

(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2019-200698

(P2019-200698A)

(43) 公開日 令和1年11月21日(2019.11.21)

(51) Int.Cl.		F I		テーマコード (参考)
G06F 12/06	(2006.01)	G06F 12/06	550B	5B160
H04L 12/771	(2013.01)	H04L 12/771		5K030

審査請求 未請求 請求項の数 8 O L (全 25 頁)

(21) 出願番号 特願2018-96227 (P2018-96227)
 (22) 出願日 平成30年5月18日 (2018.5.18)

(71) 出願人 000004226
 日本電信電話株式会社
 東京都千代田区大手町一丁目5番1号
 (71) 出願人 504132272
 国立大学法人京都大学
 京都府京都市左京区吉田本町36番地1
 (74) 代理人 110001863
 特許業務法人アテンダ国際特許事務所
 (72) 発明者 郡川 智洋
 東京都千代田区大手町一丁目5番1号 日
 本電信電話株式会社内
 (72) 発明者 川端 明生
 東京都千代田区大手町一丁目5番1号 日
 本電信電話株式会社内

最終頁に続く

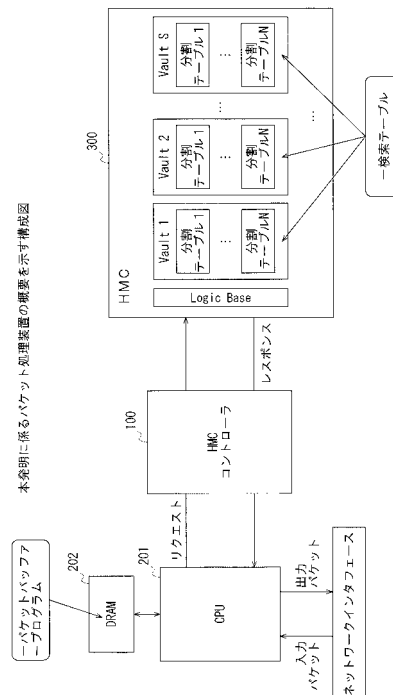
(54) 【発明の名称】 パケット処理装置及びそのメモリアクセス制御方法

(57) 【要約】

【課題】 仮想化環境での使用を前提とした、汎用デバイスから構成される汎用装置で、検索テーブルへの高いメモリアクセス性能を有するパケット処理装置を提供する。

【解決手段】 パケット処理においてCPU201からアクセスされるテーブルを記憶したHMC300と、CPU201からのHMC300300の前記テーブルへのメモリアクセスを制御するHMCコントローラ100とを備えたパケット処理装置であって、HMCは並列アクセス可能なS個のブロック(Vault)に区画されており、前記各ブロックは並列アクセス可能なN個のバンクに区画されており、前記テーブルは分割されてHMC300の前記バンクに分散して記憶されており、HMCコントローラ300は、前記アクセスリクエストに対してアクセス対象データが格納されているブロック及びバンクを特定して前記HMC300にアクセスする。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

パケット処理において演算装置からアクセスされるテーブルを記憶した記憶装置と、前記演算装置からの前記記憶装置の前記テーブルへのアクセスリクエストに基づき前記記憶装置へのメモリアクセスを制御する制御装置とを備えたパケット処理装置であって、

前記記憶装置の記憶領域は互いに並列アクセス可能な S 個（ S は2以上の自然数）のブロックに区画されており、

前記各ブロックの記憶領域は互いに並列アクセス可能な N 個（ N は2以上の自然数）のバンクに区画されており、

前記テーブルは分割されて前記記憶装置の前記バンクに分散して記憶されており、

前記制御装置は、前記アクセスリクエストに対してアクセス対象データが格納されているブロック及びバンクを特定して前記記憶装置にアクセスする

ことを特徴とするパケット処理装置。

【請求項 2】

前記テーブルを N 個の分割テーブルに等分割し、前記 S 個のブロックのそれぞれにおいて、前記 N 個の分割テーブルを前記 N 個のバンクに対応させて記憶した

ことを特徴とする請求項 1 記載のパケット処理装置。

【請求項 3】

前記制御装置は、さらに、ブロック識別子及びバンク識別子の組により特定される記憶領域へのアクセス状態を管理するアクセス状態管理部と、同一データが分散記憶されたブロック識別子及びバンク識別子の組により特定される複数の記憶領域からアクセス状態がアクセス可能である記憶領域を選択し、選択した記憶領域へのアクセスを行うアクセス制御部とを備えた

ことを特徴とする請求項 1 又は 2 記載のパケット処理装置。

【請求項 4】

前記制御装置は、さらに、前記演算装置からのアクセスリクエストの種別を識別するリクエスト識別部と、前記リクエスト識別部により識別された種別が前記テーブルの更新リクエストの場合に当該更新リクエストに基づき前記 S 個のブロックのそれぞれにおいてテーブルの更新処理を行うデータ更新制御部とを備えた

ことを特徴とする請求項 1 乃至 3 何れか 1 項記載のパケット処理装置。

【請求項 5】

前記制御装置は、さらに、前記演算装置からのアクセスリクエストによる前記記憶装置の負荷を監視する負荷監視部と、前記負荷監視部により計測された負荷に基づき1つ以上のブロックにおいて第1のバンクに記憶されている分割テーブルを第2のバンクにコピーする分割変動制御部とを備えた

ことを特徴とする請求項 1 乃至 4 何れか 1 項記載のパケット処理装置。

【請求項 6】

前記アクセスリクエストはパケットに付随するパケット付随情報を含み、

前記テーブルは前記パケット付随情報に基づき分割されている

ことを特徴とする請求項 1 乃至 5 何れか 1 項記載のパケット処理装置。

【請求項 7】

前記記憶装置は、複数のデータ記憶素子層とメモリコントロール機能層とを互いに接続するように積層するとともに、各データ記憶素子層を平面上において S 個の区画に分割するとともに各データ記憶素子層の同一区画間を互いに接続することによりブロックを形成した

ことを特徴とする請求項 1 乃至 6 何れか 1 項記載のパケット処理装置。

【請求項 8】

パケット処理において演算装置からアクセスされるテーブルを記憶した記憶装置と、前記演算装置からの前記記憶装置の前記テーブルへのアクセスリクエストに基づき前記記憶装置へのメモリアクセスを制御する制御装置とを備えたパケット処理装置におけるメモリ

10

20

30

40

50

アクセス制御方法であって、

前記記憶装置の記憶領域は互いに並列アクセス可能な S 個(S は2以上の自然数)のブロックに区画されており、

前記各ブロックの記憶領域は互いに並列アクセス可能な N 個(N は2以上の自然数)のバンクに区画されており、

前記テーブルは分割されて前記記憶装置の前記バンクに分散して記憶されており、

前記制御装置は、前記アクセスリクエストに対してアクセス対象データが格納されているブロック及びバンクを特定して前記記憶装置にアクセスする

ことを特徴とするパケット処理装置のメモリアクセス制御方法。

【発明の詳細な説明】

10

【技術分野】

【0001】

本発明は、通信ネットワークにおける大規模トラフィックフローを対象とするパケット処理装置及びそのメモリアクセス制御方法に関する。

【背景技術】

【0002】

近年のInternet of Things (IoT)やエッジコンピューティング、第5世代モバイルネットワーク(5G)の登場により、ネットワークを流れるトラフィック量や遅延低減化の要求、ネットワークに接続されるデバイス数、さらには通信の多様性は急速に増加している。通信事業者やサービスプロバイダのネットワークは、その規模や信頼性由来の要件から、従来は用途に特化した専用デバイスや独自のアーキテクチャからなる装置により構成されてきた。

20

【0003】

しかし、近年の急激なトラフィック需要変動に対する柔軟かつ迅速な装置増減設やネットワーク機能の容易な追加実装を可能にするために、通信事業者ネットワークやサービスプロバイダネットワークのような大規模ネットワークにおいても、ネットワーク仮想化(Network Function Virtualization; NFV)やソフトウェア定義ネットワーク(Software Defined Networking; SDN)などの仮想化技術の活用が期待されている。

【0004】

このような仮想化技術活用の機運到来の背景には、従来に比べてより汎用的なデバイスの性能向上がある。CPUやDynamic Random Access Memory (DRAM)といった、汎用的で安価なデバイスからなる汎用コンピュータの性能が向上したことにより、従来は専用装置を用いないと実現困難であった数十ギガビット毎秒級のパケット処理が汎用コンピュータ上のソフトウェアにより実現可能になってきている。したがって、今後、大規模ネットワークにおいても、汎用コンピュータを活用したネットワーク構築により、急激な需要変動や新サービスのための機能追加実装を柔軟・迅速・安価に実現することが可能になると期待される。

30

【0005】

しかし、このような大規模ネットワークにおいては、以降で議論するように、パケット処理のためのテーブル検索等の処理で、現在の汎用コンピュータアーキテクチャではメモリアクセス性能が支配的な性能ボトルネックとなり、これが大規模ネットワークにおける仮想化技術導入の性能観点での障壁になる。

40

【0006】

従来は検索に特化した専用デバイスであるTernary Content-Addressable Memory (TCAM。以降、本表記を使用)を活用することで、パケット処理テーブル検索処理の性能を担保できた。しかし、これは専用デバイスで高価・高消費電力・小容量という課題もあるため、仮想化技術を用いた柔軟かつ低コストな大規模ネットワーク実現に向けて汎用コンピュータに専用デバイスを組み込むというアプローチは望ましくない。

【0007】

50

一方、このメモリアクセス性能高めるデバイスとしてHybrid Memory Cube (HMC。以降、本表記を使用)が2013年4月に仕様が開示され、既にスーパーコンピュータ等の領域で使用されている。HMCは、3次元形状を持つ半導体の層が4~8枚積層され、各層がシリコン貫通電極によって接続されている。その積層した縦の列を“Vault”と呼び、各Vaultは、独立したDRAMベースのメモリであり独立にアクセス可能で並列動作が可能である。また、Vault内には、各層ごとに数個のBankと呼ばれる領域がある。同一Vault内でこれらBankは、共有バスにより接続されているが、共有バス衝突が発生しない範囲内で並列に動作(Bank間interleaving。以降、本表記を使用)可能。このため、汎用メモリデバイスながらきわめて高い性能を実現できる可能性を有している。

10

【0008】

前述のように、パケット処理におけるルーティングやフィルタリング等のためのテーブル検索処理は、特に高いメモリアクセス性能を要求するため、従来のネットワーク装置においては、TCAMのような専用の高速なメモリを使用されている。しかし、上記した仮想化技術を用いた柔軟かつ低コストな大規模ネットワーク実現に向けてTCAMのような専用デバイスを汎用コンピュータに組み込むというアプローチから望ましくないとともに、検索処理に限らず今後より多くのネットワーク機能が仮想化されていくうえでは、汎用コンピュータにおけるメモリアクセスの高性能化が必要である。

