開始する(ステップS2100)。

【0206】

ステップS2101において、パケット検出部2000は、DMAライトパケットのヘッダからSIDとGPAを取得する。また、パケット検出部2000は、GPAからGPAページを取得する。

20

【0207】

ステップS2102において、パケット検出部2000は、ステップS2101で取得したSIDと、SID記憶部2011に記憶されているSIDと、を比較する。

ステップS2101で取得したSIDと、SID記憶部2011に記憶されているSIDと、が異なる場合(ステップS2102 NO)、パケット検出部2000は、ダーティ・ページ情報をFIFOに登録することなく、処理をステップS2104に移行する。

【0208】

また、ステップS2101で取得したSIDと、SID記憶部2011に記憶されているSIDと、が一致する場合(ステップS2102 YES)、パケット検出部2000は、処理をステップS2103に移行する。この場合、パケット検出部2000は、ステップS2101で取得したSIDとGPAページを含むダーティ・ページ情報をFIFOに登録する(ステップS2103)。そして、パケット検出部2000は、処理をステップS2104に移行する。

30

以上の処理が終了すると、パケット検出部2000は、図12に示したステップS1207から処理を開始する。

【0209】

図22は、ライブ・マイグレーション時におけるSID記憶部2011へのSID設定処理を示すフローチャートである。なお、図22に示すステップS2202~SS2210は、図15に示したステップS1501~S1509と同様なので、説明は省略する。

【0210】

ステップS2201において、VMM#0は、ライブ・マイグレーションの指示を受けると、ライブ・マイグレーションによる移動の対象であるVM#0と接続しているデバイス、すなわち、移動対象デバイスのSIDを、SID記憶部2011に記憶する。

40

【0211】

VMM#0は、例えば、以下のようにして、移動対象デバイスのSIDを、SID記憶部2011に記憶する。

まず、VMM#0は、制御I/F部810のデータレジスタ812にSIDを設定する。そして、VMM#0は、コントロールレジスタ813に所定の値を入力する。すると、FIFO制御部901は、データレジスタ812に記憶されているSIDを、SID記憶部2011する。

50

【 0 2 1 2 】

ステップ S 2 2 0 1 ~ S 2 2 1 0 の処理が完了すると、VMM # 0 は、ステップ S 2 2 0 1 の処理で S I D 記憶部 2 0 1 1 に記憶した S I D をクリアする。この場合、例えば、VMM # 0 が、コントロールレジスタ 8 1 3 に所定の値を入力すると、F I F O 制御部 9 0 1 が、データレジスタ 8 1 2 に記憶されている S I D をクリアするようにすることができる。

【 0 2 1 3 】

ステップ S 2 2 1 2 において、VMM # 0 からメモリ・コピーの完了通知を受けると、VMM # 1 は、VMM # 1 上で V M # 0 の業務を開始させる。VMM # 1 上で V M # 0 が業務を再開すると、VMM # 1 は、処理をステップ S 2 2 1 3 に移行する。そして、V M M # 0 および # 1 は、ライブ・マイグレーションを終了する。

10

【 0 2 1 4 】

(パケット検出部 8 2 0 のその他の具体例)

図 2 3 は、条件 3 を用いたパケット検出部 8 2 0 の具体例を示す図である。以下、条件 3 を用いたパケット検出部 8 2 0 を「パケット検出部 2 3 0 0」という。

【 0 2 1 5 】

パケット検出部 2 3 0 0 を実現する場合、F I F O に登録するダーティ・ページ情報には、S I D と、G P A ページと、カウント値と、を含んだダーティ・ページ情報を使用する必要がある。カウント値は、ダーティ・ページ情報に含まれる G P A ページから何ページまで連続してアクセスがあったかを示す値である。カウント値を使用することにより、複数ページをまたぐ連続アクセスを、1 つのエントリで F I F O に記憶することができる。

20

【 0 2 1 6 】

パケット検出部 2 3 0 0 は、F I F O 制御部 2 3 1 0 と、判定部 2 3 2 0 と、を備える。

F I F O 制御部 2 3 1 0 は、図 9 で説明した機能に加えて、判定部 2 3 1 0 の指示にしたがって、最後に F I F O に登録されたダーティ・ページ情報に含まれるカウント値を、1 だけインクリメントした値に更新する。

【 0 2 1 7 】

判定部 2 3 2 0 は、ページ記憶部 2 3 2 1 と、比較部 2 3 2 2 と、比較部 2 3 2 3 と、を備える。

30

ページ記憶部 2 3 2 1 は、D M A ライトパケット 8 5 0 のヘッダから取得した G P A ページを記憶するレジスタである。

【 0 2 1 8 】

比較部 2 3 2 2 は、D M A ライトパケット 8 5 0 のヘッダから取得した G P A ページと、ページ記憶部 2 3 2 1 に記憶されている G P A ページと、を比較する。

D M A ライトパケット 8 5 0 のヘッダから取得した G P A ページと、ページ記憶部 2 3 2 1 に記憶されている G P A ページと、が一致する場合、比較部 2 3 2 2 は、F I F O にダーティ・ページ情報を登録する必要がないと判断する。この場合、比較部 2 3 2 2 は、F I F O 制御部 2 3 1 0 に対して、F I F O にダーティ・ページ情報を登録する必要がない旨を通知する。

40

【 0 2 1 9 】

また、D M A ライトパケット 8 5 0 のヘッダから取得した G P A ページと、ページ記憶部 2 3 2 1 に記憶されている G P A ページと、が一致しない場合、比較部 2 3 2 2 は、その旨を比較部 2 3 2 3 に通知する。

【 0 2 2 0 】

比較部 2 3 2 3 は、D M A ライトパケット 8 5 0 のヘッダから取得した G P A ページと、ページ記憶部 2 3 2 1 に記憶されている G P A ページを 1 ページだけインクリメントした値と、を比較する。

【 0 2 2 1 】

50

DMAライトパケット850のヘッダから取得したGPAページと、ページ記憶部2321に記憶されているGPAページに1ページだけインクリメントした値と、が一致する場合、比較部2323は、FIFOに登録したダーティ・ページ情報の更新が必要と判断する。この場合の更新とは、最後にFIFOに登録したダーティ・ページ情報に含まれるカウント値を1だけインクリメントすることである。比較部2323は、FIFOに登録したダーティ・ページ情報の更新をFIFO制御部2310に指示する。そして、比較部2323は、DMAライトパケット850のヘッダから取得したGPAページをページ記憶部2321に記憶する。

【0222】

また、DMAライトパケット850のヘッダから取得したGPAページと、ページ記憶部2321に記憶されているGPAページを1ページだけインクリメントした値と、が一致しない場合、比較部2323は、FIFOにダーティ・ページ情報を登録する必要があると判断する。ただし、比較部2322から、DMAライトパケット850のヘッダから取得したGPAページと、ページ記憶部2321に記憶されているGPAページと、が一致しない旨の通知を受けた場合に限る。この場合、比較部2323は、FIFO制御部2310に対して、FIFOにダーティ・ページ情報を登録する必要がある旨を通知する。また、比較部2323は、DMAライトパケット850のヘッダから取得したGPAページをページ記憶部2321に記憶する。

10

【0223】

図24は、パケット検出部2300の処理を示すフローチャートである。図24に示す処理は、図12に示したステップS1204～S1206に対応する処理である。

20

図12に示したステップS1204に処理が移行すると、パケット検出部2300は、以下の処理を開始する(ステップS2400)。

【0224】

ステップS2401において、パケット検出部2300は、DMAライトパケットのヘッダからSIDとGPAを取得する。そして、パケット検出部2300は、GPAからGPAページを取得する。

【0225】

ステップS2402において、パケット検出部2300は、ステップS2401で取得したGPAページと、ページ記憶部2321に記憶されているGPAページと、を比較する。

30

【0226】

ステップS2401で取得したGPAページと、ページ記憶部2321に記憶されているGPAページと、が一致する場合(ステップS2402 YES)、パケット検出部2300は、ダーティ・ページ情報をFIFOに登録することなく、処理をステップS2407に移行する。

【0227】

また、ステップS2401で取得したGPAページと、ページ記憶部2321に記憶されているGPAページと、が一致しない場合(ステップS2402 NO)、パケット検出部2300は、処理をステップS2403に移行する。この場合、パケット検出部2300は、パケット検出部2300は、ステップS2401で取得したGPAページと、ページ記憶部2321に記憶されているGPAページを1ページだけインクリメントした値と、を比較する。

40

【0228】

ステップS2401で取得したGPAページと、ページ記憶部2321に記憶されているGPAページを1ページだけインクリメントした値と、が一致する場合(ステップS2403 YES)、パケット検出部2300は、処理をステップS2404に移行する。この場合、パケット検出部2300は、FIFOに最後に登録されたダーティ・ページ情報に含まれるカウント値を1だけインクリメントする(ステップS2404)。そして、パケット検出部2300は、処理をステップS2406に移行する。

50

【0229】

また、ステップS2401で取得したGPAページと、ページ記憶部2321に記憶されているGPAページを1ページだけインクリメントした値と、が一致しない場合（ステップS2403 NO）、パケット検出部2300は、処理をステップS2405に移行する。この場合、パケット検出部2300は、ステップS2401で取得したSIDおよびGPAページと、初期値が0に設定されたカウント値と、を含むダーティ・ページ情報をFIFOに登録する（ステップS2405）。そして、パケット検出部2300は、処理をステップS2406に移行する。

【0230】

ステップS2406において、パケット検出部2300は、ステップS2401で取得したGPAページをページ記憶部2321に記憶する。そして、パケット検出部2300は、処理をステップS2407に移行する。

以上の処理が終了すると、パケット検出部2300は、図12に示したステップS1207から処理を開始する。

【0231】

（パケット検出部820のその他の具体例）

図25は、条件4を用いたパケット検出部820の具体例を示す図である。以下、条件4を用いたパケット検出部820を「パケット検出部2500」という。図25には、説明の都合で、FIFO#NにGPAページのみを記載しているが、図25に示す構成にFIFO#Nを限定する趣旨ではない。

【0232】

パケット検出部2500は、FIFO制御部901と、判定部2510と、を備える。

判定部2510は、キャッシュメモリ2511と、キャッシュ制御回路2512と、を備える。

【0233】

キャッシュメモリ2511は、DMAライトパケット850のヘッダから取得したGPAページを記憶するキャッシュメモリである。キャッシュメモリ2511には、複数のGPAページを格納することができる。なお、図25に示すキャッシュメモリ2511には、3つのGPAページX、YおよびZを記憶するキャッシュメモリを示しているが、図25に示す構成にキャッシュメモリ2511の構成を限定する趣旨ではない。

【0234】

キャッシュ制御回路2512は、キャッシュメモリ2511を参照し、DMAライトパケット850のヘッダから取得したGPAページと一致するGPAページを検索する。

DMAライトパケット850のヘッダから取得したGPAページと一致するGPAページを検出すると、キャッシュ制御回路2512は、FIFOにダーティ・ページ情報を登録する必要がないと判断する。この場合、キャッシュ制御回路2512は、FIFO制御部901に対して、FIFOにダーティ・ページ情報を登録する必要がない旨を通知する。

【0235】

また、DMAライトパケット850のヘッダから取得したGPAページと一致するGPAページを検出しない場合、キャッシュ制御回路2512は、FIFOにダーティ・ページ情報を登録する必要があると判断する。この場合、キャッシュ制御回路2512は、FIFO制御部901に対して、FIFOにダーティ・ページ情報を登録する必要がある旨を通知する。そして、キャッシュメモリ2511に空きエントリがあれば、キャッシュ制御回路2512は、空きエントリに、DMAライトパケット850のヘッダから取得したGPAページを記憶する。また、キャッシュメモリ2511に空きエントリがなければ、キャッシュ制御回路2512は、参照されていない期間が最も長いエントリに、DMAライトパケット850のヘッダから取得したGPAページを記憶する。

【0236】

図26は、パケット検出部2500の処理を示すフローチャートである。図26に示す

10

20

30

40

50

処理は、図 12 に示したステップ S 1204 ~ S 1206 に対応する処理である。

図 12 に示したステップ S 1204 に処理が移行すると、パケット検出部 2500 は、以下の処理を開始する（ステップ S 2600）。

【0237】

ステップ S 2601 において、パケット検出部 2500 は、DMA ライトパケットのヘッダから SID と GPA を取得する。そして、パケット検出部 2500 は、GPA から GPA ページを取得する。

【0238】

ステップ S 2602 において、パケット検出部 2500 は、キャッシュメモリ 2511 を参照し、ステップ S 2601 で取得した GPA ページと一致する GPA ページを検索する。

10

【0239】

ステップ S 2601 で取得した GPA ページと一致する GPA ページを検出した場合（ステップ S 2602 YES）、パケット検出部 2500 は、ダーティ・ページ情報を FIFO に登録することなく、処理をステップ S 2607 に移行する。

【0240】

また、ステップ S 2601 で取得した GPA ページと一致する GPA ページを検出しない場合（ステップ S 2602 NO）、パケット検出部 2500 は、処理をステップ S 2603 に移行する。この場合、パケット検出部 2500 は、キャッシュメモリ 2511 に空きエントリがあるか否かを判別する（ステップ S 2603）。

20

【0241】

キャッシュメモリ 2511 に空きエントリがある場合（ステップ S 2603 YES）、パケット検出部 2500 は、ステップ S 2601 で取得した GPA ページを、空きエントリに記憶する（ステップ S 2604）。そして、パケット検出部 2500 は、処理をステップ S 2606 に移行する。

【0242】

キャッシュメモリ 2511 に空きエントリがない場合（ステップ S 2603 NO）、パケット検出部 2500 は、処理をステップ S 2605 に移行する。この場合、パケット検出部 2500 は、ステップ S 2601 で取得した GPA ページを、キャッシュメモリ 2511 に含まれるエントリのうち、参照されていない期間が最も長いエントリに記憶する（ステップ S 2604）。そして、パケット検出部 2500 は、処理をステップ S 2606 に移行する。

30

【0243】

ステップ S 2606 において、パケット検出部 2500 は、ダーティ・ページ情報を FIFO に登録する。そして、パケット検出部 2500 は、処理をステップ S 2607 に移行する。

以上の処理が終了すると、パケット検出部 2500 は、図 12 に示したステップ S 1207 から処理を開始する。

【0244】

（パケット検出部 820 のその他の具体例）

40

図 27 は、条件 5 を用いたパケット検出部 820 の具体例を示す図である。以下、条件 5 を用いたパケット検出部 820 を「パケット検出部 2700」という。図 27 には、説明の都合で、FIFO # N には SID の記載を省略しているが、図 27 に示す構成に FIFO # N を限定する趣旨ではない。

【0245】

パケット検出部 2700 を実現する場合、FIFO に登録するダーティ・ページ情報には、SID と、GPA ページと、拡張サイズと、を含むダーティ・ページ情報を使用する。この拡張サイズは、ダーティ・ページ情報に含まれる GPA ページのページサイズの拡張に使用する情報である。GPA ページのページサイズを拡張することで、1つのダーティ・ページ情報に、より広い範囲の、変更されたメモリ領域を含ませることができる。し

50

たがって、図7に示したメモリ620の管理などのメモリシステムに使用するページサイズを拡張するものではない。拡張サイズの指定には、「0」、「1」、「2」、・・・などの数値を使用することができる。この場合、「0」は、ダーティ・ページ情報の含まれるGPAページのページサイズを4KBに拡張することを示している。同様に、「1」は、ダーティ・ページ情報に含まれるGPAページのページサイズを8KB、「2」は、ダーティ・ページ情報に含まれるGPAページのページサイズを16KBに拡張することを示している。

【0246】

パケット検出部2700は、FIFO制御部2710と、判定部2720と、を備える。

10

FIFO制御部2710は、図9で説明した機能に加えて、判定部2720の指示にしたがって、FIFOに登録されているダーティ・ページ情報に含まれる拡張サイズを変更する。

【0247】

判定部2720は、キャッシュメモリ2721と、キャッシュ制御回路2722と、を備える。

キャッシュメモリ2721は、DMAライトパケット850のヘッダから取得したGPAページを、その拡張サイズとともに記憶するキャッシュメモリである。キャッシュメモリ2721には、複数のGPAページを格納することができる。例えば、図27には、GPAページ「X」が拡張サイズの設定値「2」とともにキャッシュメモリ2721に記憶されている。これは、GPAページ「X」のページサイズを8KBに拡張することを示している。

20

【0248】

なお、図27に示すキャッシュメモリ2721には、3つのGPAページX、YおよびZを記憶するキャッシュメモリを示しているが、図27に示す構成にキャッシュメモリ2721の構成を限定する趣旨ではない。

【0249】

キャッシュ制御回路2722は、キャッシュメモリ2721を参照する。そして、キャッシュ制御回路2722は、DMAライトパケット850のヘッダから取得したGPAページの下位nビットをマスクした値と、キャッシュメモリ2721に記憶されているGPAページの下位nビットをマスクした値と、を比較する。そして、キャッシュ制御回路2722は、DMAライトパケット850のヘッダから取得したGPAページの下位nビットをマスクした値と、キャッシュメモリ2721に記憶されているGPAページの下位nビットをマスクした値と、が一致したGPAページを検出する。このとき、キャッシュ制御回路2722は、DMAライトパケット850のヘッダから取得したGPAページがキャッシュメモリ2721にヒットしたと判断する。nを大きくすると、キャッシュメモリ2721にヒットしたと判断されるGPAページの範囲が大きくなる。なお、nは1以上の整数とする。nは、拡張サイズとの関係を考慮して決定することができる。例えば、4KB、8KBおよび16KBを拡張サイズに使用する場合、nを「2」などとすることができる。