【0009】

一方、新しいメモリのアーキテクチャを持つHMCについては、HMC内に検索テーブルを配置して、HMCのもつ広帯域を活用した高速な読み出しに関する検討も存在するが(非特許文献1参照)、TCAMのようなテーブル検索の専用メモリでないため、検索処理に要求されるメモリアクセス性能に達していないとともに、後述する本発明のようにHMCのもつ並列構造を積極的に活用する方式はまだ検討されていない。

20

【0010】

NFVを考慮した汎用コンピュータを適用した従来の技術のパケット処理装置構成には、図20に示すような(1)のDDR \times DRAM及び(2)のHMCを使用したアーキテクチャがある。

【0011】

図20の(1)では、上記したようにDouble-Data-Rate 3(DDR3)DRAMや速度がこの2倍となるDouble-Data-Rate 4(DDR4)DRAMを採用している。最近は、更にDDR4の2倍程度高速なDouble-Data-Rate 5(DDR5)等が次世代メモリとして登場してきている。このような、Double-Data-Rate \times DRAM(DDR \times DRAM。以降、本表記を使用)は、パケット処理においてパケットバッファやアドレス検索テーブル等に使用される。CPUは、マルチコア化されたマルチスレッドでの処理技術が一般化しており、並列処理が可能となっている。また、マルチコアCPUは、各CPUコア内や各CPUコアで共通に使用する低容量で高速動作可能なキャッシュメモリを内蔵しており、キャッシュメモリに納まる範囲内の処理であれば高い処理性能を発揮する。しかしながら、これらキャッシュメモリは、容量が小さく容量不足によりメインメモリであるDDR \times DRAMへのアクセスが頻発した場合、性能のボトルネックが生じる。これは、DDR \times DRAMは、アクセス速度がキャッシュメモリと比較して遅いとともに、アクセスの並列度がないかもしくは並列度があっても低いため、複数のCPUコア側が同時に多くのアクセス要求を出す場合、DDR \times DRAM側がアクセス中でビジー状態となり、CPUコア側で待ち合わせ状態となるためである。

30

40

【0012】

図20の(2)では、メモリとしてHMCを使用し、これをパケットバッファや検索テーブルとして使用している例を示している。HMCアクセス速度は、DDR \times DRAMより高速ではあるが、TCAMのようなテーブル検索の専用メモリでないため、検索処理に要求されるメモリアクセス性能に達しないとともに、後述する本発明のようにHMC

50

のもつ並列構造を積極的に活用する方式はまだ検討されていないため、前述したように、パケット処理装置の仮想化及び将来的にさらなるパケット通信速度の高速化、トラヒックの爆発的な増加、低遅延化等によるメモリアクセス性能不足による性能劣化が主要な性能ボトルネックとなることが予想される。

【先行技術文献】

【非特許文献】

【0013】

【非特許文献1】Packet Matching on FPGAs Using HMC Memory: Towards One Million Rules, Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays

10

【発明の概要】

【発明が解決しようとする課題】

【0014】

大規模な通信事業者ネットワークを汎用コンピュータにより実現し、将来的な大容量トラヒックに対応するため、上述した従来アーキテクチャの延長によるパケット処理方式では、いずれは限界がくと想定される。これは、パケット処理の中でも特に、ルーティング、QoS (Quality of Service)、パケットフィルタリングのようなテーブル検索処理を伴う処理においてメモリアクセス性能不足が顕在化するためである。図20の従来アーキテクチャでのDDRx DRAMやHMC、また専用メモリであるTCAMでは、具体的には、以下が問題となってくる。

20

【0015】

(1) DDRx DRAMを使用したアーキテクチャでは、メモリのアクセス並列度がなくもしくは低い。マルチコアCPUの複数のCPUコアからDDRx DRAMへのアクセスが頻発した場合、アクセス待ち状態によりパケット処理性能のボトルネックになる。

【0016】

(2) HMCは、高速アクセスが可能であるが、単純にHMCを従来のDRAMの代わりに接続したとしてもCPUとのリンク帯域が向上するが、メモリアクセスの並列度が向上するわけではないため、メモリ内のテーブルへのアクセス性能はテーブル検索処理等に求められる水準までは向上しない。このため、従来はTCAMのような専用デバイスを用いる必要があったが、下記のような仮想化における課題が生じてくる。

30

【0017】

(3) 仮想化適用による柔軟な運用や低コスト化のメリットを享受するためには汎用コンピュータ等汎用装置でネットワークが作れることが重要だが、専用デバイスであるTCAMを使わないといけなかった高速テーブル検索などの領域の汎用デバイス化が課題となってくる。また、TCAMは、高価・高消費電力・小容量という課題もある。

【0018】

これらの問題を解決するためには、従来のパケット処理装置アーキテクチャでなく、新しいパケット処理装置アーキテクチャが必要となる。特に、仮想化環境での使用を前提とした、汎用デバイスから構成される汎用装置で、検索テーブルへの高いメモリアクセス性能を実現するパケット処理の具体的な方式の考案が必要である。

40

【課題を解決するための手段】

【0019】

上記目的を達成するために、本願発明は、パケット処理において演算装置からアクセスされるテーブルを記憶した記憶装置と、前記演算装置からの前記記憶装置の前記テーブルへのアクセスリクエストに基づき前記記憶装置へのメモリアクセスを制御する制御装置とを備えたパケット処理装置であって、前記記憶装置の記憶領域は互いに並列アクセス可能なS個(Sは2以上の自然数)のブロックに区画されており、前記各ブロックの記憶領域は互いに並列アクセス可能なN個(Nは2以上の自然数)のバンクに区画されており、前記テーブルは分割されて前記記憶装置の前記バンクに分散して記憶されており、前記制御

50

装置は、前記アクセスリクエストに対してアクセス対象データが格納されているブロック及びバンクを特定して前記記憶装置にアクセスすることを特徴とする。

【発明の効果】

【0020】

本発明によれば、記憶装置に記憶されたテーブルに対する演算装置からのアクセスを制御装置が制御するので、記憶装置として並列アクセス可能な複数のブロック及びバンクに区画された汎用的なものを用いることができる。これにより、仮想化環境での使用を前提とした、汎用デバイスから構成される汎用装置で、検索テーブルへの高いメモリアクセス性能を有するパケット処理を実現できる。

【図面の簡単な説明】

【0021】

【図1】本発明に係るパケット処理装置の概要を示す構成図

【図2】HMC内のテーブル分散配置方式を説明する図

【図3】第1の実施の形態に係るHMCコントローラの機能ブロック図

【図4】(Vault, Bank)対アクセス履歴部の構成例

【図5】HMCコントローラ内の振り分け機構の処理フロー例

【図6】HMCコントローラ内の振り分け機構の処理フロー例

【図7】第2の実施の形態に係るHMCコントローラの機能ブロック図

【図8】HMCコントローラ内の振り分け機構およびテーブル更新制御機構の処理フロー例

【図9】HMCコントローラ内の振り分け機構およびテーブル更新制御機構の処理フロー例

【図10】HMCコントローラ内のテーブル更新制御機構の処理フロー例

【図11】テーブル更新制御機構におけるリクエスト識別部の状態遷移図

【図12】HMC内の負荷追従型テーブル分散配置方式を説明する図

【図13】第3の実施の形態に係るHMCコントローラの機能ブロック図

【図14】分割テーブル配置管理部がもつ配置管理表の一例

【図15】HMCコントローラ内の振り分け機構および負荷追従型テーブル分割変動機構の処理フロー例

【図16】HMCコントローラ内の振り分け機構および負荷追従型テーブル分割変動機構の処理フロー例

【図17】分割変動実施時の処理フロー例

【図18】分割変動実施時における分割テーブルコピー処理の処理フロー例

【図19】分割変動実施時における分割変動リセット時の処理フロー例

【図20】従来のパケット処理装置の構成図

【発明を実施するための形態】

【0022】

まず、本発明の概要について図面を参照して説明する。図1は本発明に係るパケット処理装置の概要を示す構成図である。

【0023】

本発明では、上記の課題を解決するため、図1に示すように、検索テーブルデータの保存に、並列アクセス可能なHybrid Memory Cube (HMC) 300を用いるとともに、CPU 201 + DRAM 202とHMC 300との間に、HMC 300への並列アクセスを可能とするためのHMCコントローラ100を配置するアーキテクチャを提案する。このHMCコントローラ100は、Field Programmable Gate Array (FPGA。以降、本表記を使用)等の再プログラム可能な汎用デバイスで実装する。これにより通信事業者ネットワークのような大規模ネットワークにおけるパケット処理等、高いメモリアクセス性能が求められるアプリケーションにおいて仮想化環境下での使用を想定した汎用デバイスからなる汎用システムにより高性能を実現可能となる。

10

20

30

40

50

【0024】

図1において、CPU201は、複数のCPUコアを有するマルチコアCPUで構成され、内部にキャッシュメモリを内蔵し、これと主メモリ用のDRAM202と接続している。

【0025】

HMC300は、前述したように、Vaultを複数有し(Vault 1~Vault SのS個)、各Vaultは、CPU201側から並列アクセス可能な構造をもつ。より具体的には、HMCは、データ記憶素子層である複数のDRAM層と、メモリコントロール機能を実装した層であるロジックベースとを、Through-Silicon Via(TSV/シリコン貫通電極)と呼ばれる層間接続導体により互いに接続するように積層したものである。HMCは、各データ記憶素子層を平面上において複数の区画に分割するとともに各データ記憶素子層の同一区画間を互いに接続することによりVaultが形成されている。また、一つのVaultは複数(N個)のBankにより構成され、Bank間共有バスにおいて衝突が発生しない範囲で、各Bankは並列アクセスが可能(Bank interleaving。以下、本表記を使用)である。なお、Vaultは、特許請求の範囲の「ブロック」に相当する。

10

【0026】

パケット処理においては、パケット処理プログラム及びパケットバッファは、CPU201に接続されたDRAM202内に設け、パケット処理時間に特に影響する検索テーブルをHMC300内に設ける方式を示している。なお、本発明においては「パケット」とは、例えばInternet Protocol(IP)パケットなどOpen Systems Interconnection(OSI)参照モデルのレイヤー3のパケットを意味するものとする。

20

【0027】

本発明では、以下に述べるように、検索テーブルデータをHMC300内で並列アクセス可能な単位である複数のVaultとBankに分割・分散して配置し、さらに、それらの分散配置された検索テーブルへのCPU201からのメモリアクセスを振り分けるためのHMCコントローラ100をCPU201とHMC300間に設けることにより高速パケット処理方式を実現する。