30

40

【0250】

DMAライトパケット850のヘッダから取得したGPAページがキャッシュメモリ2721にヒットすると、キャッシュ制御回路2722は、DMAライトパケット850のヘッダから取得したGPAページと、キャッシュメモリ2721にヒットしたGPAページと、の差を求める。求めた差が、キャッシュメモリ2721にヒットしたGPAページの拡張サイズ以上に大きい場合、キャッシュ制御回路2722は、キャッシュメモリ2721にヒットしたGPAページの拡張サイズをより大きいページサイズに更新する。同様に、キャッシュ制御回路2722は、FIFOに記憶されているダーティ・ページ情報であって、キャッシュメモリ2721にヒットしたGPAページを含むダーティ・ページ情報に含まれる拡張サイズも更新する。

50

【0251】

また、DMAライトパケット850のヘッダから取得したGPAページがキャッシュメモリ2721にヒットしない場合、キャッシュ制御回路2722は、DMAライトパケット850のヘッダから取得したGPAページをキャッシュメモリ2721に記憶する。このとき、キャッシュ制御回路2722は、DMAライトパケット850のヘッダから取得したGPAページと、拡張サイズと、を関連付けて記憶する。本実施例では、拡張サイズの初期値として4KBを用いるので、拡張サイズの設定値として「0」が設定される。

【0252】

図28は、パケット検出部2700の処理を示すフローチャートである。図28に示す処理は、図12に示したステップS1204～S1206に対応する処理である。

図12に示したステップS1204に処理が移行すると、パケット検出部2700は、以下の処理を開始する(ステップS2800)。

【0253】

ステップS2801において、パケット検出部2700は、DMAライトパケット850のヘッダからSIDとGPAを取得する。そして、パケット検出部2700は、GPAからGPAページを取得する。

【0254】

ステップS2802において、パケット検出部2700は、ステップS2801で取得したGPAページの下位nビットをマスクした値と、キャッシュメモリ2721に記憶されているGPAページの下位nビットをマスクした値と、を比較する。そして、パケット検出部2700は、ステップS2801で取得したGPAページの下位nビットをマスクした値と、キャッシュメモリ2721に記憶されているGPAページの下位nビットをマスクした値と、が一致したGPAページを検出する。このとき、パケット検出部2700は、ステップS2801で取得したGPAページがキャッシュメモリ2721にヒットしたと判断する。

【0255】

ステップS2801で取得したGPAページがキャッシュメモリ2721にヒットした場合(ステップS2802 YES)、パケット検出部2700は、処理をステップS2803に移行する。この場合、パケット検出部2700は、ステップS2801で取得したGPAページと、キャッシュメモリ2721にヒットしたGPAページと、の差を算出する(ステップS2803)。

【0256】

ステップS2801で取得したGPAページと、キャッシュメモリ2721にヒットしたGPAページと、の差が、キャッシュメモリ2721にヒットしたGPAページの拡張サイズより小さい場合(ステップS2804 YES)、パケット検出部2700は、処理をステップS2811に移行する。

【0257】

ステップS2801で取得したGPAページと、キャッシュメモリ2721にヒットしたGPAページと、の差が、キャッシュメモリ2721にヒットしたGPAページの拡張サイズ以上の場合(ステップS2804 NO)、パケット検出部2700は、処理をステップS2805に移行する。この場合、パケット検出部2700は、ヒットしたGPAページと関連付けられてキャッシュメモリ2721に記憶されている拡張サイズを、ステップS2803で算出した差より大きい拡張サイズに更新する(ステップS2805)。本実施例では、4KB、8KB、16KB、・・・を拡張サイズに使用する。したがって、例えば、ステップS2803で算出した差が7KBの場合、拡張サイズを8KBに更新する。

【0258】

ステップS2806において、パケット検出部2700は、FIFOに記憶されている、キャッシュメモリ2721にヒットしたGPAページを含むダーティ・ページ情報に含まれる拡張サイズを、ステップS2805で更新した拡張サイズと同じ拡張サイズに更新

10

20

30

40

50

する。そして、パケット検出部 2700 は、処理をステップ S 2811 に移行する。

【0259】

一方、ステップ S 2801 で取得した GPA ページがキャッシュメモリ 2721 にヒットしない場合（ステップ S 2802 NO）、パケット検出部 2700 は、処理をステップ S 2807 に移行する。この場合、パケット検出部 2700 は、キャッシュメモリ 2721 に空きエントリがあるか否かを判別する（ステップ S 2807）。

【0260】

キャッシュメモリ 2721 に空きエントリがある場合（ステップ S 2807 YES）、パケット検出部 2700 は、ステップ S 2801 で取得した GPA ページを空きエントリに記憶する（ステップ S 2808）。このとき、パケット検出部 2700 は、ステップ S 2801 で取得した GPA ページと、初期値として 4KB に設定された拡張サイズと、を関連付けて記憶する。そして、パケット検出部 2700 は、処理をステップ S 2810 に移行する。

10

【0261】

また、キャッシュメモリ 2721 に空きエントリがない場合（ステップ S 2807 NO）、パケット検出部 2700 は、処理をステップ S 2809 に移行する。この場合、パケット検出部 2700 は、ステップ S 2801 で取得した GPA ページを、キャッシュメモリ 2721 に含まれるエントリのうち、参照されていない期間が最も長いエントリに記憶する（ステップ S 2809）。この場合も、パケット検出部 2700 は、ステップ S 2801 で取得した GPA ページと、初期値として 4KB に設定された拡張サイズと、を関連付けて記憶する。そして、パケット検出部 2700 は、処理をステップ S 2810 に移行する。

20

【0262】

ステップ S 2810 において、パケット検出部 2700 は、ダーティ・ページ情報を FIFO に登録する。このとき、パケット検出部 2700 は、ダーティ・ページ情報に含まれる拡張サイズに、初期値として 4KB を設定する。そして、パケット検出部 2700 は、処理をステップ S 2811 に移行する。

以上の処理が終了すると、パケット検出部 2700 は、図 12 に示したステップ S 1207 から処理を開始する。

【0263】

30

（パケット検出部 820 のその他の具体例）

図 29 は、条件 1 および 2 を用いたパケット検出部 820 の具体例を示す図である。以下、条件 1 および 2 を用いたパケット検出部 820 を「パケット検出部 2900」という。

【0264】

パケット検出部 2900 は、FIFO 制御部 2910 と、判定部 2920 と、を備える。また、判定部 2920 は、図 18 に示した判定部 1810 と、図 20 に示した判定部 2010 と、を備える。

【0265】

FIFO 制御部 2910 は、判定部 2010 から、DMA ライトパケット 850 のヘッダから取得した SID と、SID 記憶部 2011 に記憶されている SID と、が一致しない旨の通知を受ける。この場合、FIFO 制御部 2910 は、図 18 で説明したように、判定部 1810 の判定結果にしたがって、FIFO にダーティ・ページ情報を登録する。

40

【0266】

図 30 は、パケット検出部 2900 の処理を示すフローチャートである。図 30 に示す処理は図 12 に示したステップ S 1204 ~ S 1206 に対応する。

図 12 に示したステップ S 1204 に処理が移行すると、パケット検出部 2900 は、以下の処理を開始する（ステップ S 3000）。

【0267】

ステップ S 3001 において、パケット検出部 2900 は、DMA ライトパケットのへ

50

ッダからSIDとGPAを取得する。そして、パケット検出部2900は、GPAからGPAページを取得する。

【0268】

ステップS3002において、パケット検出部2900は、ステップS3001で取得したSIDと、SID記憶部2011に記憶されているSIDと、を比較する。

ステップS3001で取得したSIDと、SID記憶部2011に記憶されているSIDと、が異なる場合(ステップS3002 NO)、パケット検出部2900は、ダーティ・ページ情報をFIFOに登録することなく、処理をステップS3006に移行する。

【0269】

また、ステップS3001で取得したSIDと、SID記憶部2011に記憶されているSIDと、が一致する場合(ステップS3002 YES)、パケット検出部2900は、処理をステップS3003に移行する。ステップS3003～S3006の処理は、図19に示したステップS1902～S1905と同じなので説明を省略する。

以上の処理が終了すると、パケット検出部2900は、図12に示したステップS1207から処理を開始する。

【0270】

(パケット検出部820のその他の具体例)

図31は、条件2および3を用いたパケット検出部820の具体例を示す図である。以下、条件2および3を用いたパケット検出部820を「パケット検出部3100」という。

【0271】

パケット検出部3100は、FIFO制御部3110と、判定部3120と、を備える。また、判定部3120は、図20に示した判定部2010と、図23に示した判定部2320と、を備える。

【0272】

FIFO制御部3110は、判定部2010から、DMAライトパケット850のヘッダから取得したSIDと、SID記憶部2011に記憶されているSIDと、が一致しない旨の通知を受ける。この場合、FIFO制御部3110は、図23で説明したように、判定部2320の判定結果にしたがって、FIFOにダーティ・ページ情報を登録する。

【0273】

図32は、パケット検出部3100の処理を示すフローチャートである。図32に示す処理は図12に示したステップS1204～S1206に対応する。

図12に示したステップS1204に処理が移行すると、パケット検出部3100は、以下の処理を開始する(ステップS3200)。

【0274】

ステップS3201において、パケット検出部3100は、DMAライトパケットのヘッダからSIDとGPAを取得する。そして、パケット検出部3100は、GPAからGPAページを取得する。

【0275】

ステップS3202において、パケット検出部3100は、ステップS3201で取得したSIDと、SID記憶部2011に記憶されているSIDと、を比較する。

ステップS3201で取得したSIDと、SID記憶部2011に記憶されているSIDと、が異なる場合(ステップS3202 NO)、パケット検出部3100は、ダーティ・ページ情報をFIFOに登録することなく、処理をステップS3208に移行する。

【0276】

また、ステップS3201で取得したSIDと、SID記憶部2011に記憶されているSIDと、が一致する場合(ステップS3202 YES)、パケット検出部3100は、処理をステップS3203に移行する。ステップS3203～S3208の処理は、図24に示したステップS2402～S2407と同じなので説明を省略する。

以上の処理が終了すると、パケット検出部3100は、図12に示したステップS12

10

20

30

40

50

07から処理を開始する。

【0277】

(パケット検出部820のその他の具体例)

図33は、条件2および4を用いたパケット検出部820の具体例を示す図である。以下、条件2および4を用いたパケット検出部820を「パケット検出部3300」という。

【0278】

パケット検出部3300は、FIFO制御部3310と、判定部3320と、を備える。また、判定部3320は、図20に示した判定部2010と、図25に示した判定部2510と、を備える。

10

【0279】

FIFO制御部3310は、判定部2010から、DMAライトパケット850のヘッダから取得したSIDと、SID記憶部2011に記憶されているSIDと、が一致しない旨の通知を受ける。この場合、FIFO制御部3310は、図25で説明したように、判定部2510の判定結果にしたがって、FIFOにダーティ・ページ情報を登録する。

【0280】

図34は、パケット検出部3300の処理を示すフローチャートである。図34に示す処理は図12に示したステップS1204～S1206に対応する。

図12に示したステップS1204に処理が移行すると、パケット検出部3300は、以下の処理を開始する(ステップS3400)。

20

【0281】

ステップS3401において、パケット検出部3300は、DMAライトパケットのヘッダからSIDとGPAを取得する。そして、パケット検出部3300は、GPAからGPAページを取得する。

【0282】

ステップS3402において、パケット検出部3300は、ステップS3401で取得したSIDと、SID記憶部2011に記憶されているSIDと、を比較する。

ステップS3401で取得したSIDと、SID記憶部2011に記憶されているSIDと、が異なる場合(ステップS3402 NO)、パケット検出部3300は、ダーティ・ページ情報をFIFOに登録することなく、処理をステップS3408に移行する。

30

【0283】

また、ステップS3401で取得したSIDと、SID記憶部2011に記憶されているSIDと、が一致する場合(ステップS3402 YES)、パケット検出部3300は、処理をステップS3403に移行する。ステップS3403～S3408の処理は、図26に示したステップS2602～S2607と同じなので説明を省略する。

以上の処理が終了すると、パケット検出部3300は、図12に示したステップS1207から処理を開始する。

【0284】

(パケット検出部820のその他の具体例)

図35は、条件2および5を用いたパケット検出部820の具体例を示す図である。以下、条件2および5を用いたパケット検出部820を「パケット検出部3500」という。図35には、説明を簡単にするために、FIFO#NにはSIDの記載を省略しているが、図35に示す構成にFIFO#Nを限定する趣旨ではない。

40

【0285】

パケット検出部3500は、FIFO制御部3510と、判定部3520と、を備える。また、判定部3520は、図20に示した判定部2010と、図27に示した判定部2720と、を備える。

【0286】

FIFO制御部3510は、判定部2010から、DMAライトパケット850のヘッダから取得したSIDと、SID記憶部2011に記憶されているSIDと、が一致しな

50

い旨の通知を受ける。この場合、F I F O制御部 3 5 1 0 は、図 2 7 で説明したように、判定部 2 7 2 0 の判定結果にしたがって、F I F O にダーティ・ページ情報を登録する。

【 0 2 8 7 】

図 3 6 は、パケット検出部 3 5 0 0 の処理を示すフローチャートである。図 3 6 に示す処理は図 1 2 に示したステップ S 1 2 0 4 ~ S 1 2 0 6 に対応する。

図 1 2 に示したステップ S 1 2 0 4 に処理が移行すると、パケット検出部 3 5 0 0 は、以下の処理を開始する（ステップ S 3 6 0 0 ）。

【 0 2 8 8 】

ステップ S 3 6 0 1 において、パケット検出部 3 5 0 0 は、DMA ライトパケットのヘッダから S I D と G P A を取得する。そして、パケット検出部 3 5 0 0 は、G P A から G P A ページを取得する。

10

【 0 2 8 9 】

ステップ S 3 6 0 2 において、パケット検出部 3 5 0 0 は、ステップ S 3 6 0 1 で取得した S I D と、S I D 記憶部 2 0 1 1 に記憶されている S I D と、を比較する。

ステップ S 3 6 0 1 で取得した S I D と、S I D 記憶部 2 0 1 1 に記憶されている S I D と、が異なる場合（ステップ S 3 6 0 2 N O ）、パケット検出部 3 5 0 0 は、ダーティ・ページ情報を F I F O に登録することなく、処理をステップ S 3 6 1 2 に移行する。

【 0 2 9 0 】

また、ステップ S 3 6 0 1 で取得した S I D と、S I D 記憶部 2 0 1 1 に記憶されている S I D と、が一致する場合（ステップ S 3 6 0 2 Y E S ）、パケット検出部 3 5 0 0 は、処理をステップ S 3 6 0 3 に移行する。ステップ S 3 6 0 3 ~ S 3 6 1 2 の処理は、図 2 8 に示したステップ S 2 8 0 2 ~ S 2 6 1 1 と同じなので説明を省略する。

20

【 0 2 9 1 】

以上の処理が終了すると、パケット検出部 3 5 0 0 は、図 1 2 に示したステップ S 1 2 0 7 から処理を開始する。

なお、条件 3 を用いたパケット検出部 8 2 0、例えば、パケット検出部 2 3 0 0 やパケット検出部 3 1 0 0 を実現する場合、F I F O には、図 3 7 に示す構成の F I F O を用いることができる。図 3 7 に示す F I F O は、図 1 0 に示した F I F O に含まれる各エントリに、カウント値を記憶する 8 ビットの領域を追加したものである。ただし、図 3 7 に示す F I F O の構成は、単なる一例であって、F I F O の構成を図 3 7 に示す構成に限定する趣旨ではない。

30

【 0 2 9 2 】

また、条件 5 を用いたパケット検出部 8 2 0、例えば、パケット検出部 2 7 0 0 やパケット検出部 3 5 0 0 を実現する場合、F I F O には、図 3 8 に示す構成の F I F O を用いることができる。図 3 8 に示す F I F O は、図 1 0 に示した F I F O に含まれる各エントリに、拡張サイズを記憶する 4 ビットの領域を追加したものである。ただし、図 3 8 に示す F I F O の構成は、単なる一例であって、F I F O の構成を図 3 8 に示す構成に限定する趣旨ではない。

【 0 2 9 3 】

以上に説明したように、ダーティ・ページ管理ユニット 6 5 1 は、I O アダプタ 6 4 0、6 4 1 および 6 4 2 からメモリ 6 2 0 に転送される DMA パケットが P C I e スイッチ 6 5 0 を通過すると、その DMA パケットをパケット検出部 8 2 0 によって検出する。

40

【 0 2 9 4 】

そして、検出した DMA パケットが DMA ライトパケットである場合、ダーティ・ページ管理ユニット 6 5 1 は、DMA パケットのヘッダからゲスト物理メモリと S I D を取得し、ダーティ・ページ情報として F I F O # 0 または # 1 に登録する。

【 0 2 9 5 】

また、ダーティ・ページ管理ユニット 6 5 1 は、C P U 6 1 0 または 6 1 1、すなわち C P U 6 1 0 または 6 1 1 が実現する V M M との制御 I / F 部 8 1 0 を備え、V M M からの要求に応じて F I F O # 0 または # 1 に登録されているダーティ・ページ情報を出力す

50

る。

【0296】

その結果、情報処理装置600で動作しているVMMは、ダーティ・ページ情報から、IOMMU631を利用したDMAライトによって追加、変更または更新等されたメモリ620上のデータを取得することが可能となる。

【0297】

そして、DMAライトによって追加、変更または更新等されたメモリ620上のデータを取得することが可能となったことにより、図15～図17に示したプレ・コピーおよびストップ・アンド・コピーが可能となる。その結果、IOMMUを使用する場合であってもライブ・マイグレーションを実行することが可能となる。

10

【0298】

また、情報処理装置600は、ダーティ・ページ記憶部652として、FIFO#0および#1の2つのFIFOを備える。そして、一方のFIFOが書込み対象の場合、他方のFIFOを読み出し対象とすることにより、1つのFIFOに対して読み込み処理と書込み処理が競合することを回避することが可能となる。その結果、FIFOへのダーティ・ページ情報の登録処理と、FIFOからのダーティ・ページ情報の読み出し処理と、を効率的に行なうことができる。

【0299】

また、パケット検出部1800は、最後にFIFOに記憶した、ダーティ・ページ情報に含まれるGPAページを、ページ記憶部1811に保持する。そして、パケット検出部1800は、ページ記憶部1811に記憶されているGPAページと一致しないGPAページをDMAライトパケット850から取得した場合にだけ、FIFOにダーティ・ページ情報を記憶する。

20

【0300】

これにより、パケット検出部1800は、DMAライトが同一ページに対して連続して行われている間の、FIFOの消費を抑えることができる。

例えば、大容量データの転送では、複数のDMAアクセスが連続アドレスに対して行われる。この場合、一般に、DMAアクセスの範囲は、ページサイズより小さいので、同一ページに対するDMAアクセスが連続して発生することになる。このような場合であっても、パケット検出部1800は、FIFOの消費量を抑えることが可能となる。

30

【0301】

また、パケット検出部2000は、ライブ・マイグレーションの対象となるユニットのSIDをSID記憶部2011に保持する。そして、パケット検出部2000は、SID記憶部2011に記憶されているSIDと一致するSIDが設定されたDMAパケットを検出した場合だけ、FIFOにダーティ・ページ情報を記憶する。

【0302】

これにより、パケット検出部2000は、複数のユニットから同時にDMAアクセスが行われた場合でも、ライブ・マイグレーションの対象となるユニットからのDMAアクセスだけを抽出することができる。その結果、パケット検出部2000は、ライブ・マイグレーションの対象となるユニットからのDMAアクセスについてのダーティ・ページ情報だけをFIFOに記憶することができる。この場合、パケット検出部2000は、ライブ・マイグレーションの対象以外のユニットからのDMAアクセスがあった場合には、FIFOが消費されないことになる。

40

【0303】

また、パケット検出部2300は、SIDと、GPAページと、そのGPAページから何ページDMAアクセスが連続して行われたかを示すカウント値と、を含むダーティ・ページ情報をFIFOに記憶する。

【0304】

そして、パケット検出部2300は、最後にFIFOに記憶されたダーティ・ページ情報に含まれるGPAページにカウント値を加えたページの次のページへのDMAアクセス

50

を検出する。この場合、パケット検出部 2300 は、新たなダーティ・ページ情報を F I F O に記憶することはせずに、ダーティ・ページ情報に含まれるカウント値をインクリメントする。

【0305】

これにより、パケット検出部 2300 は、DMA ライトが複数のページ境界をまたいで行われる連続アクセスの場合でも、F I F O の消費を抑えることができる。

例えば、I / O アダプタに割当てられるメモリ上の DMA 空間は、一般に、連続した複数のページに割り当てられる。このような DMA 空間に対して複数のページ境界をまたいで連続アクセスが行われる場合であっても、パケット検出部 2300 は、F I F O の消費量を抑えることが可能となる。

10

【0306】

また、パケット検出部 2500 は、F I F O に記憶した複数のダーティ・ページ情報それぞれに含まれる G P A ページを、キャッシュメモリ 2511 に保持する。そして、パケット検出部 2500 は、キャッシュメモリ 2511 に記憶されているいずれの G P A ページとも一致しない G P A ページを、DMA ライトパケット 850 から取得した場合にだけ、F I F O にダーティ・ページ情報を記憶する。

【0307】

これにより、パケット検出部 2500 は、複数の I / O アダプタを介して、複数の DMA アクセスが行われる場合でも、それぞれの DMA アクセスについて、DMA ライトが同一ページに連続して行われている間の、F I F O の消費を抑えることができる。

20

【0308】

DMA アクセスには局所性がある。特に G P G P U (G e n e r a l P u r p o s e c o m p u t i n g o n G r a p h i c s P r o c e s s i n g U n i t s) などの演算装置によるメモリ転送では、プログラムの性質により局所性が顕著になる場合が多い。G P G P U など複数のユニットが、各ユニットと接続する I / O アダプタを介して DMA アクセスを行う場合でも、パケット検出部 2500 は、それぞれの DMA アクセスについて、DMA ライトが同一ページに連続して行われている間の、F I F O の消費を抑えられる。

【0309】

また、DMA ライトパケット 850 から取得した G P A ページが、キャッシュメモリ 2721 に記憶されている G P A ページと一定範囲内にある場合、パケット検出部 2700 は、ダーティ・ページ情報を F I F O に記憶しない。

30

【0310】

ただし、パケット検出部 2700 は、DMA ライトパケット 850 から取得した G P A ページと、キャッシュメモリ 2721 に記憶されている G P A ページと、の差が、キャッシュメモリ 2721 に記憶されている G P A ページの拡張サイズ以上の場合を検出する。この場合、パケット検出部 2700 は、キャッシュメモリ 2721 に記憶されている G P A ページの拡張サイズを、DMA ライトパケット 850 から取得した G P A ページが含まれるサイズに更新する。そして、この拡張サイズの更新を、パケット検出部 2700 は、F I F O に記憶されているダーティ・ページ情報にも反映する。

40

【0311】

このように、パケット検出部 2700 は、DMA ライトパケット 850 から取得した G P A ページが、F I F O に記憶したダーティ・ページ情報に含まれる G P A ページと、一定範囲内にある場合には、ダーティ・ページ情報を F I F O に記憶しない。その結果、パケット検出部 2700 は、F I F O に記憶しているダーティ・ページ情報に含まれる G P A ページを含む一定範囲内で DMA ライトが行われる場合には、F I F O の消費を抑えることができる。

【0312】

また、DMA ライトパケット 850 から取得した G P A ページと、キャッシュメモリ 2721 に記憶されている G P A ページと、の差に応じて、拡張サイズが更新されるので、

50

G P G P U などアクセスパターンが不規則な D M A アクセスでも、F I F O の消費を抑える。

【 0 3 1 3 】

また、パケット検出部 2 9 0 0 は、ライブ・マイグレーションの対象となるユニットからの D M A アクセスについてのダーティ・ページ情報だけを F I F O への記憶の対象とする。そして、パケット検出部 2 9 0 0 は、ページ記憶部 1 8 1 1 に記憶されている G P A ページと一致しない G P A ページを D M A ライトパケット 8 5 0 から取得した場合にだけ、F I F O にダーティ・ページ情報を記憶する。

【 0 3 1 4 】

これにより、ライブ・マイグレーションの対象となるユニット以外のユニットからの D M A ライトが行われる場合、パケット検出部 2 9 0 0 は、F I F O の消費を抑えることができる。また、ライブ・マイグレーションの対象となるユニットからの D M A ライトが、同一ページに対して連続して行われない間、パケット検出部 2 9 0 0 は、F I F O の消費を抑えることができる。

10

【 0 3 1 5 】

また、パケット検出部 3 1 0 0 は、ライブ・マイグレーションの対象となるユニットからの D M A アクセスについてのダーティ・ページ情報だけを F I F O への記憶の対象とする。ただし、パケット検出部 3 1 0 0 は、最後に F I F O に記憶されたダーティ・ページ情報に含まれる G P A ページにカウント値を加えたページの次のページへの D M A アクセスを検出する。このとき、パケット検出部 3 1 0 0 は、新たなダーティ・ページ情報を F I F O に記憶することはせずに、ダーティ・ページ情報に含まれるカウント値をインクリメントする。

20

【 0 3 1 6 】

これにより、ライブ・マイグレーションの対象となるユニット以外のユニットからの D M A ライトが行われる場合、パケット検出部 3 1 0 0 は、F I F O の消費を抑えることができる。また、ライブ・マイグレーションの対象となるユニットからの D M A ライトが複数のページ境界をまたいで行われる連続アクセスの場合、パケット検出部 3 1 0 0 は、F I F O の消費を抑えることができる。

【 0 3 1 7 】

また、パケット検出部 3 3 0 0 は、ライブ・マイグレーションの対象となるユニットからの D M A アクセスについてのダーティ・ページ情報だけを F I F O への記憶の対象とする。そして、パケット検出部 3 3 0 0 は、キャッシュメモリ 2 5 1 1 に記憶されているいずれの G P A ページとも一致しない G P A ページを、D M A ライトパケット 8 5 0 から取得した場合にだけ、F I F O にダーティ・ページ情報を記憶する。

30

【 0 3 1 8 】

これにより、ライブ・マイグレーションの対象となるユニット以外のユニットからの D M A ライトが行われる場合、パケット検出部 3 3 0 0 は、F I F O の消費を抑えることができる。また、ライブ・マイグレーションの対象となる複数のユニットそれぞれから、複数の D M A アクセスが行われる場合でも、それぞれの D M A アクセスについて、D M A ライトが同一ページに連続して行われている間、パケット検出部 3 3 0 0 は、F I F O の消費を抑える。

40

【 0 3 1 9 】

また、パケット検出部 3 5 0 0 は、ライブ・マイグレーションの対象となるユニットからの D M A アクセスについてのダーティ・ページ情報だけを F I F O への記憶の対象とする。ただし、D M A ライトパケット 8 5 0 から取得した G P A ページが、キャッシュメモリ 2 7 2 1 に記憶されている G P A ページと一定範囲内にある場合、パケット検出部 3 5 0 0 は、ダーティ・ページ情報を F I F O に記憶しない。

【 0 3 2 0 】

これにより、ライブ・マイグレーションの対象となるユニット以外のユニットからの D M A ライトが行われる場合、パケット検出部 3 5 0 0 は、F I F O の消費を抑えることができる。

50

できる。また、ライブ・マイグレーションの対象となるユニットからのDMAライトが、FIFOに記憶しているダーティ・ページ情報に含まれるGPAページを含む一定範囲内で行われる場合、パケット検出部3500は、FIFOの消費を抑えることができる。

【0321】

以上に説明したように、本実施例に係るパケット検出部820を使用すると、情報処理装置600は、IOMMUを使用する場合であっても、少ないハードウェア資源でライブ・マイグレーションを実行することが可能となる。

【0322】

以上の実施例を含む実施形態に関し、さらに以下の付記を開示する。

(付記1)

10

ハードウェア資源を管理することにより、1以上の仮想機械を実現する仮想機械実現手段と、

前記仮想機械に割り当てられる第1のメモリ領域のアドレスと、前記第1のメモリ領域の実メモリである第2のメモリ領域のアドレスと、を相互に変換することにより、外部装置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送を行なうデータ転送手段と、

前記入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送されるデータを検出する検出手段と、

検出したデータが一定の条件を満たす場合に、該検出したデータにより変更される前記第1のメモリ領域に関する更新情報を生成し、前記更新情報を第1記憶部に記憶する登録手段と、

20

前記第1記憶部に記憶される更新情報を出力する出力手段と、

を備えることを特徴とする情報処理装置。

(付記2)

検出したデータにより変更される前記第1のメモリ領域と、最後に変更された前記第1のメモリ領域を記憶する第2記憶部に記憶した前記第1のメモリ領域と、が異なる場合、前記更新情報を生成して該更新情報を前記第1記憶部に記憶するとともに、前記変更される前記第1のメモリ領域を前記第2記憶部に記憶する、

ことを特徴とする付記1に記載の情報処理装置。

(付記3)

30

前記登録手段は、前記検出手段により検出したデータを出力した前記外部装置と、前記外部装置を識別する識別情報を記憶する外部装置記憶手段に記憶された識別情報が示す前記外部装置と、が同じ場合、前記更新情報を生成して該更新情報を前記第1記憶部に記憶する、

ことを特徴とする付記1に記載の情報処理装置。

(付記4)

前記登録手段は、前記検出手段により検出したデータにより変更される前記第1のメモリ領域を含むページと、前記第2記憶部に記憶した前記第1のメモリ領域を含むページの次のページと、が異なる場合、前記更新情報を生成して該更新情報を前記第1記憶部に記憶するとともに、前記変更される前記第1のメモリ領域を前記第2記憶部に記憶する、

40

ことを特徴とする付記2に記載の情報処理装置。

(付記5)

前記更新情報は、前記検出手段により検出したデータにより変更された前記第1のメモリ領域のページ数を計数した計数情報を含み、

前記登録手段は、前記検出手段により検出したデータにより変更される前記第1のメモリ領域を含むページと、前記第2記憶部に記憶した前記第1のメモリ領域を含むページの次のページと、が同一の場合、前記更新情報に含まれる計数情報が示す計数値に1を加算するとともに、前記変更される前記第1のメモリ領域を前記第2記憶部に記憶する、

ことを特徴とする付記2に記載の情報処理装置。

(付記6)

50

前記登録手段は、過去に変更された複数の前記第 1 のメモリ領域を前記第 2 記憶部に記憶し、前記検出手段により検出したデータにより変更される前記第 1 のメモリ領域が、前記第 2 記憶部に記憶した前記第 1 のメモリ領域のいずれかとも異なる場合に、前記更新情報を生成して該更新情報を前記第 1 記憶部に記憶するとともに、前記変更される前記第 1 のメモリ領域を前記第 2 記憶部に記憶する、

ことを特徴とする付記 2 に記載の情報処理装置。

(付記 7)

検出したデータにより変更される前記第 1 のメモリ領域と一定範囲内の前記第 1 のメモリ領域が、前記第 2 記憶部に含まれない場合、前記更新情報を生成して該更新情報を前記第 1 記憶部に記憶するとともに、前記変更される前記第 1 のメモリ領域を前記第 2 記憶部に記憶する、

ことを特徴とする付記 6 に記載の情報処理装置。

(付記 8)

前記更新情報は、変更された前記第 1 のメモリ領域を含むページのサイズを拡張する拡張情報を含み、

検出したデータにより変更される前記第 1 のメモリ領域と一定範囲内の前記第 1 のメモリ領域が前記第 2 記憶部に含まれる場合、該第 2 記憶部に含まれる前記第 1 のメモリ領域についての前記更新情報の拡張情報を、前記検出手段により検出したデータにより変更される前記第 1 のメモリ領域を含むページのサイズに更新する、

ことを特徴とする付記 6 に記載の情報処理装置。

(付記 9)

前記登録手段は、前記検出手段により検出したデータを出力した前記外部装置と、前記外部装置を識別する識別情報を記憶する外部装置記憶手段に記憶された識別情報が示す前記外部装置と、が同じ場合であって、前記検出手段により検出したデータにより変更される前記第 1 のメモリ領域と、最後に変更された前記第 1 のメモリ領域を記憶する第 2 記憶部に記憶した前記第 1 のメモリ領域と、が異なる場合、前記更新情報を生成して該更新情報を前記第 1 記憶部に記憶するとともに、前記変更される前記第 1 のメモリ領域を前記第 2 記憶部に記憶する、

ことを特徴とする付記 1 に記載の情報処理装置。

(付記 10)

前記登録手段は、前記検出手段により検出したデータを出力した前記外部装置と、前記外部装置を識別する識別情報を記憶する外部装置記憶手段に記憶された識別情報が示す前記外部装置と、が同じ場合であって、前記検出手段により検出したデータにより変更される前記第 1 のメモリ領域を含むページと、前記第 2 記憶部に記憶した前記第 1 のメモリ領域を含むページの次のページと、が異なる場合、前記更新情報を生成して該更新情報を前記第 1 記憶部に記憶するとともに、前記変更される前記第 1 のメモリ領域を前記第 2 記憶部に記憶する、

ことを特徴とする付記 9 に記載の情報処理装置。

(付記 11)

前記更新情報は、前記検出手段により検出したデータにより変更された前記第 1 のメモリ領域のページ数を計数した計数情報を含み、

前記登録手段は、前記検出手段により検出したデータを出力した前記外部装置と、前記外部装置を識別する識別情報を記憶する外部装置記憶手段に記憶された識別情報が示す前記外部装置と、が同じ場合であって、前記検出手段により検出したデータにより変更される前記第 1 のメモリ領域を含むページと、前記第 2 記憶部に記憶した前記第 1 のメモリ領域を含むページの次のページと、が同一の場合、前記更新情報に含まれる計数情報が示す計数値に 1 を加算するとともに、前記変更される前記第 1 のメモリ領域を前記第 2 記憶部に記憶する、

ことを特徴とする付記 9 に記載の情報処理装置。

(付記 12)

10

20

30

40

50

前記登録手段は、前記検出手段により検出したデータを出力した前記外部装置と、前記外部装置を識別する識別情報を記憶する外部装置記憶手段に記憶された識別情報が示す前記外部装置と、が同じ場合であって、過去に変更された複数の前記第1のメモリ領域を前記第2記憶部に記憶し、前記検出手段により検出したデータにより変更される前記第1のメモリ領域が、前記第2記憶部に記憶した前記第1のメモリ領域のいずれかとも異なる場合に、前記更新情報を生成して該更新情報を前記第1記憶部に記憶するとともに、前記変更される前記第1のメモリ領域を前記第2記憶部に記憶する、

ことを特徴とする付記9に記載の情報処理装置。

(付記13)

前記登録手段は、前記検出手段により検出したデータを出力した前記外部装置と、前記外部装置を識別する識別情報を記憶する外部装置記憶手段に記憶された識別情報が示す前記外部装置と、が同じ場合であって、前記検出手段により検出したデータにより変更される前記第1のメモリ領域と一定範囲内の前記第1のメモリ領域が、前記第2記憶部に含まれない場合、前記更新情報を生成して該更新情報を前記第1記憶部に記憶するとともに、前記変更される前記第1のメモリ領域を前記第2記憶部に記憶する、

ことを特徴とする付記12に記載の情報処理装置。

(付記14)

前記更新情報は、変更された前記第1のメモリ領域を含むページのサイズを拡張する拡張情報を含み、

前記登録手段は、前記検出手段により検出したデータを出力した前記外部装置と、前記外部装置を識別する識別情報を記憶する外部装置記憶手段に記憶された識別情報が示す前記外部装置と、が同じ場合であって、前記検出手段により検出したデータにより変更される前記第1のメモリ領域と一定範囲内の前記第1のメモリ領域が前記第2記憶部に含まれる場合、該第2記憶部に含まれる前記第1のメモリ領域についての前記更新情報の拡張情報を、前記検出手段により検出したデータにより変更される前記第1のメモリ領域を含むページのサイズに更新する、

ことを特徴とする付記12に記載の情報処理装置。

(付記15)

前記検出手段が検出するデータは、前記仮想機械からの指示により、前記入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送され、前記第1のメモリ領域に書き込まれるデータである、