【0028】

(1)HMC300内の複数あるVaultの1つに、検索テーブル全体を当該Vault内のN個のBank毎にテーブル1~テーブルNに等分割(分割テーブル。以降、本表記を使用)して配置する。この1つのVaultに配置したものと同一内容かつ同一分割のテーブルを残りの全てのS-1個のVaultにコピーする。したがって、各(Vault, Bank)の組み合わせ((Vault, Bank)対。以降、本表記を使用)それぞれが分割テーブルを持ち、Vault間では独立に、同一Vault内のBank間ではBank interleavingとして並列アクセスにアクセスが可能である。また、後述の第2の実施の形態では、アプリケーションにおけるルーティングプロトコル等のやりとりにより、HMC300内のテーブルの内容に更新が発生した場合は、テーブル更新処理を実施する。また、Vault内のBank間でのテーブル分割方法は、等分割を基本とするが、後述する第3の実施の形態では、各分割テーブルへのメモリアクセスの集中度(負荷。以下、本表記を使用)に応じて動的に変動させる。

30

40

【0029】

(2)HMCコントローラ100において、CPU201からのアクセスリクエストに基づいて、上記(1)のHMC300内に分散配置された分割テーブル番号の抽出とこれに該当するアクセス先Bank番号を特定し、さらにこのアクセス先Bank番号をもとにアイドル状態(メモリアクセス中でない状態)のVaultを見つけてアクセス先の(Vault, Bank)対を決定し、当該(Vault, Bank)対へアクセスを振り分ける。また、後述の第2の実施の形態では、HMCコントローラ100において、上記(1)のテーブル更新処理を制御する。また、後述の第3の実施の形態では、HMCコン

50

トローラ100において、上記(1)の負荷に応じたBank間テーブル分割方法の動的変動を制御する。

【0030】

上記(1)及び(2)により、HMC300の各(Vault, Bank)対に分散配置された分割テーブルに並列アクセスが可能となり、また後述の第2の実施の形態では動作中のテーブル更新に対応し、さらに第3の後述の実施の形態ではバーストラヒックのような特定の分割テーブルへの負荷増大に対しても動的な対処が可能となるため、仮想化環境下での使用に適した汎用デバイスを用いて高速なパケット処理が可能となる。

【0031】

以下に本発明の第1～第3の実施の形態について詳述する。

10

【0032】

(第1の実施の形態)

本発明の第1の実施の形態に係るパケット処理装置について図面を参照して説明する。図2は、図1のHMC内のテーブル分散配置方式を具体化したものである。

【0033】

図2に示すように、HMC300が、1つのVault内のN個のBankにおいて、ルーティングテーブルやフローテーブル等の検索テーブル全体をテーブル1からテーブルNまで等分割して配置し、この一つのVaultに配置した検索テーブルを、さらに残りの全てのVaultにコピーして配置する。これにより、同一テーブル番号のアクセスが競合しても複数のVaultに同一内容の検索テーブルがあるため、Vault間の並列動作が可能となる。また、1つのVault内では、Bank間のinterleavingによる並列動作が可能である。これら並列動作機能を高めた方式の採用により、CPUからの検索テーブルアクセス頻度が増大する、より高いレートでのパケット処理が期待できる。

20

【0034】

図3は、図1におけるHMCコントローラ100内の振り分け機構構成を示したものである。この振り分け機構は、図2のHMC内テーブル分散配置方式により分散配置された各分割テーブルへのCPU201からのメモリリクエストを振り分けることにより、検索テーブルへのメモリアクセス並列動作を実現する。

【0035】

図3に示すように、本発明は、CPU201及びDRAM202等から構成するプロセッサ部200と、HMCコントローラ100と、HMC300の3つの主要部分から構成される。

30

【0036】

プロセッサ部200は、CPU201とこれに接続されるプログラムやパケットバッファを有するDRAM202を備える。CPU201には数個～数十個のオーダの複数のCPUコアとこれに内蔵される複数個のキャッシュがある。

【0037】

HMC300は、並列動作できる32個程度のVaultから構成され、それぞれのVaultには16個程度のBankがある。HMC300内のこれらのVaultおよびBankに、図1及び図2を参照して上述した検索テーブルが配置されている。

40

【0038】

このプロセッサ部200とHMC300間に振り分け機構であるHMCコントローラ100を設ける。HMCコントローラ100は、FPGA等の再プログラム可能な汎用デバイスにより構成可能である。

【0039】

HMCコントローラ100には、CPU201からHMC300への検索テーブルアクセスによるメモリリクエストを受け付け、アクセス結果を返すCPUインタフェース部101と、これと接続してメモリリクエストから検索テーブル処理するために必要な情報を抽出するパケット付随情報抽出部102と、この抽出した宛先アドレス等からハッシュ計

50

算によりHMC300の検索テーブルの分割テーブル番号(1~N)を特定する分割テーブル特定部103と、分割テーブル番号からこれと対応するHMC300のBank番号を特定するBank番号特定部104と、Bank番号からHMC300のアクセスするVaultを決定するVault決定部105と、このアクセスするVaultを決定する際に(Vault, Bank)対がアイドル状態(メモリアクセス中でない状態)なのかビジー状態(メモリアクセス中状態)なのかを表示している(Vault, Bank)対アクセス履歴部106と、決定したアクセス先(Vault, Bank)対アドレスをもとにHMC300を実際にアクセスするインタフェース部となるHMCアクセスコントローラ部107とを備えている。

【0040】

図4は、図3の(Vault, Bank)対アクセス履歴部106におけるHMC300内の(Vault, Bank)対が現在、アイドル状態なのかビジー状態なのかを表示するアクセス表示フラグ構成を示す。

【0041】

図4に示すように、マトリックス構成(Bank番号, Vault番号)で行がBank番号を示し、Bank 1からBank Nまであり、列がVault番号を示し、Vault 1からVault Sまでである。現状のHMC300では、前述したように最大でも16×32程度の簡易なマトリックスであり、アイドル状態時が“0”でビジー状態が“1”のフラグ表示構成となっている。本フラグは、HMCアクセス開始時に“1”を立て、HMCアクセス完了時に“0”リセットする。図4では、例として、マトリックス(3, 2)においてBank 3がアクセス該当部となった場合、Vault 2が“0”でアイドル状態であり、アクセス可能な状態にあることを示す。

【0042】

以下、図2~図4の構成をもとに、図5及び図6を参照して、検索テーブル処理の流れについて、検索テーブルのメモリアクセスリクエストの入力から検索結果の出力までについて、図4のHMCコントローラとの処理部位の関連を含めて説明する。図5及び図6は本発明のHMCコントローラ100内の振り分け機構の処理フロー例である。

【0043】

図5及び図6において、図3のプロセッサ部200内のCPU201からHMC300の検索テーブルへのアクセスに伴うメモリアクセスリクエストを受け付け、振り分け機構の処理を開始する(ステップS101)。

【0044】

CPUインタフェース部101では、受け付けた検索テーブルへのメモリアクセスリクエストをパケット付随情報抽出部102に転送する(ステップS102)。

【0045】

これを受信したパケット付随情報抽出部102では、メモリアクセスリクエスト内容に応じて検索テーブル処理に必要な情報を抽出する(ステップS103)。例えば、ルーティングテーブル検索では、メモリアクセスリクエスト内での宛先IPアドレス情報を抽出する。

【0046】

この抽出した宛先アドレス等の情報をもとに分割テーブル特定部103では、ハッシュ計算によりHMCの検索テーブルの分割テーブル番号(1~N)を特定する(ステップS104)。

【0047】

この分割テーブル特定部103で特定した分割テーブル番号からBank番号特定部104では、これと対応するHMC300内のアクセスするBank番号を特定する(ステップS105)。

【0048】

次に、Bank番号を受信したVault決定部105では、Bank番号からHMC300のアクセスするVaultを決定するために該当するBank番号をアクセスして

10

20

30

40

50

いないアイドル状態の Vault を見つけるため、(Vault , Bank) 対アクセス履歴部 106 にアイドル状態の参照要求を出す (ステップ S 106)。

【0049】

この参照要求を受信した (Vault , Bank) 対アクセス履歴部 106 では、図 4 に示す (Vault , Bank) 対のメモリアクセスしていないアイドル状態かメモリアクセス中のビジー状態かを表示するアクセス表示フラグを該当アクセス Bank 部分について順次確認する (ステップ S 107 , S 108)。全部アクセス中でビジー状態である場合 (全てフラグ “ 1 ”)、一定時間 W (1 ~ 数クロック程度) 待機し (ステップ S 110)、再びフラグを順次確認する。そして、アイドル状態を最初に見つけた Vault 番号を参照番号結果として Vault 決定部 105 に返送する (ステップ S 109)。この返送直後に、このフラグを “ 1 ” としてビジー状態にする (ステップ S 111)。

10

【0050】

アクセスする Vault 番号を参照結果として受け取った Vault 決定部 105 では、アクセスする Bank 番号と Vault 番号の対番号を HMC アクセスコントローラ部 107 にアクセス要求する (ステップ S 112)。

【0051】

これを受信した HMC アクセスコントローラ部 107 では、この (Vault , Bank) 対番号より HMC 300 の該当アドレスを割り出して、HMC に対してアクセス要求を出す (ステップ S 113)。このアクセスにおいて HMC 300 から検索アクセス応答の状態を監視し (ステップ S 114)、アクセス応答が正常である場合には、アクセス結果を Vault 決定部 105 に返却転送する (ステップ S 115)。

20

【0052】

これを受信した Vault 決定部 105 では、(Vault , Bank) 対アクセス履歴部 106 の該当アクセス表示フラグにアクセス完了リセット指示するとともに、アクセス検索結果を CPU インタフェース部 101 側に返送する (ステップ S 116)。

【0053】

Vault 決定部 105 からのアクセス完了リセット指示により (Vault , Bank) 対アクセス履歴部 106 では、該当アクセス表示フラグを “ 0 ” リセットし、アイドル状態とする (ステップ S 117)。

【0054】

HMC 300 からのアクセス応答が異常でエラー状態であった場合 (ステップ S 114) には、アクセス結果をエラーとして Vault 決定部 105 に返却する (ステップ S 118)。Vault 決定部 105 では、これをアクセスエラーとして CPU 201 側に返送する (ステップ S 119)。CPU 201 では、エラー内容に応じてアプリケーションレベルで適宜エラー処理を行う。

30

【0055】

このようなパケット処理装置によれば以下の効果が生じる。

【0056】

(1) 高速パケットの検索テーブルの分散配置により HMC 300 のもつメモリ容量を効率的に利用できるとともに、並列処理により汎用的なデバイスのみを活用してサーバ上のパケット処理性能向上を飛躍的に大きく拡大できる。

40

【0057】

(2) HMC コントローラおよび振り分け機構を FPGA 等の再プログラム可能な汎用デバイスで実現することにより CPU 201 や HMC 300 などのデバイス自体の変更は不要である。

【0058】

(3) HMC 300 を含め汎用デバイスから成る汎用コンピュータによるシステム構成であるため、幅広い既存パケット処理ソフトウェアをより高速に動作させることが可能である。

【0059】

50

(4) TCAMに比べて低消費電力なDRAMベースのHMC300の採用により、システム全体の消費電力削減や実装面積削減によるコンパクト化が可能となる。

【0060】

(第2の実施の形態)