ことを特徴とする付記1に記載の情報処理装置。

(付記16)

前記更新情報は、前記データを転送した前記入出力装置を識別する識別情報と、前記データの転送先である前記第1のメモリ領域のアドレスと、を含む情報である、

ことを特徴とする付記1に記載の情報処理装置。

(付記17)

前記第1記憶部は、第1の記憶手段と第2の記憶手段と、を備え、

互いに独立して更新情報の記憶または読み出しが可能である、

ことを特徴とする付記1に記載の情報処理装置。

(付記18)

ハードウェア資源を管理することにより、1以上の仮想機械を実現するステップと、

前記仮想機械に割り当てられる第1のメモリ領域のアドレスと、前記第1のメモリ領域の実メモリである第2のメモリ領域のアドレスと、を相互に変換することにより、外部装置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送を行なうステップと、

前記入出力装置から前記仮想機械に割り当てられる前記第1のメモリ領域に直接的にデータ転送されるデータを検出するステップと、

検出したデータが一定の条件を満たす場合に、前記検出したデータにより変更される前記第1のメモリ領域に関する更新情報を生成し、前記更新情報を第1記憶部に記憶するス

10

20

30

40

50

テップと、

前記第 1 記憶部に記憶される更新情報を出力するステップと、
を備えることを特徴とする情報処理装置の仮想化方法。

(付記 19)

ハードウェア資源を管理することにより、1 以上の仮想機械を実現する仮想機械実現手段と、

前記仮想機械に割り当てられる第 1 のメモリ領域のアドレスと、前記第 1 のメモリ領域の実メモリである第 2 のメモリ領域のアドレスと、を相互に変換することにより、外部装置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第 1 のメモリ領域に直接的にデータ転送を行なうデータ転送手段と、

移動対象の前記仮想機械に割り当てられる第 1 のメモリ領域に格納されるデータを取得して移動先に転送するメモリ領域転送手段と、

前記入出力装置から前記仮想機械に割り当てられる前記第 1 のメモリ領域に直接的にデータ転送されるデータを検出する検出手段と、

検出したデータが一定の条件を満たす場合に、該検出したデータにより変更される前記第 1 のメモリ領域に関する更新情報を生成し、前記更新情報を第 1 記憶部に記憶する登録手段と、

前記第 1 記憶部から更新情報を取得する更新情報取得手段と、

前記更新情報取得手段により取得した更新情報に基づいて、移動対象の前記仮想機械に割り当てられる第 1 のメモリ領域において前記検出手段により検出したデータにより変更された更新データを、前記移動先に転送する更新データ転送手段と、

を備えることを特徴とする情報処理装置。

(付記 20)

前記更新データ転送手段は、

前記第 1 記憶部に記憶される更新情報の数が一定値以下になるまで、前記更新データを前記移動先に転送する第 1 の更新データ転送手段と、

前記移動対象の前記仮想機械の動作を停止した後に、前記更新データを前記移動先に転送する第 2 の更新データ転送手段と、

を備える、

ことを特徴とする付記 19 に記載の情報処理装置。

(付記 21)

第 1 の情報処理装置により実現する仮想機械を、前記第 1 の情報処理装置と通信可能に接続する第 2 の情報処理装置に移動させるマイグレーション方法において、

ハードウェア資源を管理することにより、1 または 2 以上の前記仮想機械を実現するステップと、

前記仮想機械に割り当てられる第 1 のメモリ領域のアドレスと、前記第 1 のメモリ領域の実メモリである第 2 のメモリ領域のアドレスと、を相互に変換することにより、外部装置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第 1 のメモリ領域に直接的にデータ転送を行なうステップと、

移動対象の前記仮想機械に割り当てられる第 1 のメモリ領域に格納されるデータを取得して前記第 2 の情報処理装置に転送するステップと、

前記入出力装置から前記仮想機械に割り当てられる前記第 1 のメモリ領域に直接的にデータ転送されるデータを検出するステップと、

検出したデータが一定の条件を満たす場合に、前記検出したデータにより変更される前記第 1 のメモリ領域に関する更新情報を生成し、前記更新情報を第 1 記憶部に記憶するステップと、

前記第 1 記憶部から更新情報を取得するステップと、

前記更新情報に基づいて、移動対象の前記仮想機械に割り当てられる第 1 のメモリ領域において前記検出したデータにより変更された更新データを、前記第 2 の情報処理装置に転送するステップと、

10

20

30

40

50

を備えることを特徴とするマイグレーション方法。

(付記 2 2)

演算処理装置を有する第 1 の情報処理装置により実現する仮想機械を、前記第 1 の情報処理装置と通信可能に接続する第 2 の情報処理装置に移動させるマイグレーション用プログラムにおいて、

ハードウェア資源を管理することにより、1 以上の仮想機械を実現する仮想機械実現処理と、

前記仮想機械に割り当てられる第 1 のメモリ領域のアドレスと、前記第 1 のメモリ領域の実メモリである第 2 のメモリ領域のアドレスと、を相互に変換することにより、外部装置とのデータの入出力を制御する入出力装置から前記仮想機械に割り当てられる前記第 1 のメモリ領域に直接的にデータ転送を行なうデータ転送処理と、

移動対象の前記仮想機械に割り当てられる第 1 のメモリ領域に格納されるデータを取得して前記第 1 の情報処理装置に転送するメモリ領域転送処理と、

前記入出力装置から前記仮想機械に割り当てられる前記第 1 のメモリ領域に直接的にデータ転送されるデータを検出する転送データ検出処理と、

検出したデータが一定の条件を満たす場合に、該検出したデータにより変更される前記第 1 のメモリ領域に関する更新情報を生成し、前記更新情報を第 1 記憶部に記憶する更新情報登録処理と、

前記第 1 記憶部から更新情報を取得する更新情報取得処理と、

前記更新情報取得処理により取得した更新情報に基づいて、移動対象の前記仮想機械に割り当てられる第 1 のメモリ領域において前記転送データ検出処理により検出したデータにより変更された更新データを、前記第 1 の情報処理装置に転送する更新データ転送処理と、

を前記演算処理装置に実行させることを特徴とするプログラム。

【符号の説明】

【0323】

650	P C I e - スイッチ
651	ダーティ・ページ管理ユニット
810	制御 I / F 部
820	パケット検出部
830	ライトポインタ
840	リードポインタ
850	D M A ライトパケット
901	F I F O 制御部
902	判定部

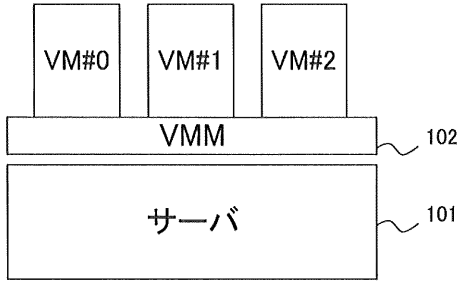
10

20

30

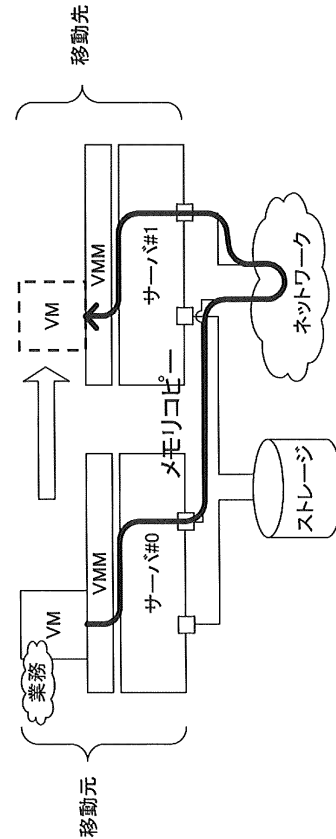
【 図 1 】

サーバ仮想化技術の概要を示す図



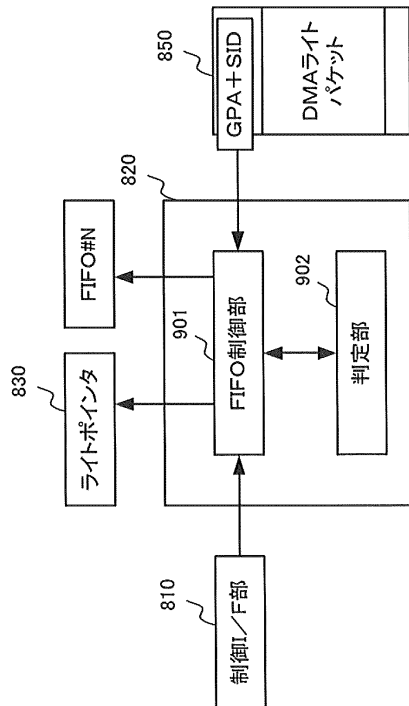
【 図 3 】

ライブ・マイグレーションの概要を示す図



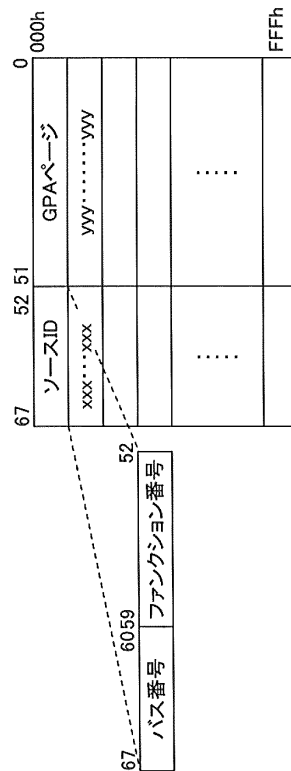
【 図 9 】

パケット検出部の構成例を示す図



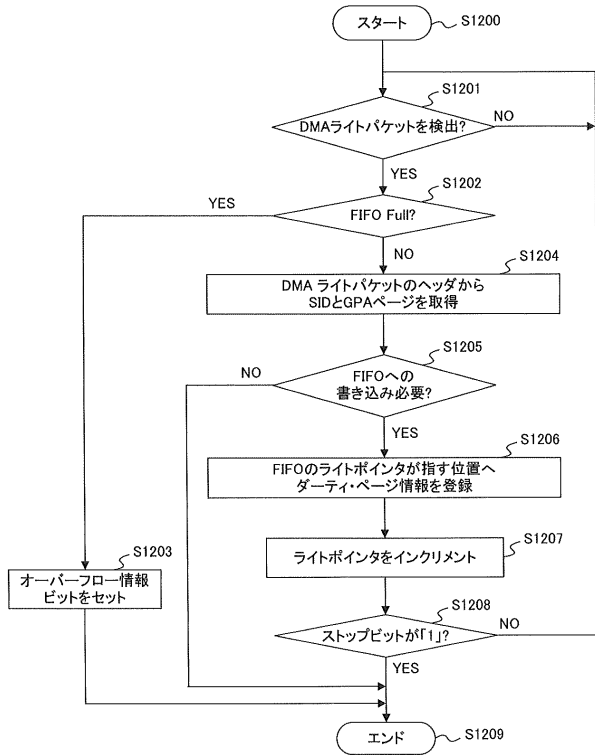
【 図 10 】

FIFO#0の構成例を示す図



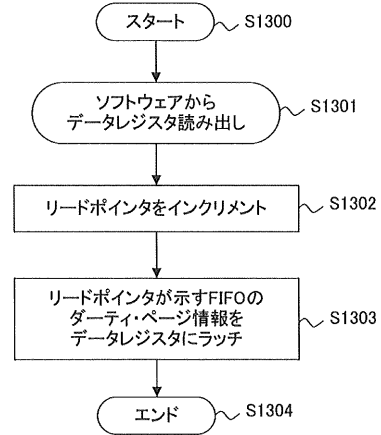
【 図 1 2 】

ダーティ・ページ管理ユニットによるダーティ・ページ情報の登録処理を示すフローチャート



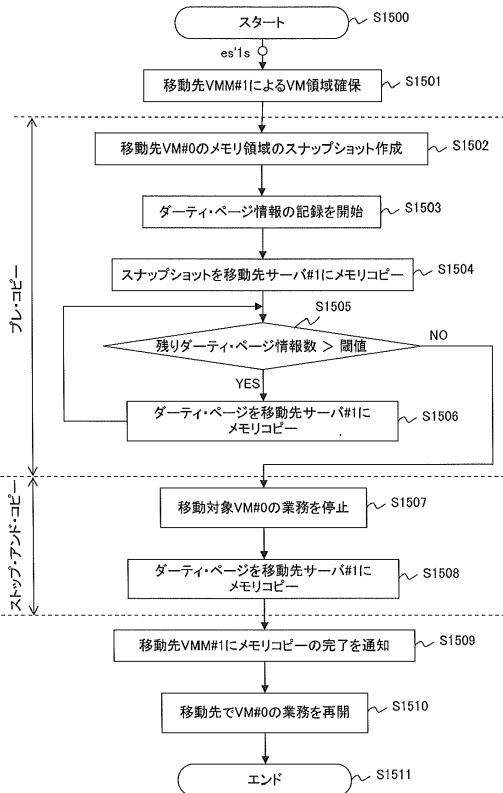
【 図 1 3 】

ダーティ・ページ管理ユニットによるダーティ・ページ情報の出力処理を示すフローチャート



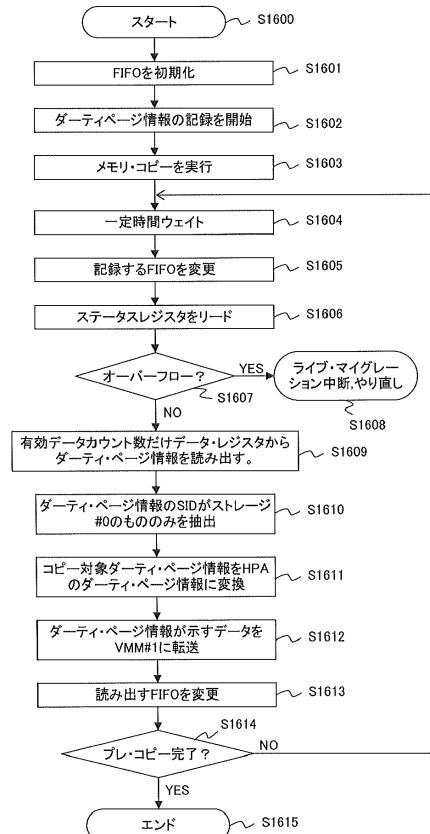
【 図 1 5 】

本実施例に係るライブ・マイグレーションの概要を示すフローチャート



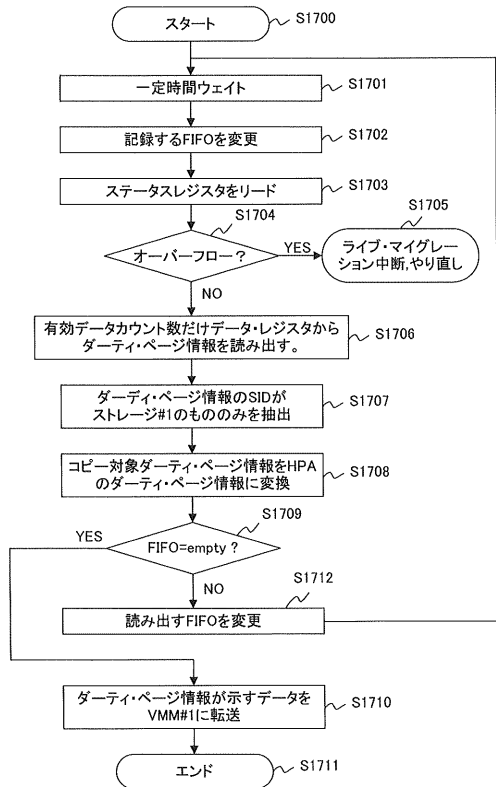
【 図 1 6 】

プレ・コピーの具体的な処理を示すフローチャート



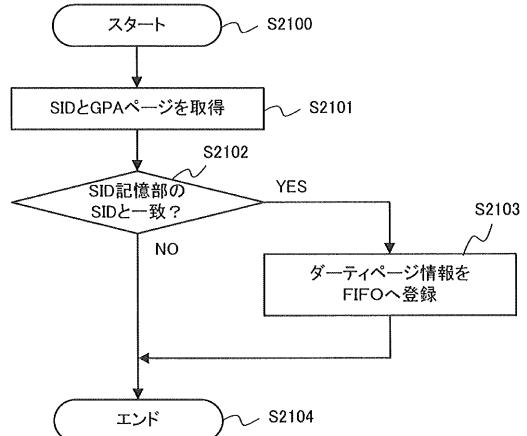