本発明の第2の実施の形態に係るパケット処理装置について図面を参照して説明する。本実施の形態が第1の実施の形態と異なる点は、第1の実施の形態に係るHMCコントローラ100に対して、さらに、HMC300に記憶されたテーブルを更新するテーブル更新制御機能を設けた点にある。他の点については第1の実施の形態と同様なので、ここでは主として相違点について説明する。なお、第1の実施の形態と同様の構成については同一の符号を付した。

10

【0061】

図7は、図1におけるHMCコントローラ100内の振り分け機構およびテーブル更新制御機構構成を示したものである。この振り分け機構については、第1の実施の形態と同様である。一方、テーブル更新制御機構は、振り分け機構と連携し、アプリケーション側からの要求に応じてHMC300内のテーブル更新処理を制御することにより、動作中のテーブル更新を可能とするものである。

【0062】

図7に示すように、HMCコントローラ100は、CPU201からHMC300への検索テーブルアクセスによるメモリリクエストを受け付け、アクセス結果を返すCPUインタフェース部101と、これと接続してメモリリクエストが、通常のHMC300内のテーブルへのメモリアクセスリクエスト(リードアクセス)であるのかテーブル更新のためのリクエストであるのかを識別するリクエスト識別部111と、リクエストがHMC300内のテーブルへのメモリアクセスリクエストの場合にメモリリクエストからテーブル検索処理するために必要な情報を抽出するパケット付随情報抽出部102と、この抽出した宛先アドレス等からハッシュ計算によりHMC300の検索テーブルの分割テーブル番号(1~N)を特定する分割テーブル特定部103と、分割テーブル番号からこれと対応するHMC300のBank番号を特定するBank番号特定部104と、Bank番号からHMC300のアクセスするVaultを決定するVault決定部105と、このアクセスするVaultを決定する際に(Vault, Bank)対がアイドル状態(メモリアクセス中でない状態)なのかビジー状態(メモリアクセス中状態)なのかを表示している(Vault, Bank)対アクセス履歴部106と、決定したアクセス先(Vault, Bank)対アドレスをもとにHMC300を実際にアクセスするインタフェース部となるHMCアクセスコントローラ部107と、リクエストがHMC300内のテーブル更新である場合にそのテーブル更新処理を制御するテーブル更新制御部112とを備えている。

20

30

【0063】

以下、図2及び図4並びに図7の構成をもとに、図8~図10を参照して、検索テーブル処理の流れについて、検索テーブルのメモリアクセスリクエストの入力から検索結果の出力までおよびテーブル更新の動作について、図7のHMCコントローラとの処理部位の関連を含めて説明する。図8及び図9は本発明のHMCコントローラ100内の振り分け機構およびテーブル更新制御機構の処理フロー例である。また、図10は、本発明のテーブル更新制御機構の処理フロー例(詳細)である。

40

【0064】

本実施の形態では、HMCコントローラ100がCPU201から受け付けるメモリリクエストは、テーブルへのアクセスリクエスト(リードリクエスト)と、テーブル更新に係るメモリリクエストとに大別される。さらに、後者のテーブル更新に係るメモリリクエストは、更新開始リクエストと、更新内容を含むテーブル更新内容リクエスト、更新終了リクエストとに別れる。テーブル更新の際には、HMCコントローラ100は、更新開始リクエスト、更新内容を含む1つ又は複数のテーブル更新内容リクエスト、更新終了リクエストを順に受信する。

50

【0065】

図8及び図9において、図7のプロセッサ部200内のCPU201からHMC300の検索テーブルへのアクセスに伴うメモリリクエストを受け付け、振り分け機構の処理を開始する(ステップS201)。

【0066】

CPUインタフェース部101では、受け付けた検索テーブルへのメモリアクセスリクエストをリクエスト識別部111に送付し(ステップS202)、リクエスト識別部111ではリクエストがHMC300内のテーブルへのアクセス(リードアクセス)であるのかテーブル更新に関するものであるのかを識別する(ステップS203)。識別の結果、HMC300内のテーブルへのアクセス(リードアクセス)である場合は、リクエストをパケット付随情報抽出部102に転送(ステップS204)し、テーブル更新に関するものである場合は図10の本発明のテーブル更新制御機構の処理フロー例(詳細)により説明するフロー例に従って処理を実施する。

10

【0067】

以降の処理(図8のS205以降及び図9の処理)は、HMC300内のテーブルへのアクセス(リードアクセス)についての処理であり、当該処理については第1の実施の形態と同様なのでここでは説明は省略する。

【0068】

テーブル更新処理のフロー例について、図10および図11の本発明のテーブル更新制御機構におけるリクエスト識別部の状態遷移図を用いて説明する。

20

【0069】

リクエスト識別部111において、HMC300内のテーブルへのメモリアクセス(リードアクセス)以外と判断されたリクエストについて、さらに、テーブル更新開始を示すリクエストであるかどうかを識別する(ステップS251)。テーブル更新開始のリクエストであると判断された場合は、リクエスト識別部111の状態を図11のようにテーブル更新モードに遷移させ、テーブル更新をテーブル更新制御部112に指示(ステップS252)、また、テーブル更新開始のリクエストでないと判断された場合は、テーブル更新内容のリクエストであるかテーブル更新終了のリクエストであるのかを識別する(ステップS257)。

【0070】

テーブル更新モードに遷移した場合、テーブル更新制御部112では、テーブル更新を実施するためにHMC300内のテーブルへのメモリアクセス状況の参照を要求する(ステップS253)。(Vault, Bank)対アクセス履歴部106では、全Vaultのメモリアクセス状況を順次確認(ステップS254)し、全(Vault, Bank)対がアイドル状態かどうかを確認(ステップS255)し、すべてアイドル状態であれば当該リクエストの処理を完了、そうでなければ一定時間 W_2 だけ待機(ステップS256)し、再度アクセス状況を確認する。

30

【0071】

前記確認(ステップS257)において、テーブル更新内容のリクエストであると判断された場合は、テーブル更新制御部112において、テーブル更新のためのHMCアクセスを指示(ステップS258)し、HMCアクセスコントローラ部107において指定されたアドレスへのHMCアクセスを実施してテーブルを更新する(ステップS259)。テーブル更新制御部112では、アクセス結果をCPU201へ通知し、当該リクエストの処理を完了する(ステップS260)。

40

【0072】

また、前記確認(ステップS257)において、テーブル更新終了のリクエストであると判断された場合は、リクエスト識別部111の状態を図11のようにVault間コピーモードに遷移させ、Vault間コピーをテーブル更新制御部112に指示する(ステップS261)。テーブル更新制御部112では、Vault間コピーを指示(ステップS262)し、HMCアクセスコントローラ部107において指定されたアドレスへのH

50

MCアクセスを実施して、Vault間のテーブルデータを同期する(ステップS263)。テーブル更新制御部112では、アクセス結果をCPU201へ通知し、リクエスト識別部111の状態を図11のように通常モードへ遷移させる(ステップS264, S265)。これにより、テーブル更新処理が完了し、当該リクエストの処理を完了する。

【0073】

本実施の形態に係るパケット処理装置では、ルーティングプロトコルなどのアプリケーションによって発生するテーブル更新処理を、HMC300内に分散配置された検索テーブルについてもシステム動作中に実施することが可能となる。他の効果については第1の実施の形態と同様である。

【0074】

(第3の実施の形態)

本発明の第3の実施の形態に係るパケット処理装置について図面を参照して説明する。本実施の形態が第1の実施の形態と異なる点は、メモリアクセスの負荷に追従してHMC300内におけるテーブルの分割形態を動的に変化させる点にある。他の点については第1の実施の形態と同様なので、ここでは主として相違点について説明する。なお、第1の実施の形態と同様の構成については同一の符号を付した。

【0075】

図12は、その初期状態を示した、HMC300内の負荷追従型テーブル分散配置方式を具体化したものである。HMC300内へのテーブルの分散配置は、負荷に応じて動的に変動させるが、初期状態では、1つのVault内のN個のBankにおいてルーティングテーブルやフローテーブル等の検索テーブル全体をテーブル1からテーブルNまで等分割して配置し、この一つのVaultに配置した検索テーブルを、さらに残りの全てのVaultにコピーして配置する。これにより、同一テーブル番号のアクセスが競合しても複数のVaultに同一内容の検索テーブルがあるため、Vault間の並列動作が可能となる。また、1つのVault内では、Bank間のinterleavingによる並列動作が可能である。これら並列動作機能を高めた方式の採用により、CPUからの検索テーブルアクセス頻度が増大する、より高いレートでのパケット処理が期待できる。

【0076】

本実施の形態では、このようにして分散配置した分割テーブルのうち、負荷が最大のものについて、その負荷があらかじめプログラムされた閾値を超えていることが検出された場合、同一Vault内で最大負荷の分割テーブルを最小負荷の分割テーブルが配置されているBankへ上書きコピー配置(Vault内分割変動。以降、本表記を使用)する。このVault内分割変動は、1つのHMC内のVault 1から実施し、全Vault数Sに対して、あらかじめプログラムされた、負荷に応じてVault内分割変動を許容するVault数 S_{var} 分だけ実施するまで繰り返す。Vault($S_{var} + 1$)からVault Sまでの($S - S_{var}$)個のVaultについては、Vault内分割変動を行わない。また、負荷検出は、CPUからのメモリアクセスをR個受信するとに実施する。ここでRはあらかじめ定められた値であり、且つ、プログラマブルである。

【0077】

図13は、図2におけるHMCコントローラ内の振り分け機構および負荷追従型テーブル分割変動機構構成を示したものである。この振り分け機構は、図12のHMC内テーブル分散配置方式により分散配置された各分割テーブルへのCPU201からのメモリリクエストを振り分けることにより、検索テーブルへのメモリアクセス並列動作を実現する。また、負荷追従型テーブル分割変動機構は、振り分け機構と連携し、分割テーブルへのメモリアクセス負荷を監視し、負荷の集中を検知した場合は、HMC300内のテーブル分割を変動させることにより、特定の分割テーブルへの負荷集中に対処することを可能とするものである。

【0078】

図13に示すように、HMCコントローラ100は、CPU201からHMC300へ

10

20

30

40

50

の検索テーブルアクセスによるメモリリクエストを受け付け、アクセス結果を返すCPUインタフェース部101と、これと接続してメモリリクエストからテーブル検索処理するために必要な情報を抽出するパケット付随情報抽出部102と、この抽出した宛先アドレス等からハッシュ計算によりHMC300の検索テーブルの分割テーブル番号(1~N)を特定する分割テーブル特定部103と、分割テーブル番号をもとに各分割テーブルの負荷を監視する負荷監視部121と、負荷情報からテーブル分割を変動すべきか判断しテーブル分割実施時にはその制御を行う分割変動制御部122と、この負荷情報からテーブル分割を変動すべきか判断を行う際に必要な閾値をあらかじめプログラムしておき必要に応じて参照する負荷閾値部123と、HMC300内における分割テーブルの配置状況を示す配置管理表を管理する分割テーブル配置管理部124と、分割テーブル番号及び配置管理表からアクセス対象候補となるHMC300の一以上のBank番号及びVault番号を特定するBank番号特定部104と、アクセス候補である一以上のBank番号及びVault番号からHMC300のアクセスするVaultを決定するVault決定部105と、このアクセスするVaultを決定する際に(Vault, Bank)対がアイドル状態(メモリアクセス中でない状態)なのかビジー状態(メモリアクセス中状態)なのかを表示している(Vault, Bank)対アクセス履歴部106と、決定したアクセス先(Vault, Bank)対アドレスをもとにHMC300を実際にアクセスするインタフェース部となるHMCアクセスコントローラ部107とから構成される。