【 図 1 7 】

ストップ・アンド・コピーの具体的な処理を示すフローチャート



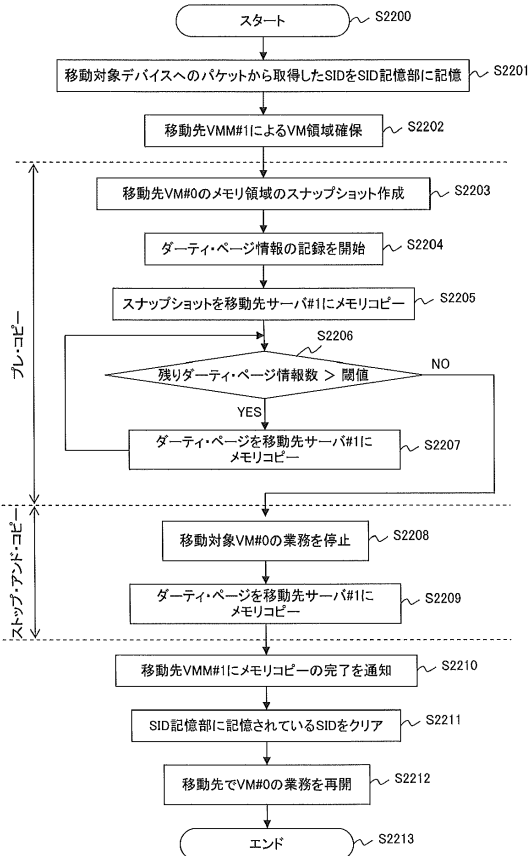
【 図 2 1 】

図20に示したパケット検出部の処理を示すフローチャート



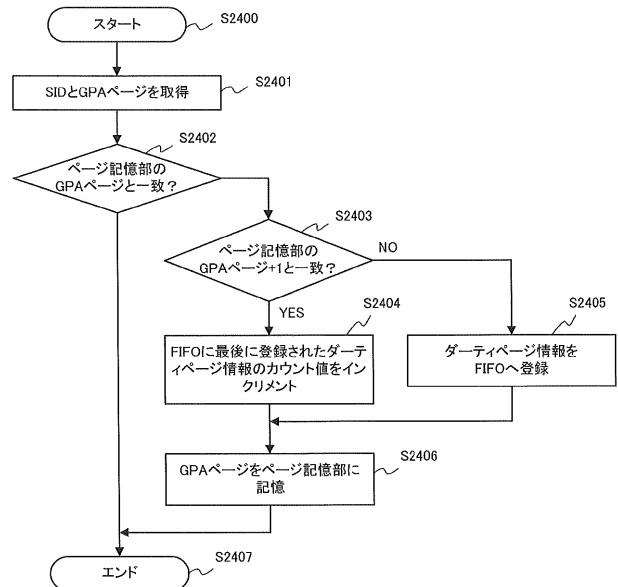
【 図 2 2 】

ライブ・マイグレーション時におけるSID記憶部へのSID設定処理を示すフローチャート



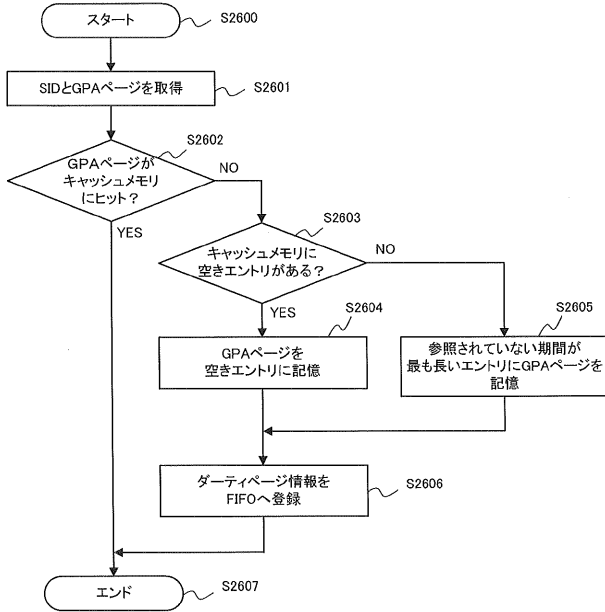
【 図 2 4 】

図23に示したパケット検出部の処理を示すフローチャート



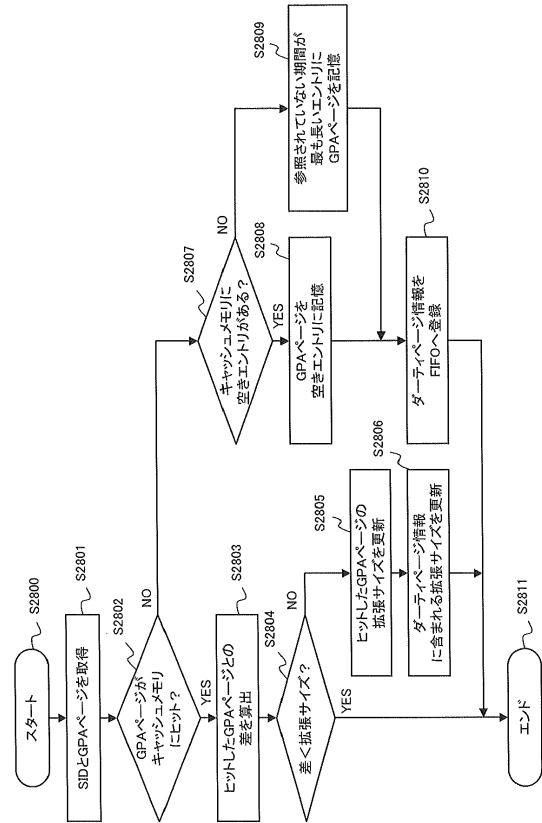
【 図 2 6 】

図25に示したパケット検出部の処理を示すフローチャート



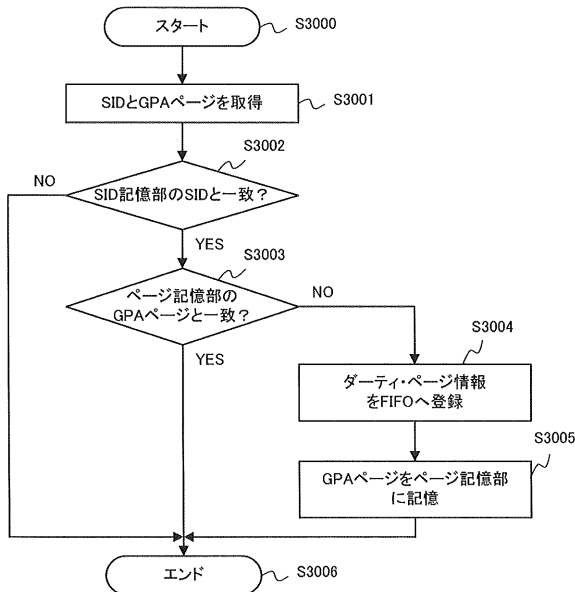
【 図 2 8 】

図27に示したパケット検出部の処理を示すフローチャート



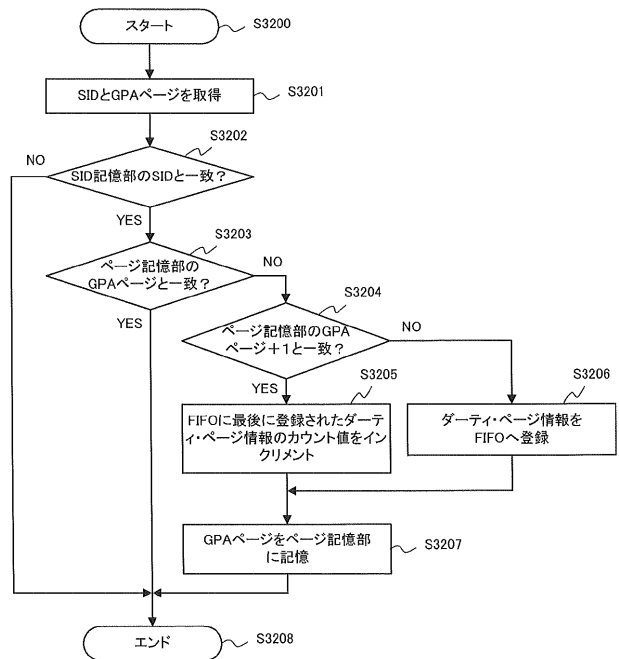
【 図 3 0 】

図29に示したパケット検出部の処理を示すフローチャート



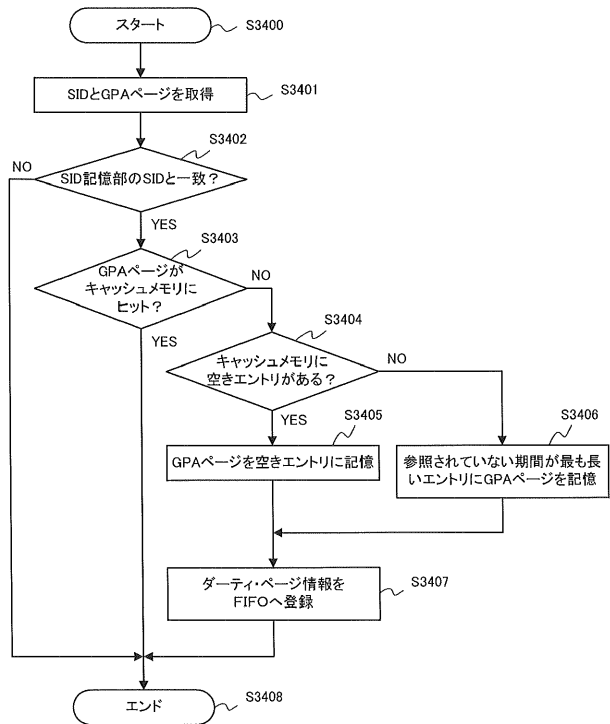
【 図 3 2 】

図31に示したパケット検出部の処理を示すフローチャート



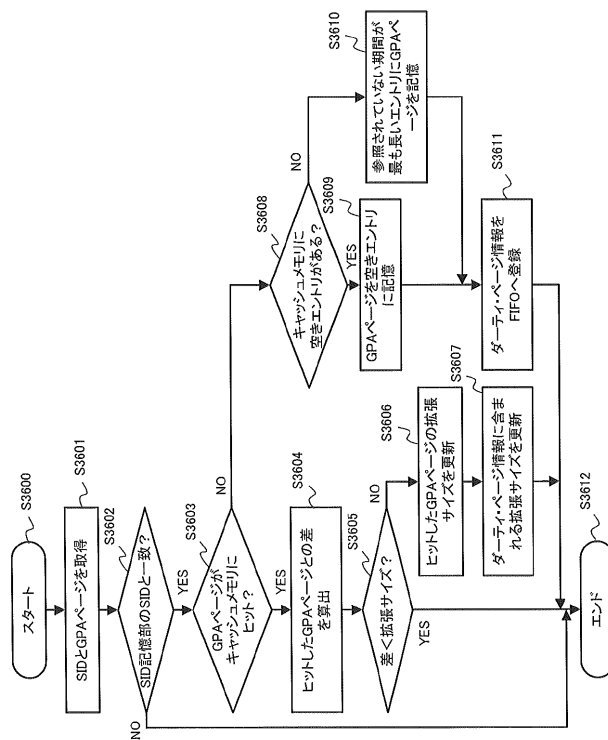
【 図 3 4 】

図33に示したパケット検出部の処理を示すフローチャート



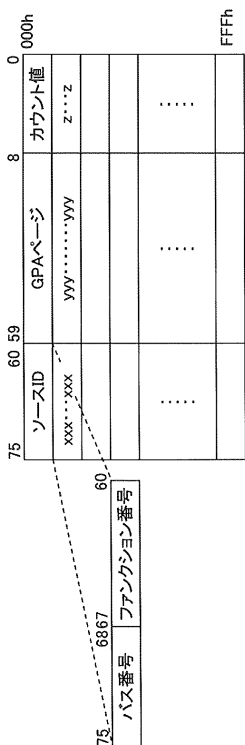
【 図 3 6 】

図35に示したパケット検出部の処理を示すフローチャート



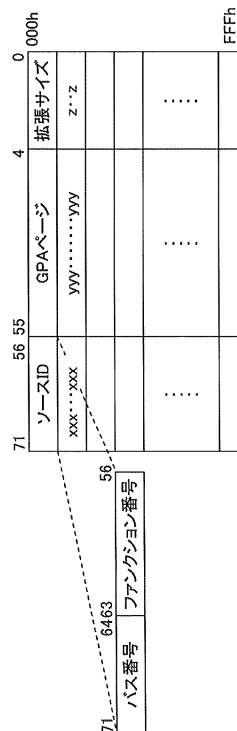
【 図 3 7 】

条件3を用いたパケット検出部を実現する場合に使用するFIFOの構成例を示す図



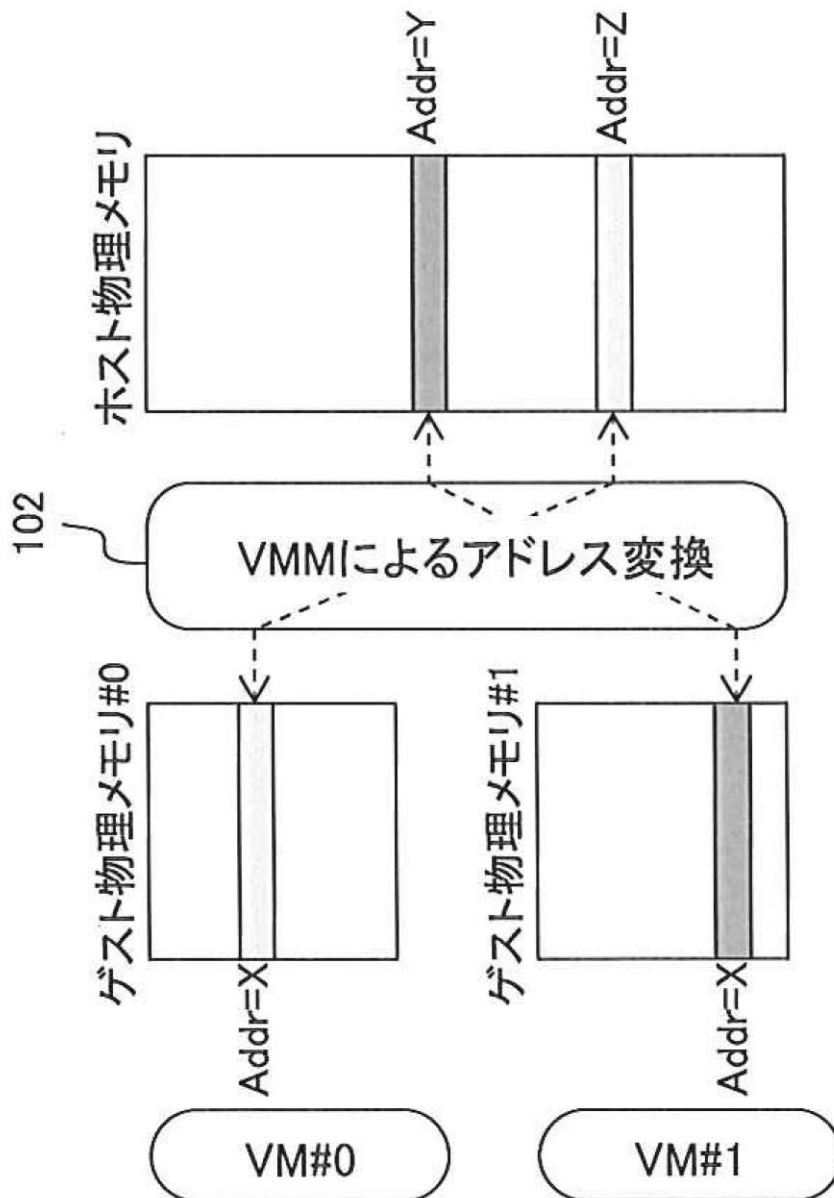
【 図 3 8 】

条件5を用いたパケット検出部を実現する場合に使用するFIFOの構成例を示す図



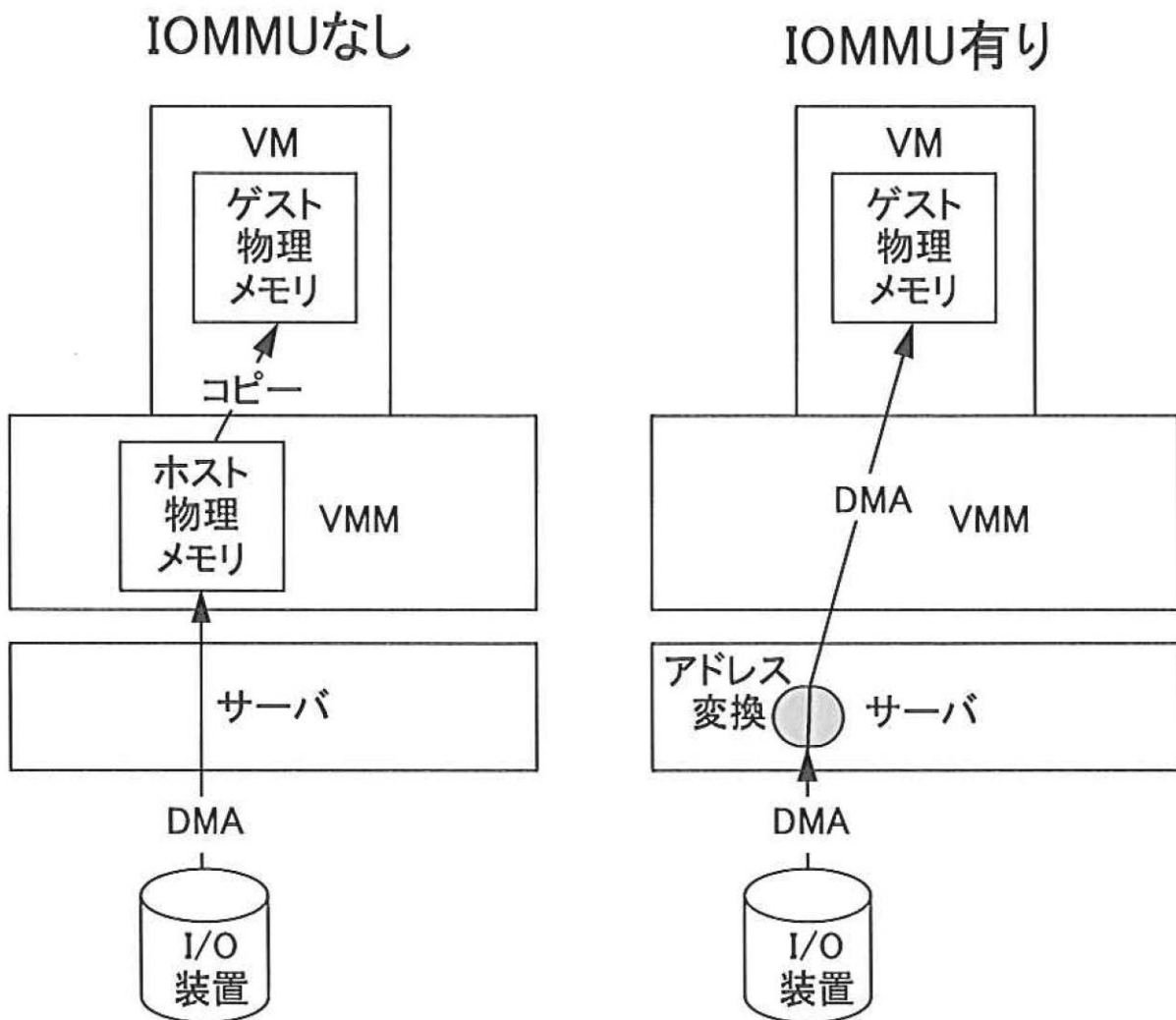
【 図 2 】

VMMによるメモリ管理の概要を示す図



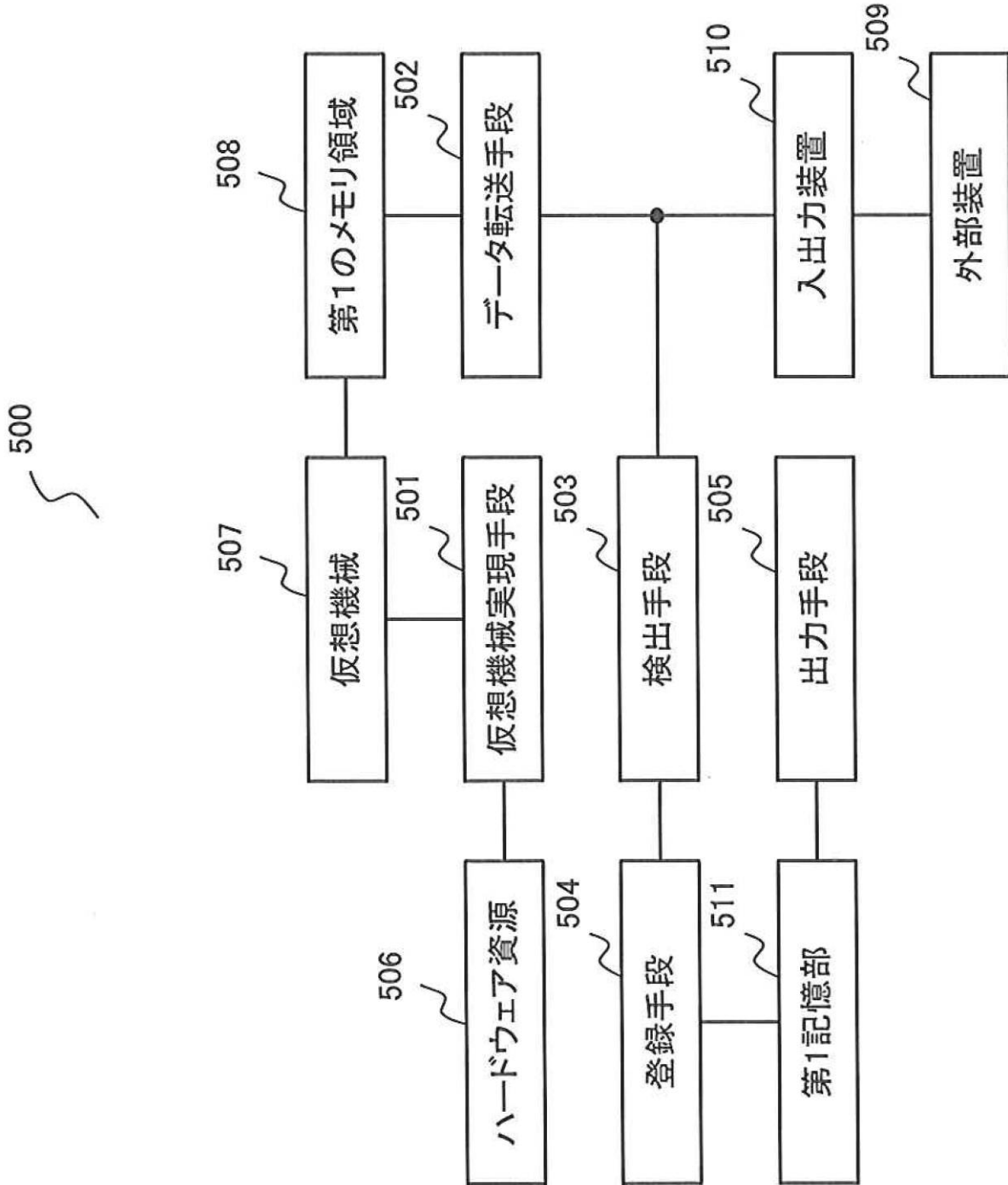
【 図 4 】

IOMMUについて説明する図



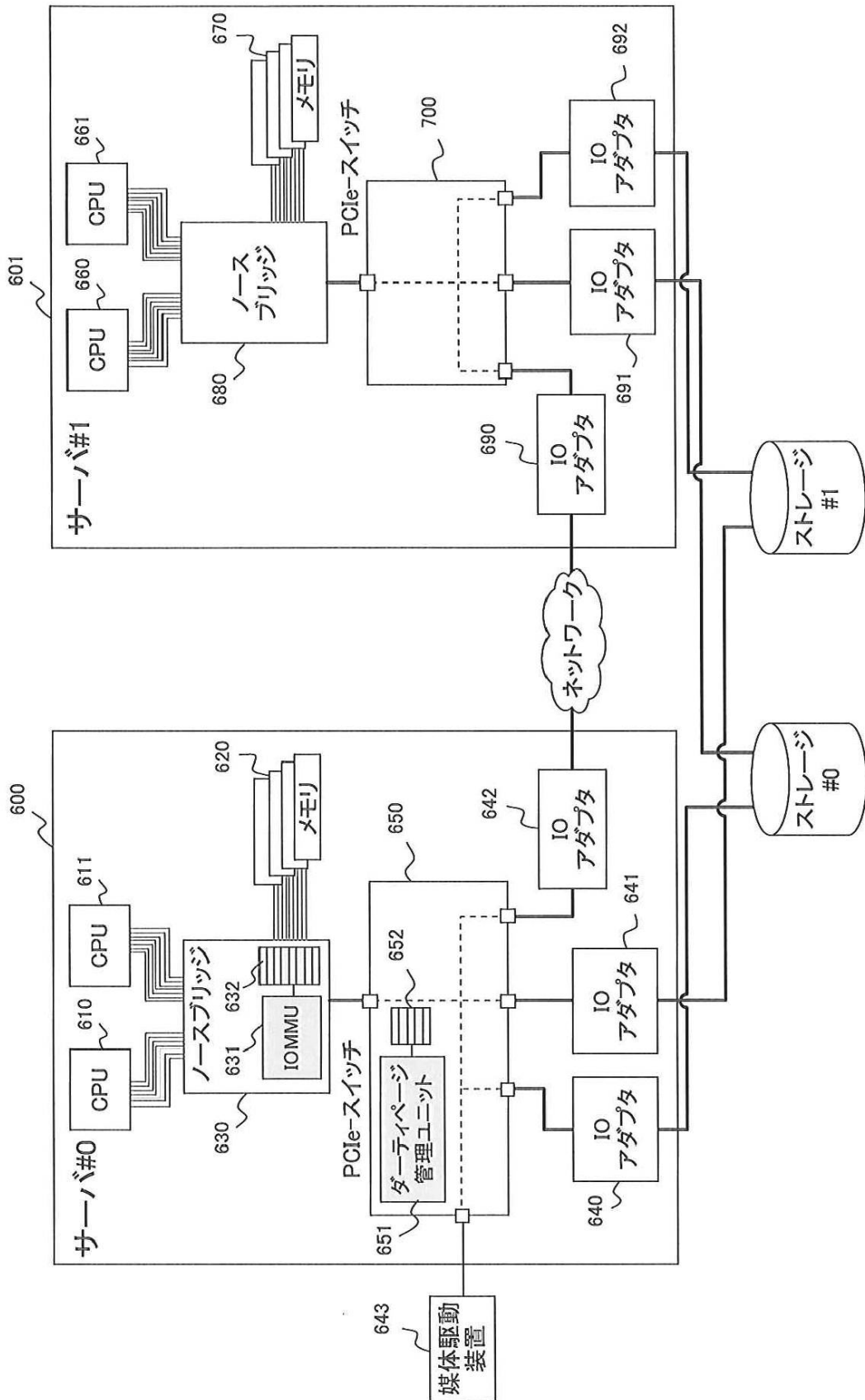
【図5】

情報処理装置を説明する図



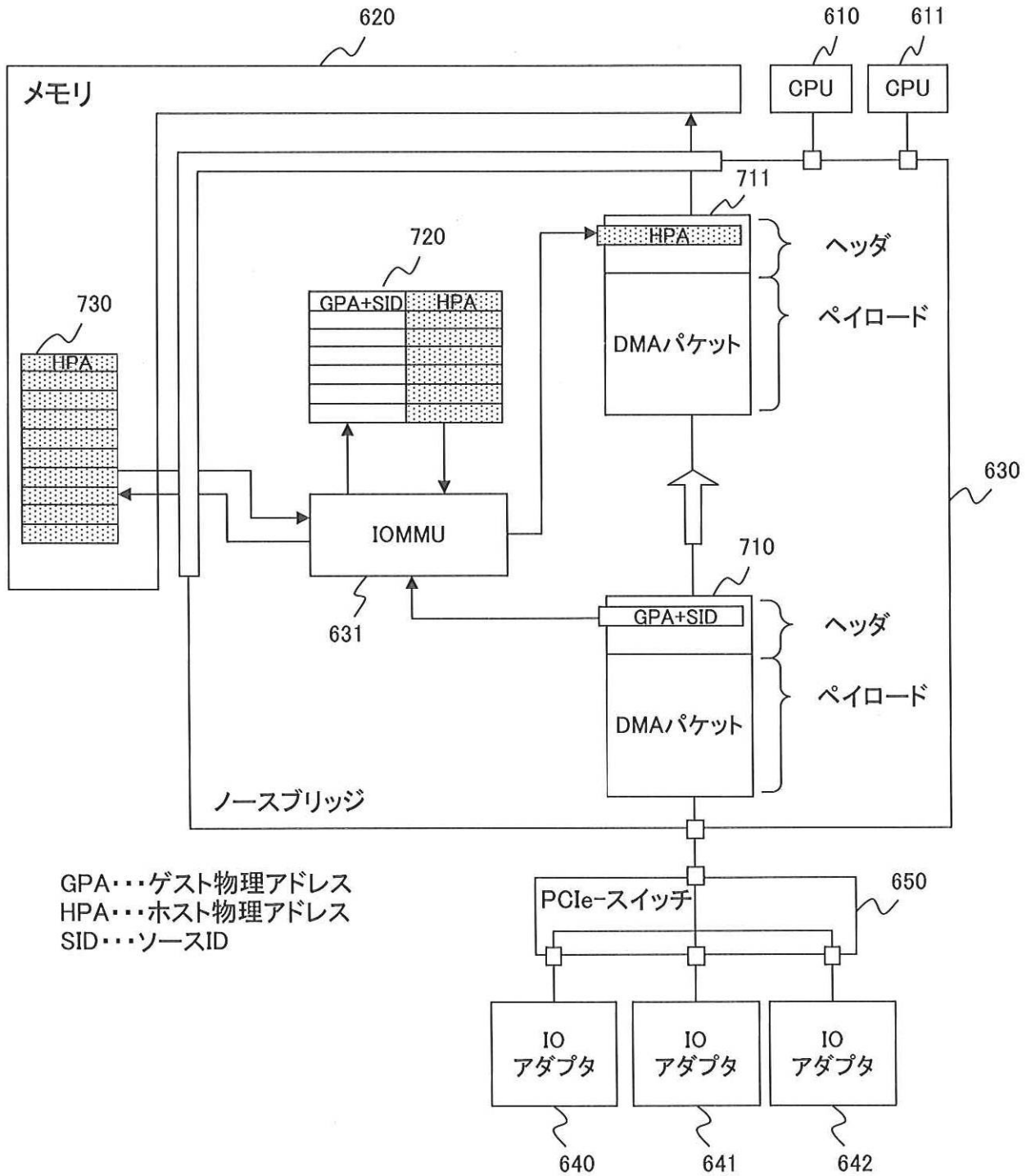
【図6】

情報処理装置の構成例を示す図



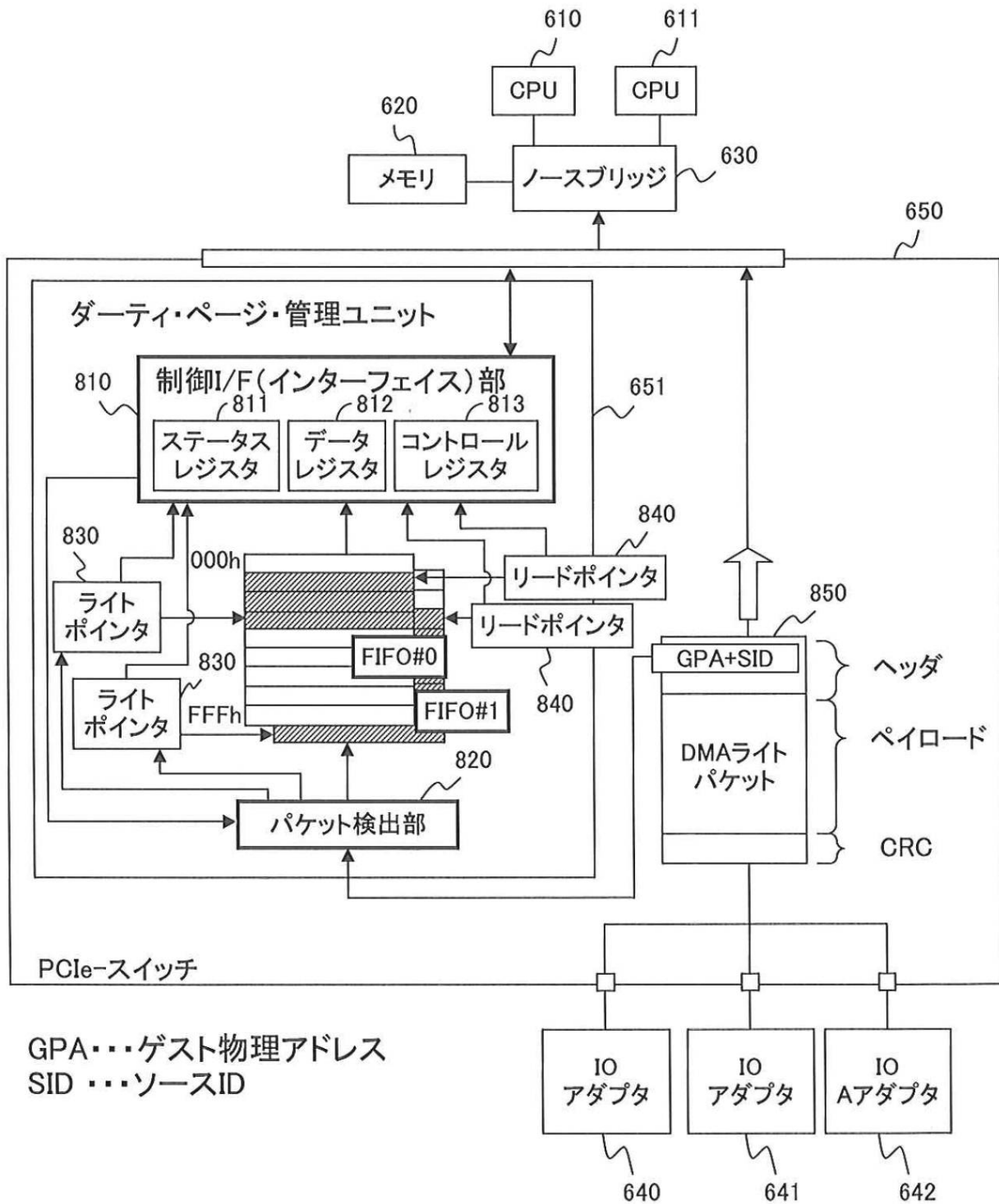
【図7】

ノースブリッジにおけるDMA処理の動作を説明する図



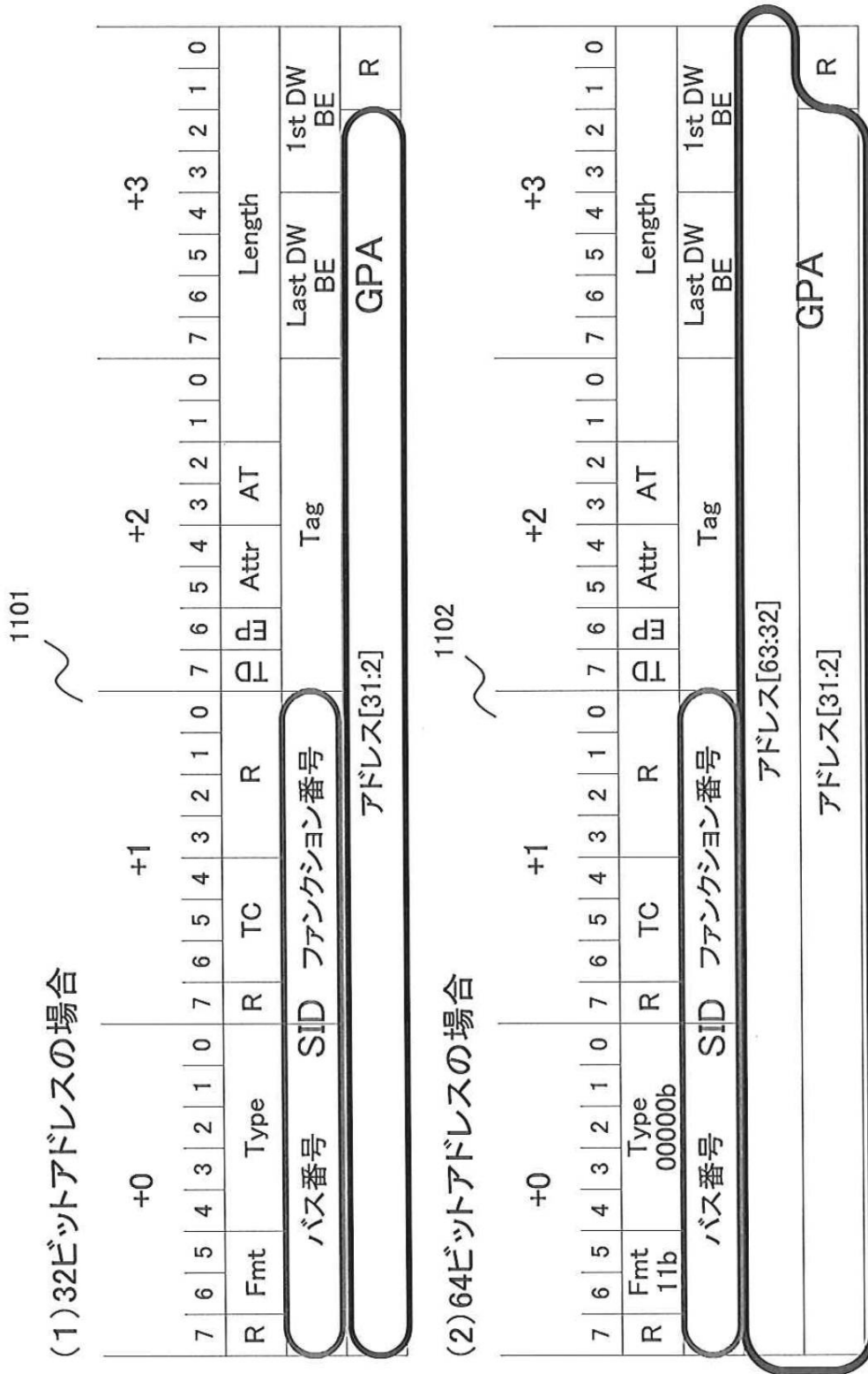
【 図 8 】

PCIeスイッチにおけるDMA処理の動作を説明する図



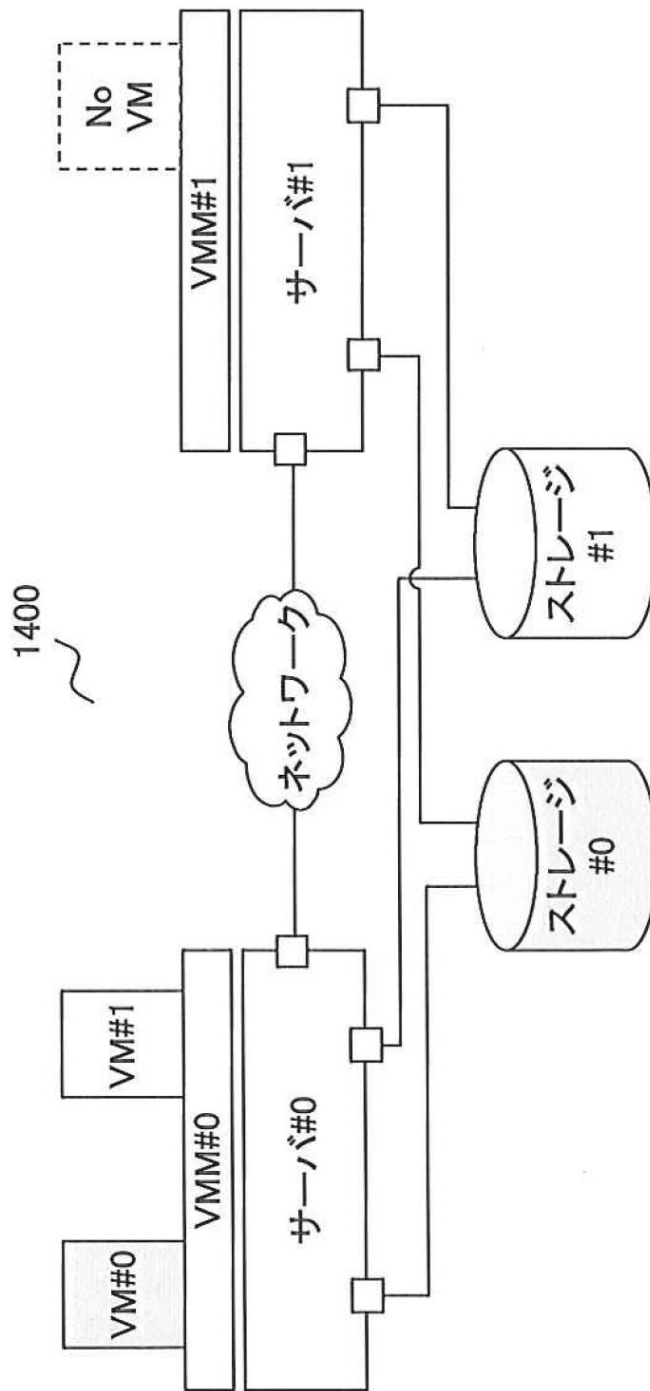