10

20

【0079】

図14に分割テーブル配置管理部124の有する配置管理表の一例を示す。配置管理表は、図14に示すように、(Vault, Bank)対で特定される記憶領域と当該記憶領域に記憶されている分割テーブルの番号及び分割変動の実施状況との対応関係を示す。本実施の形態では、分割変動の実施状況は「分割変動未実施」「分割変動実施済み」の2つの状況を含む。また、「分割変動実施済み」の場合は、付随状態として、分割変動時の負荷情報を含む。負荷情報は、「負荷最小」「負荷最大」を含む。図14の例は、分割変動によりbank iの分割テーブルをbank kに上書きコピーした場合を示している。

【0080】

以下、図4、図12~図14の構成をもとに、検索テーブル処理の流れについて図15~図19を用いて、検索テーブルへのメモリアクセスリクエストの入力から検索結果の出力およびテーブル分割変動の動作について図13のHMCコントローラとの処理部位の関連を含めて説明する。図15及び図16は本発明のHMCコントローラ内の振り分け機構および負荷追従型テーブル分割変動機構の処理フロー例であり、図17及び図18は分割変動実施時の処理フロー例、図19は分割変動リセット時の処理フロー例である。

30

【0081】

図15及び図16において、図13のプロセッサ部200内のCPU201からHMC300の検索テーブルへのアクセスに伴うメモリリクエストを受け付け、振り分け機構部の処理を開始する(ステップS301)。

【0082】

CPUインタフェース部101では、受け付けた検索テーブルへのメモリアクセスリクエストをパケット付随情報抽出部102に転送する(ステップS302)。

40

【0083】

これを受信したパケット付随情報抽出部102では、メモリアクセスリクエスト内容に応じて検索テーブル処理に必要な情報を抽出する(ステップS303)。例えば、ルーティングテーブル検索では、メモリアクセスリクエスト内での宛先IPアドレス情報を抽出する。

【0084】

この抽出した宛先アドレス等の情報をもとに分割テーブル特定部103では、ハッシュ計算によりHMCの検索テーブルの分割テーブル番号(1~N)を特定する(ステップS

50

304)。

【0085】

負荷監視部121及び分割変動制御部122では、図17～図19を用いて後述する分割変動処理を負荷閾値部123、分割テーブル配置管理部124、(Vault, Bank)対アクセス履歴部106、HMCアクセスコントローラ部107と連携して実施する(ステップS305)。

【0086】

分割テーブル特定部103で特定した分割テーブル番号及び分割テーブル配置管理部124が有する配置管理表に基づき、Bank番号特定部104では、アクセス候補となるHMC内のBank番号及びVault番号を特定する(ステップS306)。

10

【0087】

次に、アクセス候補となるBank番号及びVault番号を受信したVault決定部105では、アクセス候補のうちBank番号及びVault番号をアクセスしていないアイドル状態の(Vault, Bank)対を見つけるため、(Vault, Bank)対アクセス履歴部106にアイドル状態の参照要求を出す(ステップS307)。

【0088】

この参照要求を受信した(Vault, Bank)対アクセス履歴部106では、図4に示す(Vault, Bank)対のメモリアクセスしていないアイドル状態かメモリアクセス中のビジー状態かを表すアクセス表示フラグを該当アクセスBank部分について順次確認する(ステップS308, S309)。全部アクセス中でビジー状態である場合(全てフラグ“1”)、一定時間 W_1 (1～数クロック程度)待機し(ステップS311)、再びフラグを順次確認する。アイドル状態を最初に見つけたBank番号及びVault番号を参照番号結果としてVault決定部105に返送する(ステップS310)。

20

この返送直後に、このフラグを“1”としてビジー状態にする(ステップS312)。

【0089】

アクセスするBank番号及びVault番号を参照結果として受け取ったVault決定部105では、アクセスするBank番号とVault番号の対番号をHMCアクセスコントローラ部107にアクセス要求する(ステップS313)。

【0090】

これを受信したHMCアクセスコントローラ部107では、この(Vault, Bank)対番号よりHMC300の該当アドレスを割り出して、HMC300に対してアクセス要求を出す(ステップS314)。このアクセスにおいてHMC300から検索アクセス応答の状態を監視し(ステップS315)、アクセス応答が正常である場合には、アクセス結果をVault決定部105に返却転送する(ステップS316)。

30

【0091】

これを受信したVault決定部105では、(Vault, Bank)対アクセス履歴部106の該当アクセス表示フラグにアクセス完了リセット指示するとともに、アクセス検索結果をCPUインタフェース部101側に返送する(ステップS318)。

【0092】

Vault決定部105からのアクセス完了リセット指示により(Vault, Bank)対アクセス履歴部106では、該当アクセス表示フラグを“0”リセットし、アイドル状態とする(ステップS319)。

40

【0093】

HMC300からのアクセス応答が異常でエラー状態であった場合(ステップS315)には、アクセス結果をエラーとしてVault決定部105に返却する。Vault決定部105では、これをアクセスエラーとしてCPU201側に返送する(ステップS320)。CPU201では、エラー内容に応じてアプリケーションレベルで適宜エラー処理を行う。

【0094】

50

前記分割変動処理（ステップS305）の詳細について図17～図19を参照して説明する。なお、図17において処理主体が特に明示していない処理は、分割変動制御部122が行うものとする。

【0095】

図17に示すように、負荷監視部121は、分割テーブルごとの到着リクエスト数（負荷）を常時カウントし、全分割テーブル合計の累計到着リクエスト数がRに到達した場合、負荷最大値と最小値、対応する分割テーブル番号（負荷情報）を出力するとともに、到着リクエスト数のカウンタをリセットする（ステップS401～S403）。一方、累計到着リクエスト数がRに到達しない場合は、処理を終了する。すなわち、分割変動処理は、R個のリクエスト到着ごとに実施する。

10

【0096】

分割変動制御部122は、負荷最大値が所定の負荷閾値を超えている否かにより分割変動の実施要否を判定する（ステップS404）。

【0097】

分割変動制御部122は、負荷最大値が所定の負荷閾値を超えていない場合、分割変動実施済みか確認するために分割テーブル配置管理部124へ参照要求を行う（ステップS405）。分割テーブル配置管理部124は、配置管理表を参照して分割変動実施済みか否かの結果を分割変動制御部122に返却する（ステップS406）。分割変動制御部122は、返却結果が分割変動実施済みであった場合、後述する分割変動リセット処理を実施して処理を終了する（ステップS407）。分割変動制御部122は、返却結果が分割変動実施済みでない場合、処理を終了する。

20

【0098】

一方、分割変動制御部122は、負荷最大値が所定の負荷閾値を超えている場合、分割変動実施済みか確認するために分割テーブル配置管理部124へ参照要求を行う（ステップS408）。分割テーブル配置管理部124は、配置管理表を参照して分割変動実施済みか否かの結果を分割変動制御部122に返却する（ステップS409）。ここで、分割変動実施済みの場合には、分割テーブル配置管理部124は、前回分割変動実施した分割テーブルが今回の負荷最大となる分割テーブルか否かの判定結果を分割変動制御部122に返却する（ステップS410）。分割変動制御部122は、分割変動実施済みであり且つ前回分割変動実施した分割テーブルが今回の負荷最大となる分割テーブルである場合には、処理を終了する。

30

【0099】

一方、分割変動制御部122は、分割変動実施済みであり且つ前回分割変動実施した分割テーブルが今回の負荷最大となる分割テーブルでない場合には、分割変動リセット処理を実施し（ステップS411）、分割テーブルのコピー処理を実施する（ステップS412）。また、分割変動制御部122は、分割変動実施済みでない場合、分割テーブルのコピー処理を実施する（ステップS412）。

【0100】

分割テーブルのコピー処理（ステップS412）は、Vault 1から順にVault S_{var}まで、負荷最大となる分割テーブルの内容を負荷最小分割テーブルが配置されているbankへ上書きコピーする処理である。

40

【0101】

具体的には、図18に示すように、まず、分割変動制御部122において、コピー処理が初回かどうかの確認（ステップS501）により初回である場合はコピー処理を行うVault番号mをm=1としてVault 1から処理を開始（ステップS502）、初回ではない場合はm < S_{var}であるか確認（ステップS503）し、m < S_{var}である場合は、mを1増やして別のVaultに対してコピー処理を開始（ステップS504）、m < S_{var}でない場合は、コピー処理が完了したと判断し、本フローを終了する。

【0102】

Vault mにおいてコピー処理を実施するために、Vault mのアクセス状況

50

参照指示を出す（ステップS505）。（Vault, Bank）対アクセス履歴部106において、Vault mのアクセス状況を確認する（ステップS506）。次に、Vault mの全Bankがアイドル状態かどうか確認し（ステップS507）、アイドル状態でない場合は一定時間 W_2 待機（ステップS508）したのち、再度アクセス状態から確認する（ステップS507）。

【0103】

Vault mがアイドル状態である場合は、分割変動制御部122において、Vault mのアクセス抑止を指示する（ステップS509）。（Vault, Bank）対アクセス履歴部106において、Vault mの全Bankのアクセス状況フラグを“1”にし、アクセスを抑止する（ステップS510）。分割変動制御部122から、Vault mの内の最大負荷分割テーブルの内容を最小負荷分割テーブルが配置されているbankへコピーして配置する指示を出す（ステップS511）。

10

【0104】

HMCアクセスコントローラ部107では、指示されたアドレスへのHMCアクセスを要求する（ステップS512）。アクセス結果を確認し（ステップS513）、正常な場合は分割変動制御部122においてVault mのアクセス抑止を解除し（ステップS514）、本処理フローの最初にもどる。確認（ステップS513）の結果が異常の場合は、本処理フローを終了する。

【0105】

前述の分割変動リセット処理（ステップS407, S411）について図19を参照して説明する。分割変動制御部122において、分割変動リセット処理が初回かどうかの確認（ステップS601）により初回である場合は分割変動リセットを行うVault番号mを $m=1$ としてVault 1から処理を開始（ステップS602）、初回ではない場合は $m < S_{var}$ であるか確認（ステップS603）し、 $m < S_{var}$ である場合は、mを1増やして別のVaultに対して分割変動リセット処理を開始（ステップS604）、 $m < S_{var}$ でない場合は、分割変動処理対象である全Vaultについて分割変動リセット処理実施済みであると判断し、後述する分割変動リセット処理時の元データとするVault pのアクセス抑止を解除（ステップS616）して本処理フローを終了する。

20

【0106】

Vault mにおいて分割変動リセット処理を実施するために、Vault mおよび分割変動処理対象外であるVault ($S_{var} + 1$) ~ Sのアクセス状況参照指示を出す（ステップS605）。（Vault, Bank）対アクセス履歴部106において、Vault mのアクセス状況の確認およびVault ($S_{var} + 1$) ~ Sのアクセス状況を順次確認し、ビジー状態のBank数が最小であるVault pを選出する（ステップS606）。次に、Vault mおよびpの全Bankがアイドル状態かどうか確認し（ステップS607）、アイドル状態でない場合は一定時間 W_3 待機（ステップS608）したのち、再度アクセス状況参照指示を出す（ステップS605）。

30

【0107】

Vault mおよびpの全Bankがアイドル状態である場合は、分割変動制御部122において、Vault mおよびpのアクセス抑止を指示する（ステップS609）。（Vault, Bank）対アクセス履歴部106において、Vault mおよびpの全Bankのアクセス状況フラグを“1”にし、アクセスを抑止する（ステップS610）。分割変動制御部122から、Vault pの全Bankの内容をVault mにコピーする指示を出す（ステップS611）。

40

【0108】

HMCアクセスコントローラ部107では、指示されたアドレスへのHMCアクセスを要求する（ステップS612）。アクセス結果を確認し（ステップS613）、正常な場合は分割変動制御部122においてVault mのアクセス抑止を解除し（ステップS614）、分割変動実施済みで未リセットのVaultがあるか確認する（ステップS6

50

15)。確認(ステップS615)の結果、未リセットのV a u l tがある場合は、本処理フローの最初にもどり、無い場合は(V a u l t , B a n k)対アクセス履歴部106においてV a u l t pのアクセス状況フラグをすべて“0”にセットして、アクセス抑止を解除し(ステップS616)、本処理フローを終了する。確認(ステップS613)の結果が異常の場合は、本処理フローを終了する。

【0109】

本実施の形態に係るパケット処理装置では、バーストラヒック等により特定の分割テーブルにメモリアクセスが集中する場合において局所的な対処を行っているので、システム全体のスループットを最大化することができる。他の効果については第1の実施の形態と同様である。

10

【0110】

なお、本実施の形態は第1の実施の形態の変形例として説明したが、第2の実施の形態の変形例とすることもできる点に留意されたい。

【0111】

以上、本発明の実施の形態について詳述したが本発明はこれに限定されるものではない。HMCコントローラ100における各機能の実装形態やアルゴリズムは不問であり、他の実装形態やアルゴリズムであっても本発明を適用できる。例えば、上記各実施の形態では、アクセス対象とするアイドル状態のB a n kを特定する際に、常にV a u l t 1から検索していたが、メモリリクエスト毎に検索開始V a u l t 番号を所定のルールで又はランダムで変更するようにしてもよい。

20

【0112】

また、上記第3の実施の形態では、分割変動処理時にはS v a r個のV a u l tについて負荷最大の分割テーブルを負荷最小の分割テーブルに上書きコピーしていた。ここでS v a rは、あらかじめ定められたプログラマブルな値である。他の変形例としては、S v a rを、負荷情報に応じて可変とするようにしてもよい。より具体的には、負荷が大きいほどS v a rを大きく設定し、負荷が小さいほどS v a rを小さくするよう動的に制御してもよい。

【0113】

また、上記各実施の形態では、記憶装置の一例としてHMCについて説明したが、並列アクセス可能なブロック(V a u l t)及びバンク構成を有する他の構造・規格の記憶装置であっても本発明を適用できる。

30

【符号の説明】

【0114】

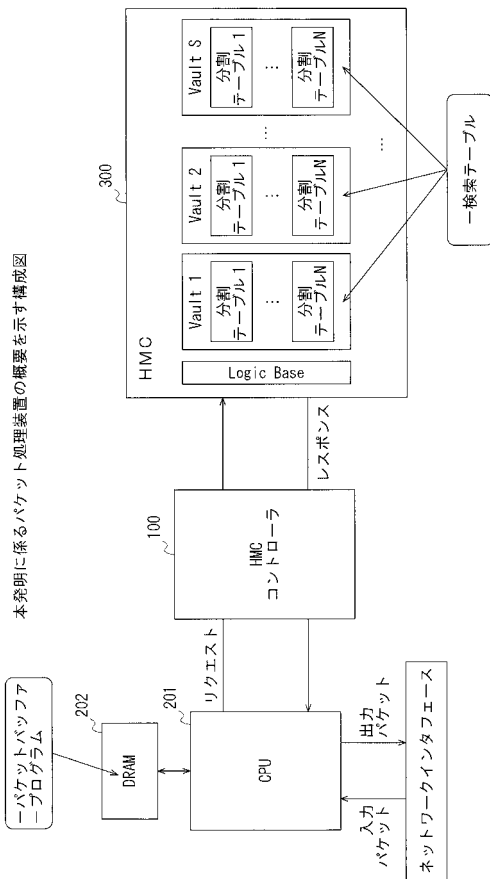
- 100 ... HMCコントローラ
- 101 ... CPUインタフェース部
- 102 ... パケット付随情報抽出部
- 103 ... 分割テーブル特定部
- 104 ... B a n k 番号特定部
- 105 ... V a u l 決定部
- 106 ... (V a u l t , B a n k) 対アクセス履歴部
- 107 ... HMCアクセスコントローラ部
- 111 ... リクエスト識別部
- 112 ... テーブル更新制御部
- 121 ... 負荷監視部
- 122 ... 分割変動制御部
- 123 ... 負荷閾値部
- 124 ... 分割テーブル配置管理部
- 200 ... プロセッサ部
- 201 ... CPU
- 202 ... DRAM