【図 1 1】

DMAパケットのヘッダの構成例を示す図



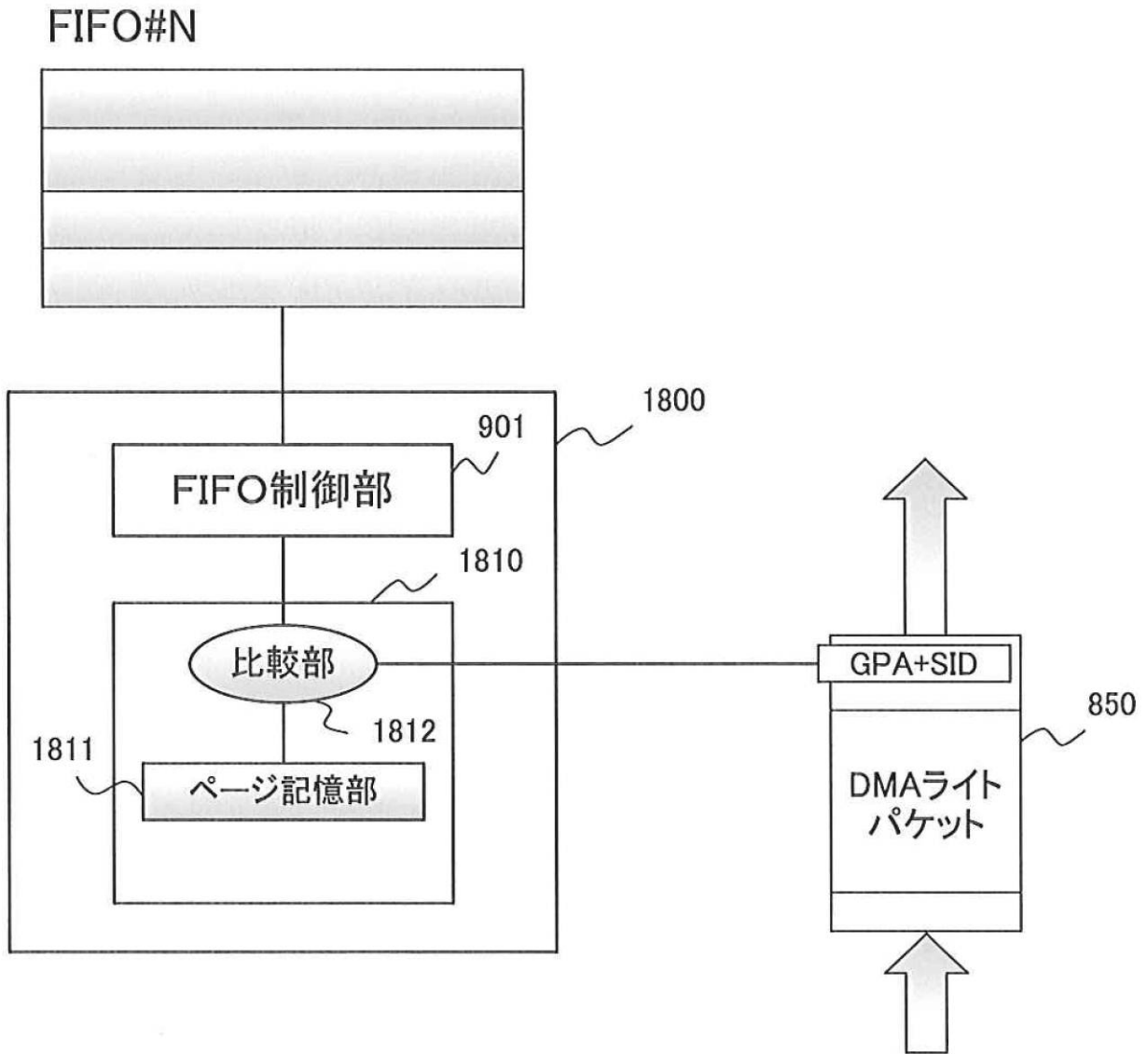
【図14】

本実施例に係るライブ・マイグレーションの概要を説明する図



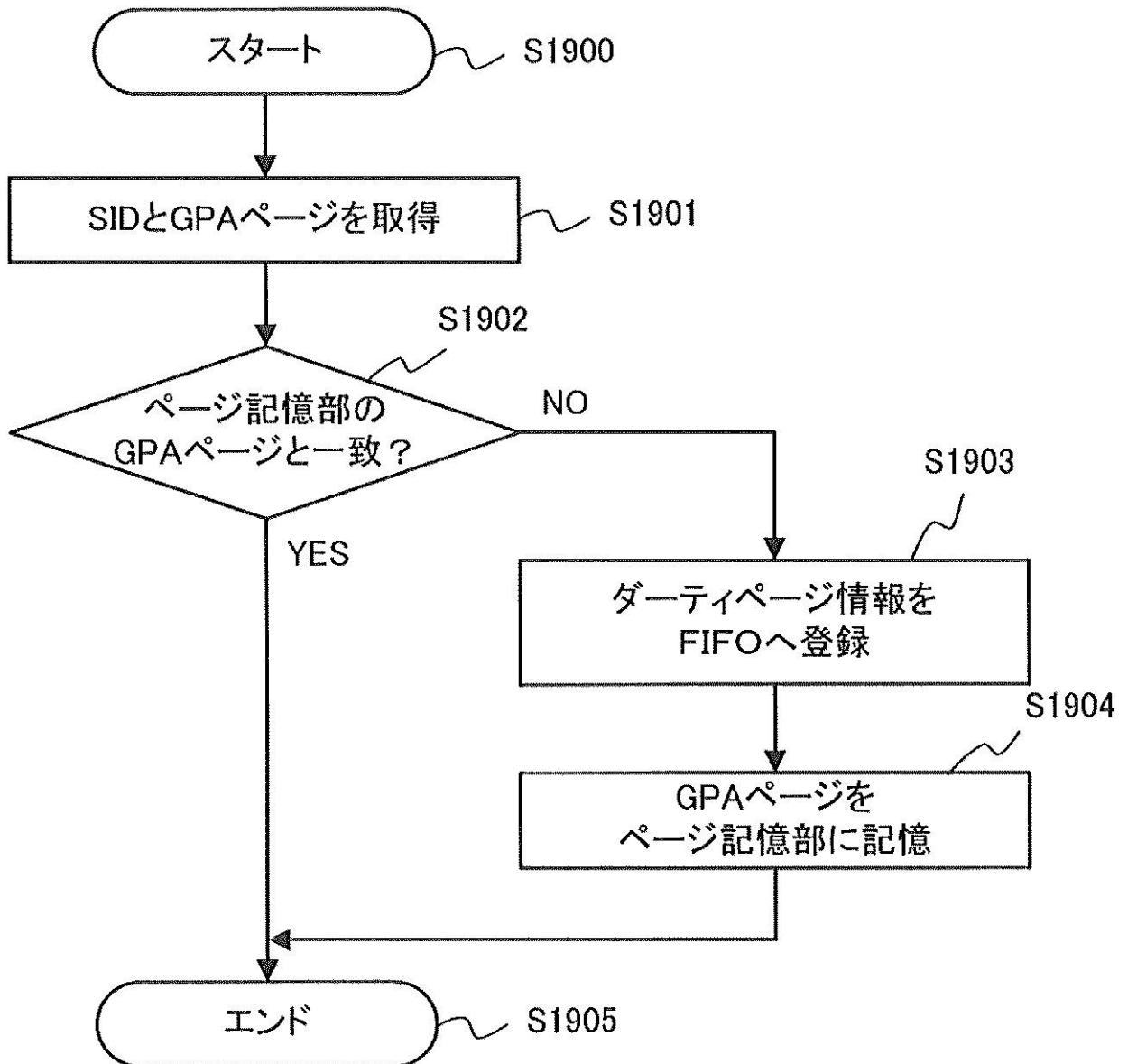
【図18】

条件1を用いたパケット検出部の具体例を示す図



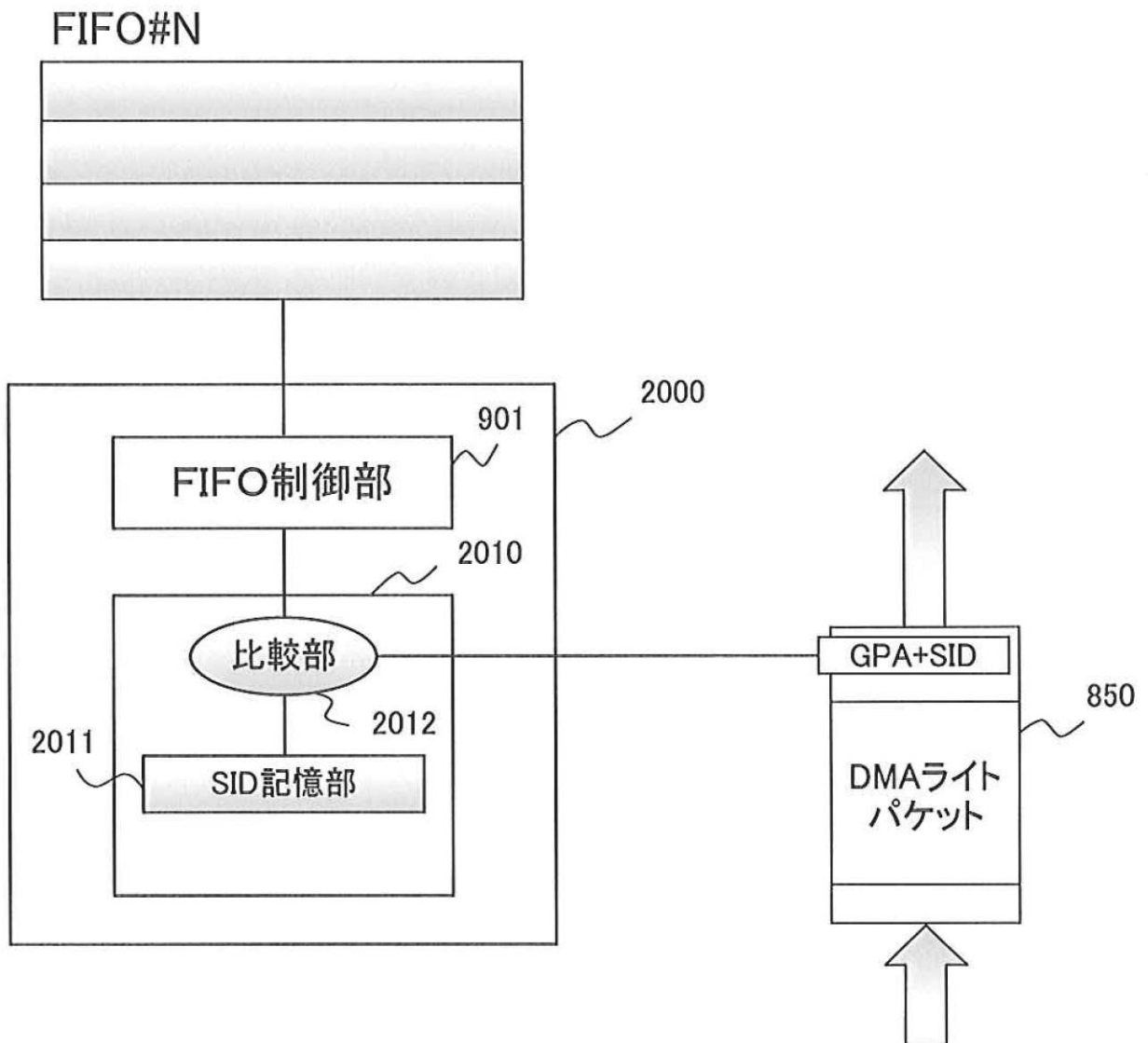
【図19】

図18に示したパケット検出部の処理を示すフローチャート



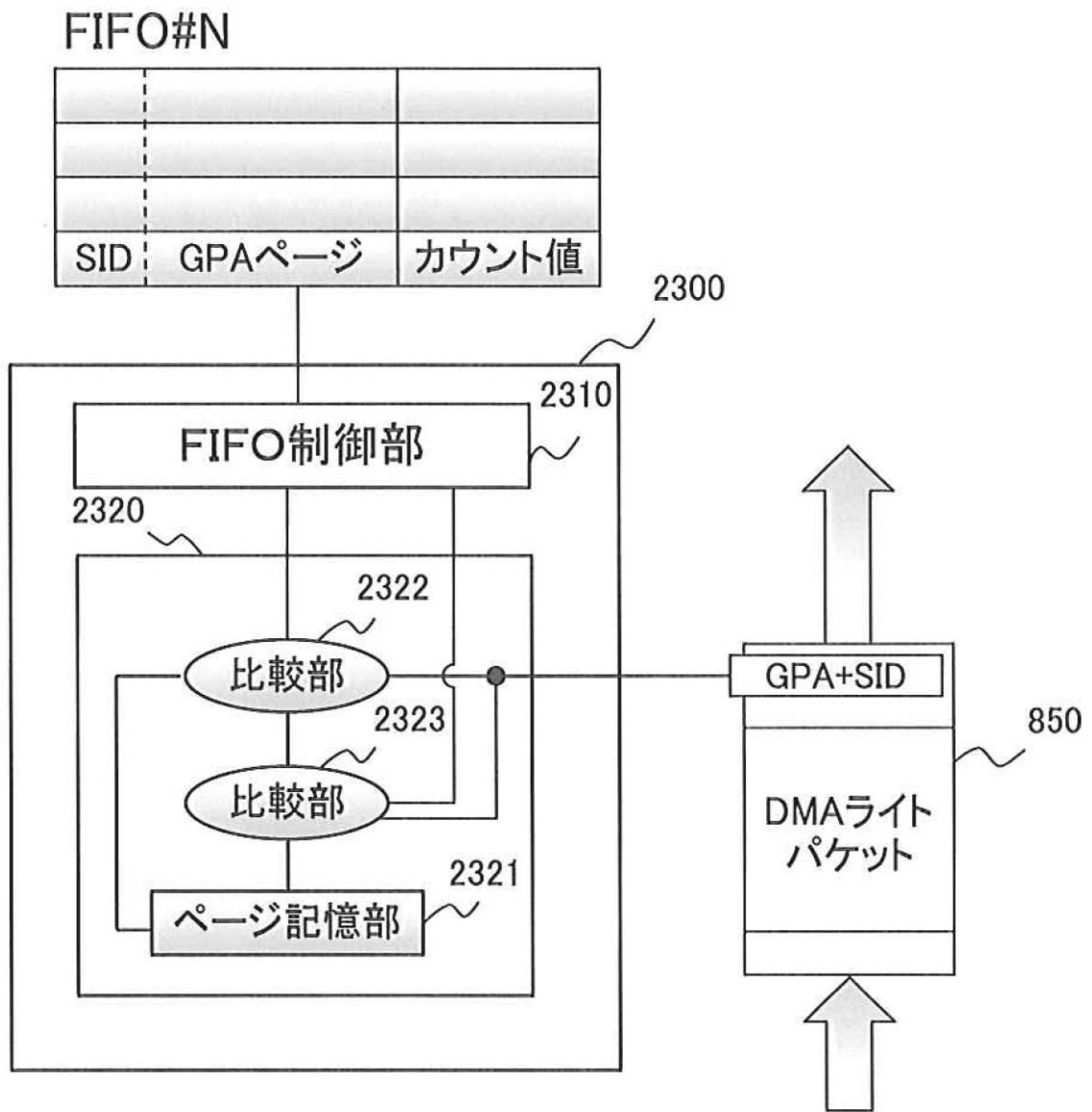
【図 20】

条件2を用いたパケット検出部の具体例を示す図



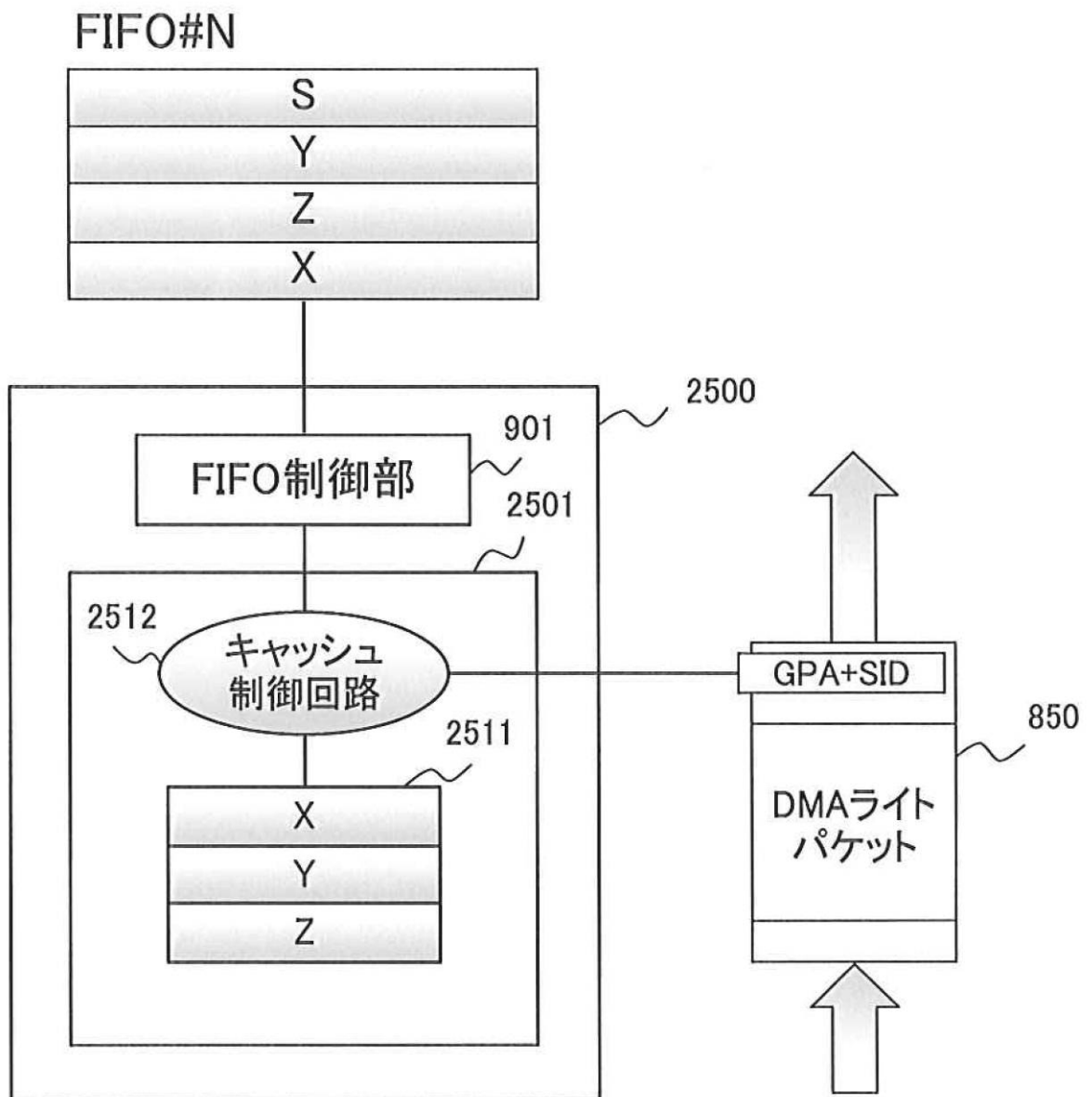
【図 23】

条件3を用いたパケット検出部の具体例を示す図



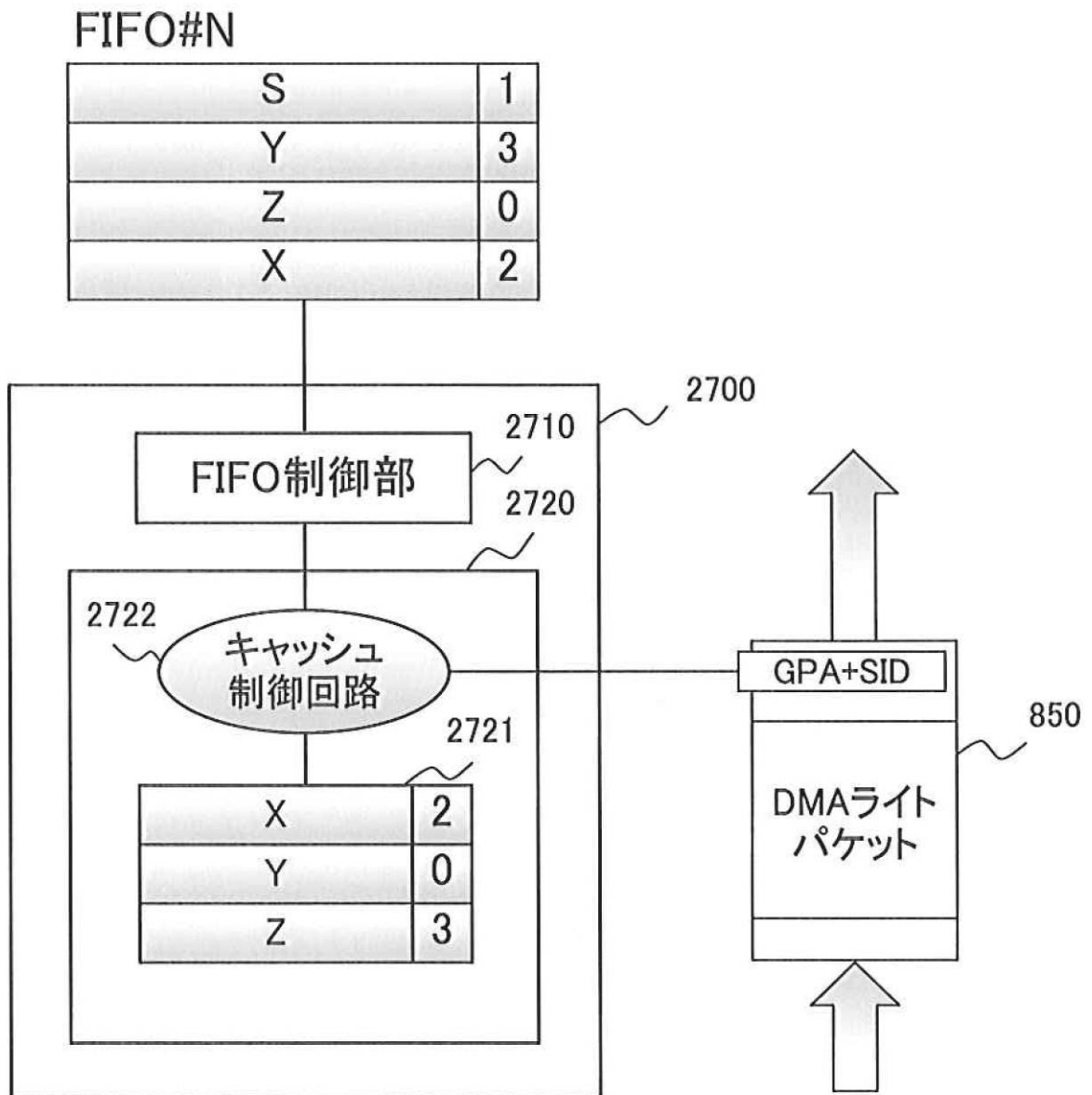
【図 25】

条件4を用いたパケット検出部の具体例を示す図



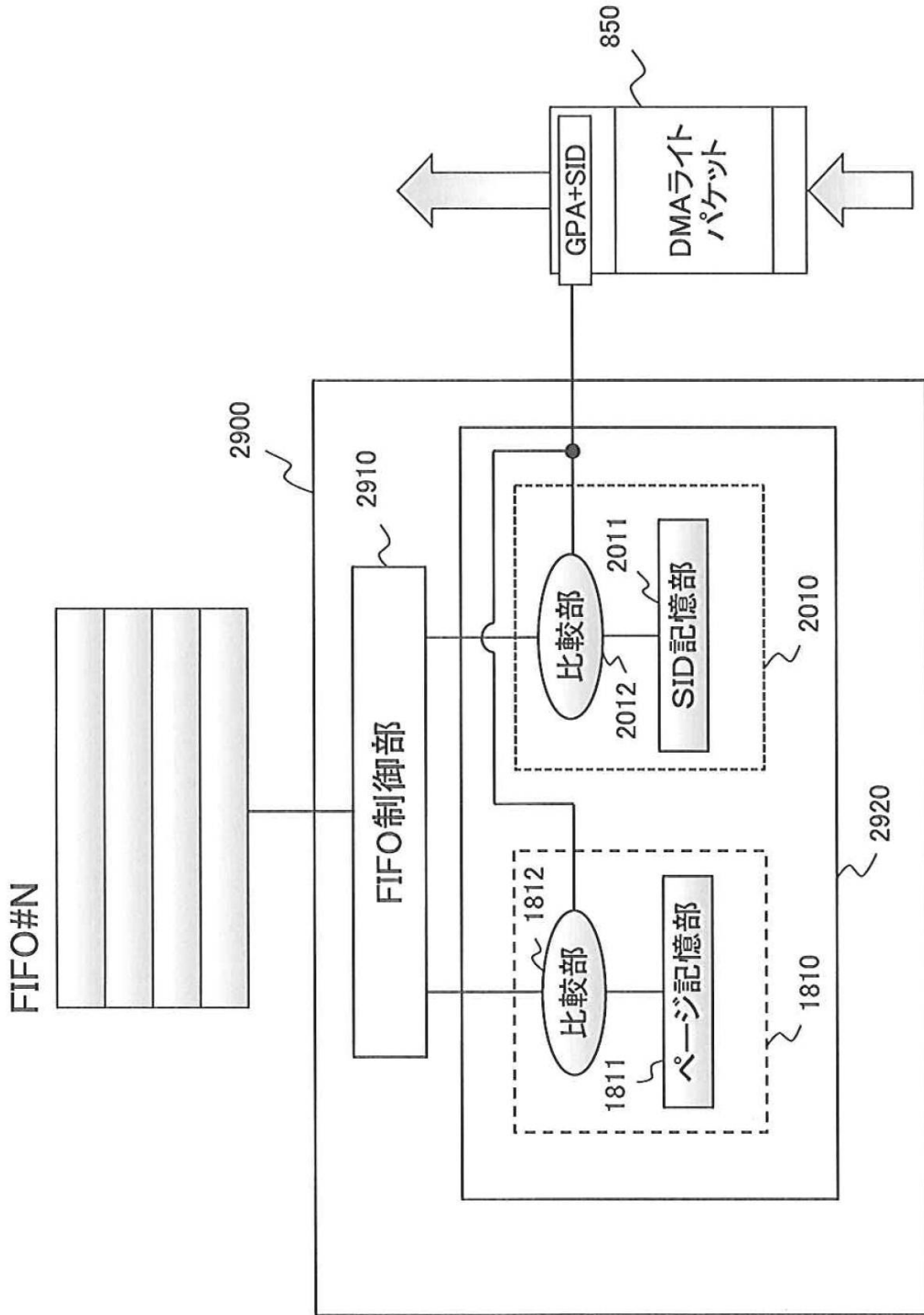
【図 27】

条件5を用いたパケット検出部の具体例を示す図



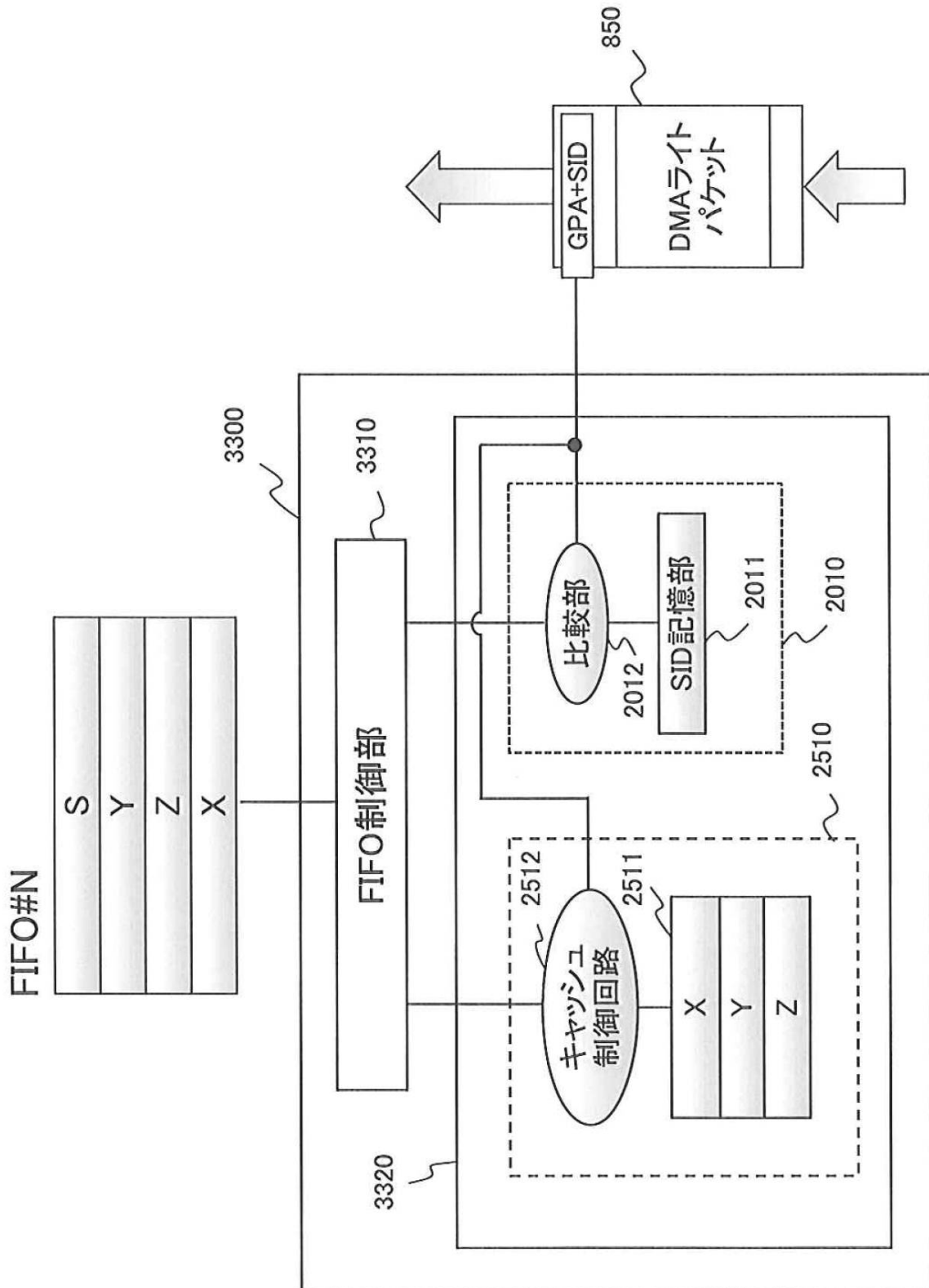
【図29】

条件1および2を用いたパケット検出部の
具体例を示す図



【図33】

条件2および4を用いたパケット検出部の
具体例を示す図



【図 35】

条件2および5を用いたパケット検出部の
具体例を示す図