40

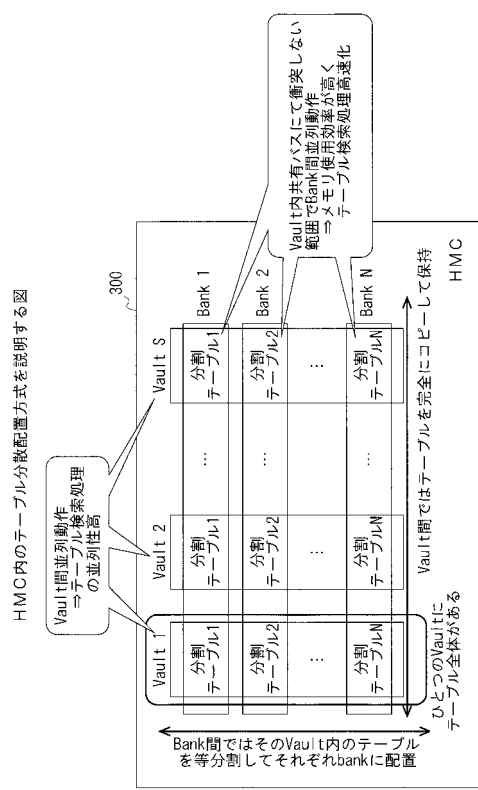
50

3 0 0 ... H M C

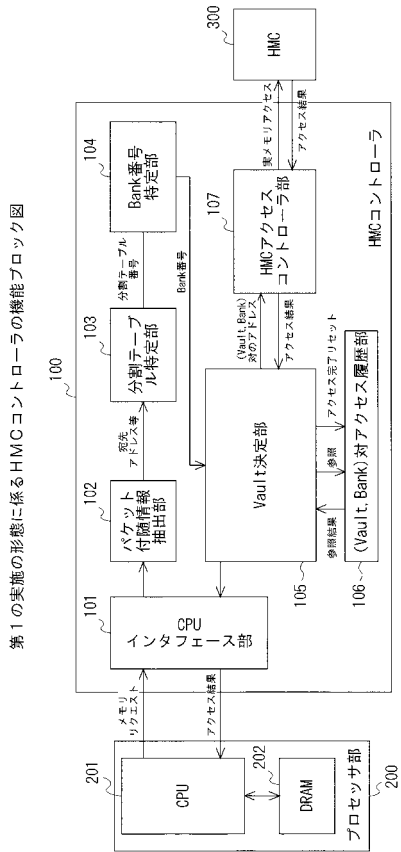
【 図 1 】



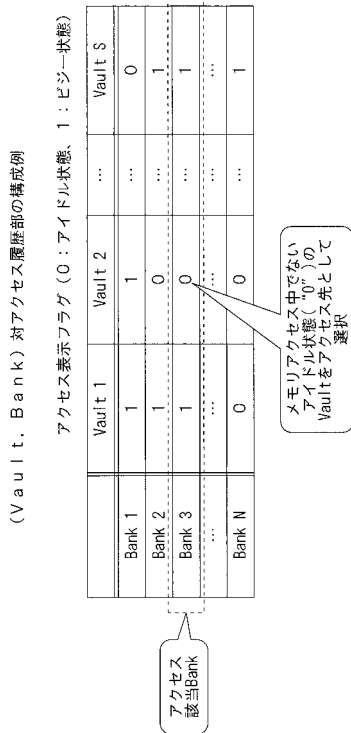
【 図 2 】



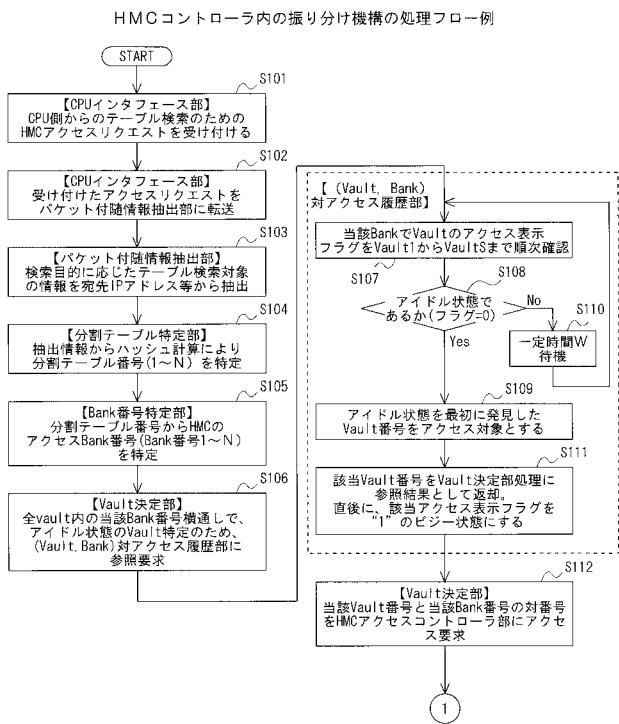
【図3】



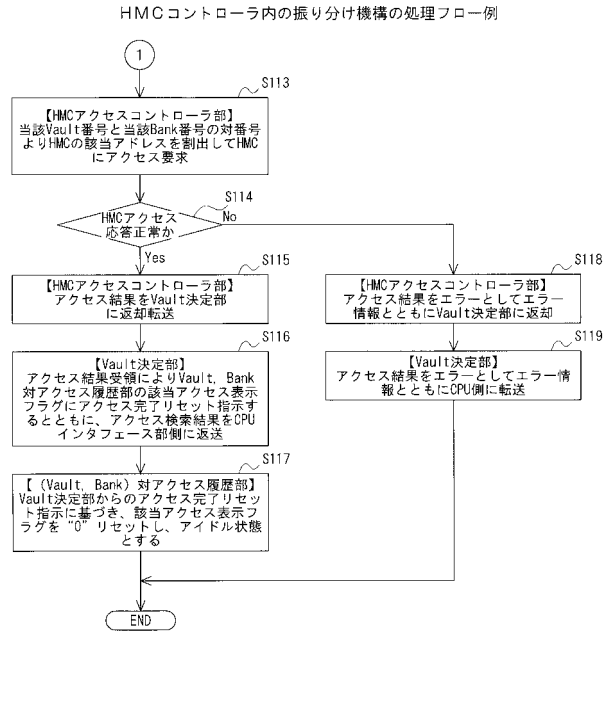
【図4】



【図5】

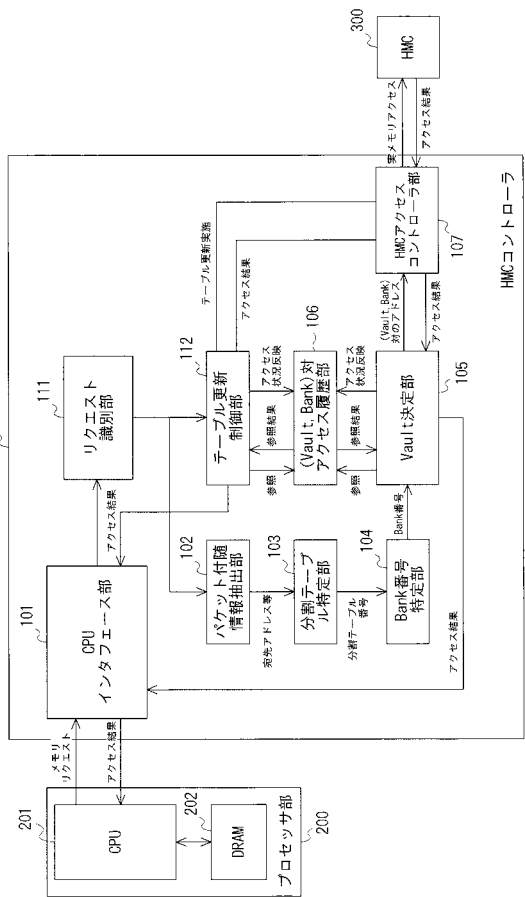


【図6】



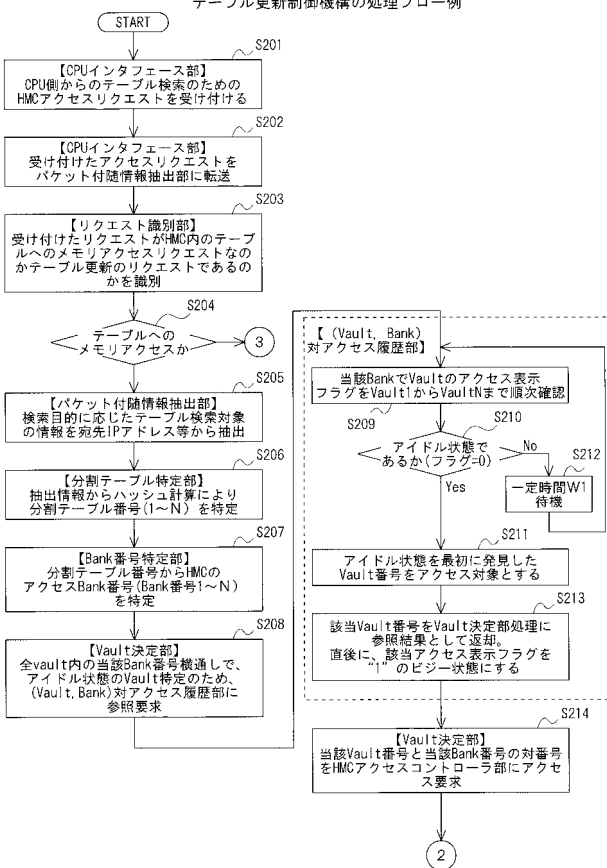
【図7】

第2の実施の形態に係るHMCコントローラの機能ブロック図



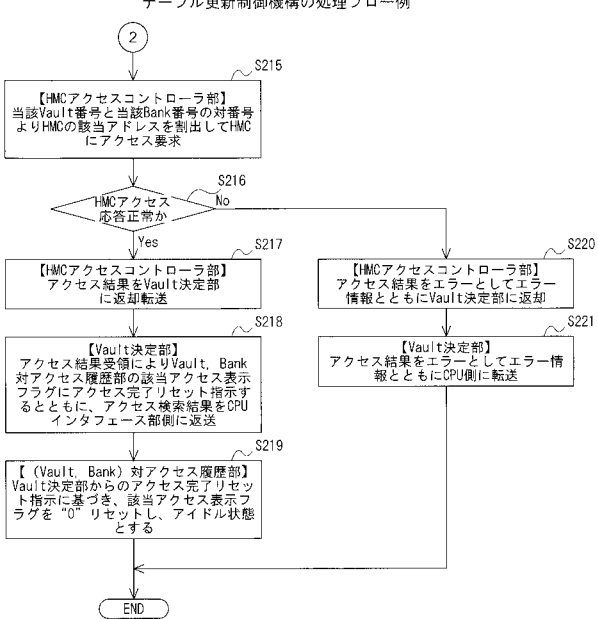
【図8】

HMCコントローラ内の振り分け機構およびテーブル更新制御機構の処理フロー例



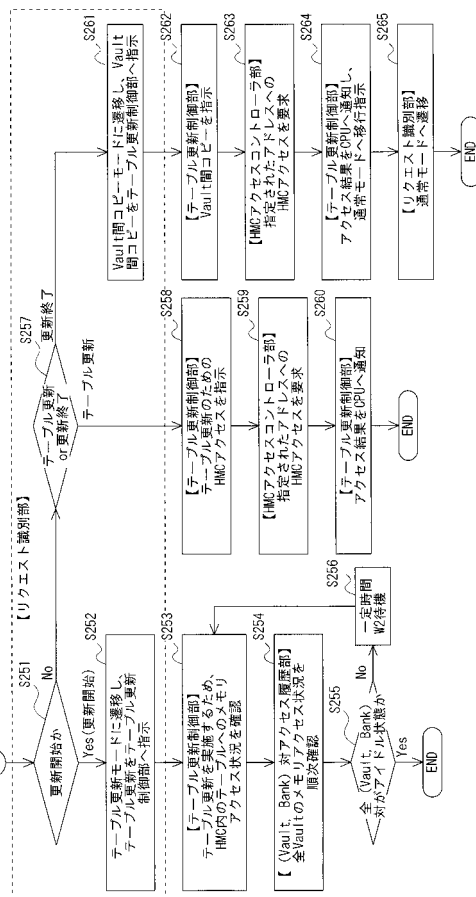
【図9】

HMCコントローラ内の振り分け機構およびテーブル更新制御機構の処理フロー例

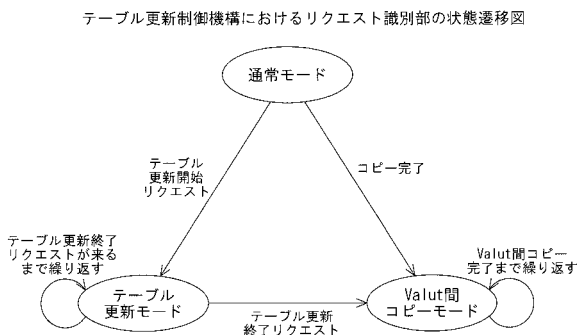


【図10】

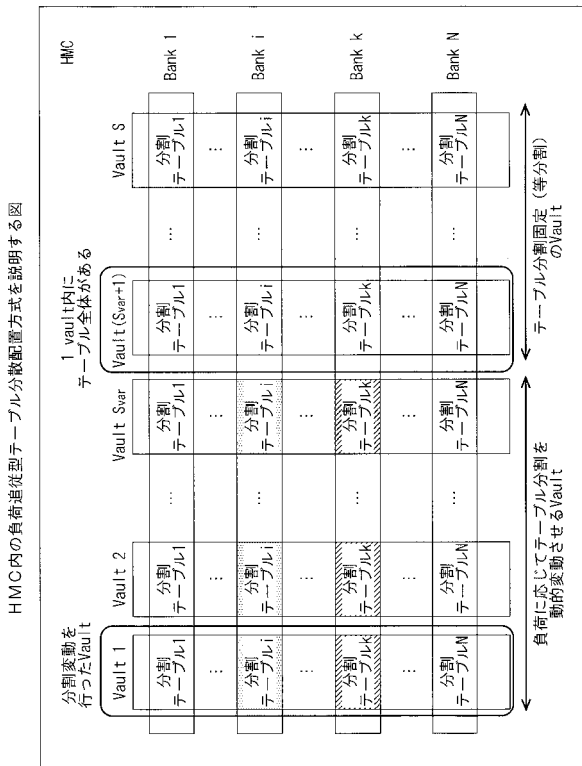
HMCコントローラ内のテーブル更新制御機構の処理フロー例



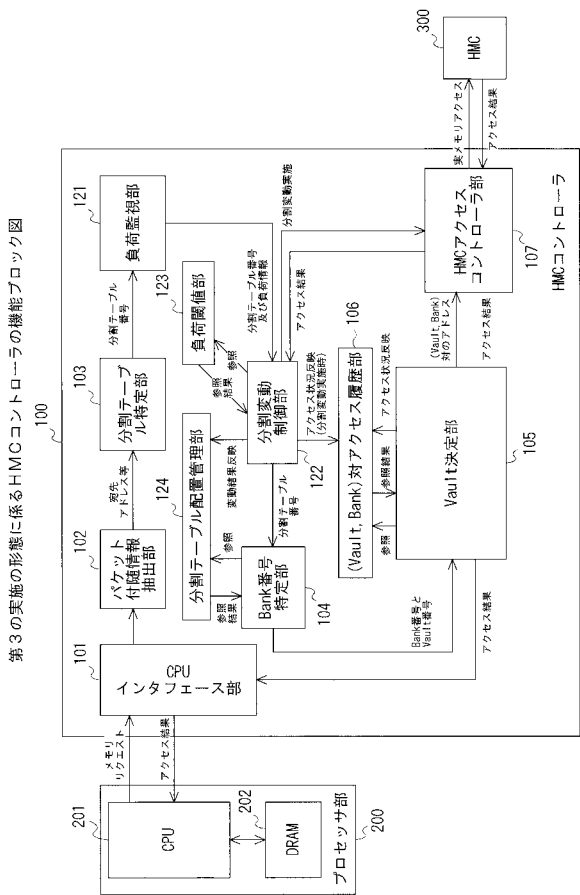
【 図 1 1 】



【 図 1 2 】



【 図 1 3 】



【 図 1 4 】

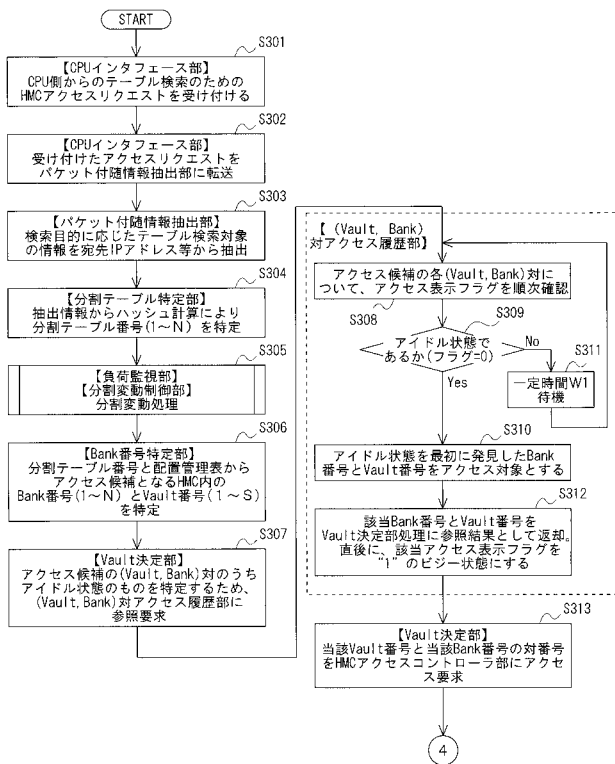
分割テーブル配置管理部がもつ配置管理表の一例

Bank	Vault	1	...	Svar	Svar+1	...	S
1		1, 0			1, 0	1, 0	1, 0
...							
i		i, 1			i, 0	i, 0	i, 0
...							
k		k, 2			k, 0	k, 0	k, 0
...							
N		N, 0			N, 0	N, 0	N, 0

(vault, bank)ごとに、(配置する分割テーブル番号、分割変動フラグ)を保持
分割変動フラグは
0: 変動未実施
1: 変動実施・負荷最大
2: 変動実施・負荷最小

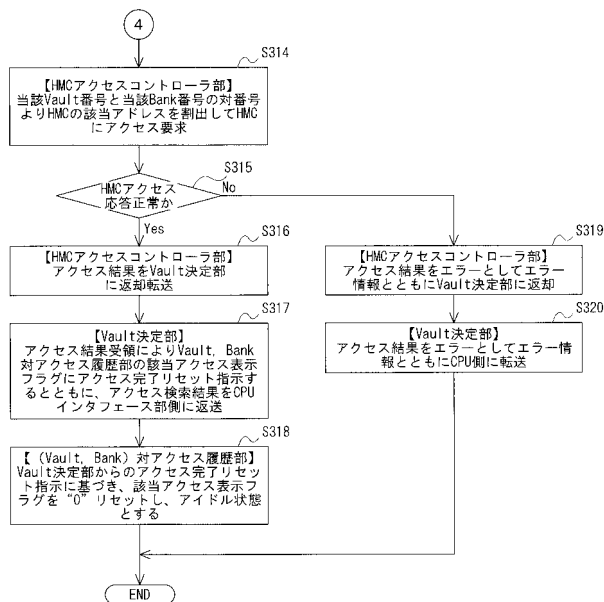
【図15】

HMCコントローラ内の振り分け機構および
負荷追従型テーブル分割変動機構の処理フロー例



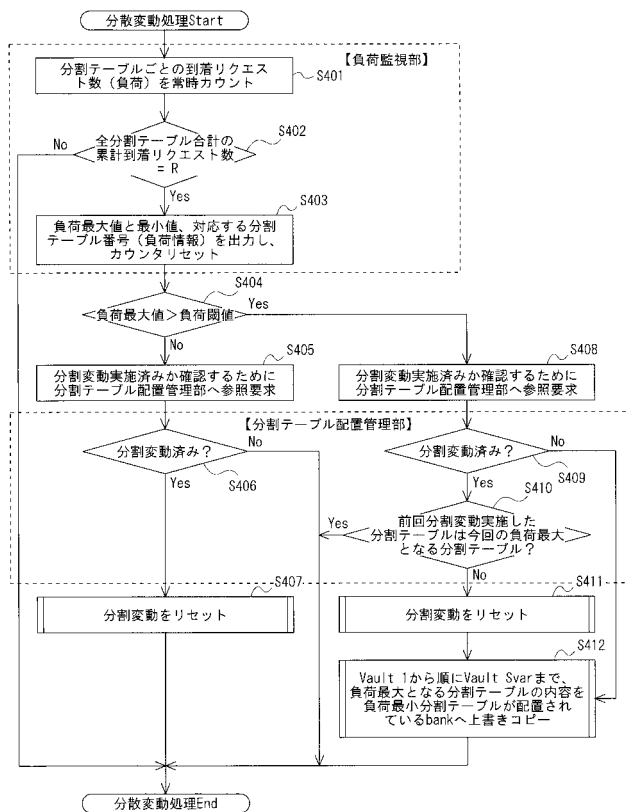
【図16】

HMCコントローラ内の振り分け機構および
負荷追従型テーブル分割変動機構の処理フロー例



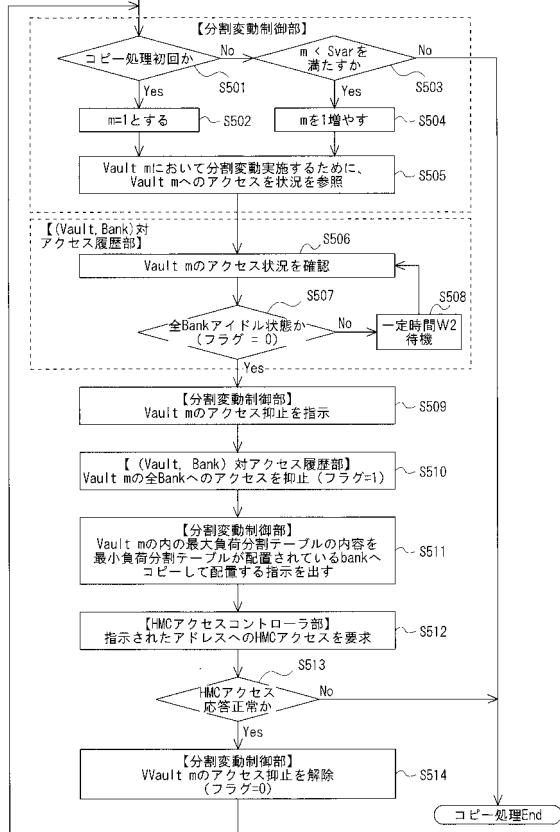
【図17】

分散変動実施時の処理フロー例

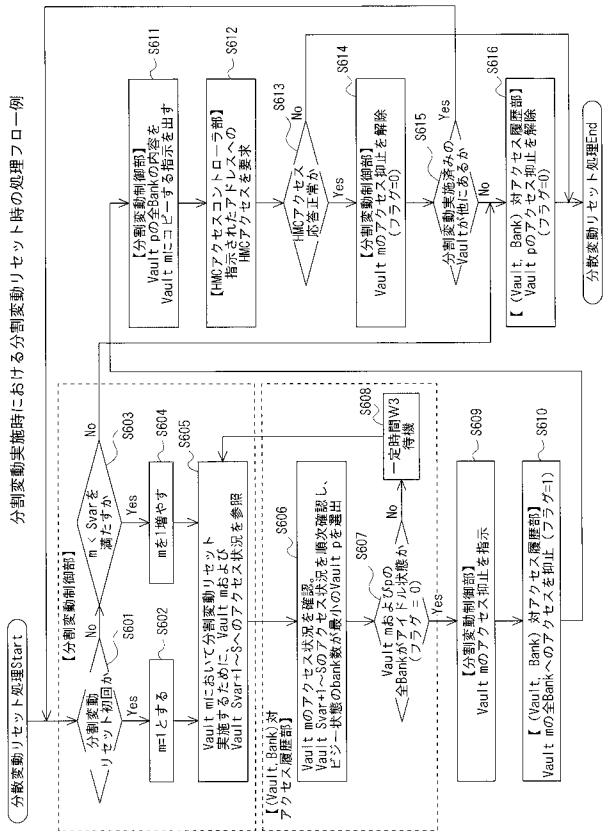


【図18】

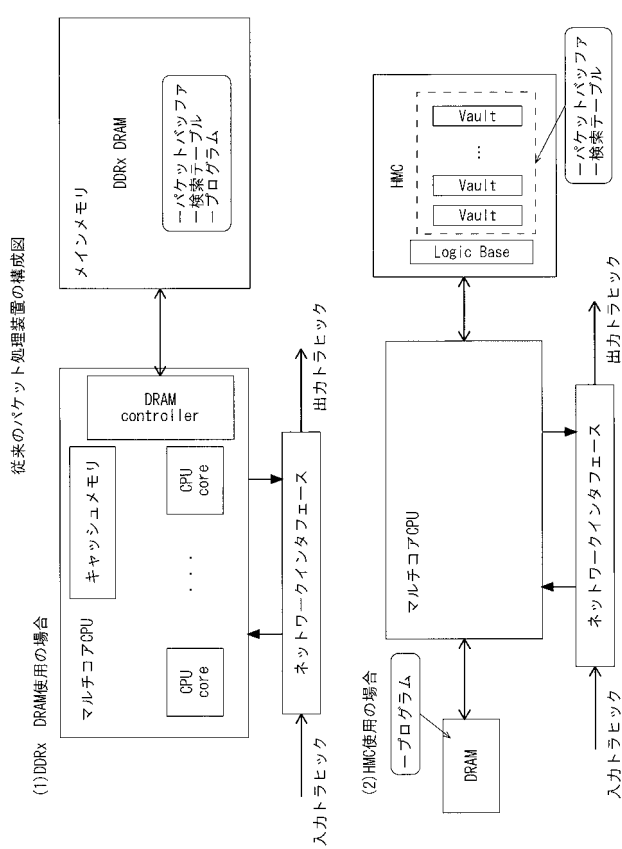
分割変動実施時における分割テーブル
コピー処理の処理フロー



【図 19】



【図 20】



フロントページの続き

(72)発明者 大木 英司

京都府京都市左京区吉田本町3番地1 国立大学法人京都大学内

(72)発明者 何 馥君

京都府京都市左京区吉田本町3番地1 国立大学法人京都大学内

Fターム(参考) 5B160 CA12

5K030 GA01 HA08 KA01 KA02 KA11