

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4574712号
(P4574712)

(45) 発行日 平成22年11月4日(2010.11.4)

(24) 登録日 平成22年8月27日(2010.8.27)

(51) Int. Cl.

F I

G06F 12/08 (2006.01)

G06F 12/08 505B

G06F 12/08 509B

請求項の数 9 (全 20 頁)

(21) 出願番号 特願2008-502585 (P2008-502585)
 (86) (22) 出願日 平成18年2月28日(2006.2.28)
 (86) 国際出願番号 PCT/JP2006/303743
 (87) 国際公開番号 W02007/099598
 (87) 国際公開日 平成19年9月7日(2007.9.7)
 審査請求日 平成20年5月2日(2008.5.2)

前置審査

(73) 特許権者 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 本藤 幹雄
 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

審査官 ▲高▼橋 正▲徳▼

最終頁に続く

(54) 【発明の名称】 演算処理装置、情報処理装置及び制御方法

(57) 【特許請求の範囲】

【請求項1】

データを保持する記憶装置に接続される演算処理装置において、
 処理を行う演算に用いるデータを要求する演算処理部と、
 第1のラインサイズの第1のキャッシュラインを複数有するとともに、前記演算処理部
 が要求したデータを、前記複数の第1のキャッシュラインのいずれにも保持しない場合、
 キャッシュミス通知を送信する第1のキャッシュメモリと、
 前記第1のラインサイズと異なる第2のラインサイズの第2のキャッシュラインを複数
 有するとともに、受信したプリフェッチ要求に対応して前記記憶装置から入力するデー
 タを保持する第2のキャッシュメモリと、
 前記キャッシュミス通知を受信した場合、前記第2のキャッシュライン毎に、前記第2
 のラインサイズのデータ量を、前記記憶装置から前記第2のキャッシュメモリにプリフェ
 ッチするプリフェッチ要求を前記キャッシュミスが発生したアドレスから前記第2のライ
 ンサイズ分先のアドレスに対して発行するプリフェッチ制御部を有することを特徴とする
 演算処理装置。

【請求項2】

前記演算処理装置において、
 前記プリフェッチ制御部は、
 前記第2のラインサイズ毎に、前記プリフェッチ要求を1回発行することを特徴とする
 請求項1記載の演算処理装置。

【請求項 3】

前記演算処理装置において、
前記プリフェッチ制御部は、

前記第 2 のキャッシュメモリにおいて、同一のアドレスから開始される前記第 2 のラインサイズのデータに対して、前記プリフェッチ要求を複数回発行することにより、前記第 2 のキャッシュメモリが有するプリフェッチ要求記憶部に前記プリフェッチ要求を複数登録することを特徴とする請求項 1 記載の演算処理装置。

【請求項 4】

前記演算処理装置において、
前記プリフェッチ制御部はさらに、

プリフェッチ要求のデータ量の設定を行うムーブインサイズ設定レジスタと、

前記ムーブイン設定レジスタに基づいて、前記第 1 のラインサイズ毎に前記プリフェッチ要求を発行するか又は前記第 2 のラインサイズ毎に前記プリフェッチ要求を発行するかを切替える切換部を有することを特徴とする請求項 1 ~ 3 のいずれか 1 項に記載の演算処理装置。

【請求項 5】

データを保持する記憶装置と、

処理を行う演算に用いるデータを要求する演算処理部と、

第 1 のラインサイズの第 1 のキャッシュラインを複数有するとともに、前記演算処理部が要求したデータを、前記複数の第 1 のキャッシュラインのいずれにも保持しない場合、

前記第 1 のラインサイズと異なる第 2 のラインサイズの第 2 のキャッシュラインを複数有するとともに、受信したプリフェッチ要求に対応して前記記憶装置から入力するデータを保持する第 2 のキャッシュメモリと、

前記キャッシュミス通知を受信した場合、前記第 2 のキャッシュライン毎に、前記第 2 のラインサイズのデータ量を、前記記憶装置から前記第 2 のキャッシュメモリにプリフェッチするプリフェッチ要求を前記キャッシュミスが発生したアドレスから前記第 2 のラインサイズ分先のアドレスに対して発行するプリフェッチ制御部を有することを特徴とする情報処理装置。

【請求項 6】

前記情報処理装置において、

前記プリフェッチ制御部は、

前記第 2 のラインサイズ毎に、前記プリフェッチ要求を 1 回発行することを特徴とする請求項 5 記載の情報処理装置。

【請求項 7】

前記情報処理装置において、

前記プリフェッチ制御部は、

前記第 2 のキャッシュメモリにおいて、同一のアドレスから開始される前記第 2 のラインサイズのデータに対して、前記プリフェッチ要求を複数回発行することにより、前記第 2 のキャッシュメモリが有するプリフェッチ要求記憶部に前記プリフェッチ要求を複数登録することを特徴とする請求項 5 記載の情報処理装置。

【請求項 8】

前記情報処理装置において、

前記プリフェッチ制御部はさらに、

プリフェッチ要求のデータ量の設定を行うムーブインサイズ設定レジスタと、

前記ムーブイン設定レジスタに基づいて、前記第 1 のラインサイズ毎に前記プリフェッチ要求を発行するか又は前記第 2 のラインサイズ毎に前記プリフェッチ要求を発行するかを切替える切換部を有することを特徴とする請求項 5 ~ 7 のいずれか 1 項に記載の情報処理装置。

【請求項 9】

データを保持する記憶装置に接続される演算処理装置の制御方法において、
前記演算処理装置が有する演算処理部が、処理を行う演算に用いるデータを要求するステップと、

第1のラインサイズの第1のキャッシュラインを複数有する第1のキャッシュメモリが、前記演算処理部が要求したデータを、前記複数の第1のキャッシュラインのいずれにも保持しない場合、キャッシュミス通知を送信するステップと、

前記演算処理装置が有するプリフェッチ制御部が、前記キャッシュミス通知を受信した場合、前記第1のラインサイズと異なる第2のラインサイズの第2のキャッシュライン毎に、前記第2のラインサイズのデータ量を、前記記憶装置から前記第2のキャッシュラインを複数有する第2のキャッシュメモリにプリフェッチするプリフェッチ要求を前記キャッシュミスが発生したアドレスから前記第2のラインサイズ分先のアドレスに対して発行するステップと、

前記第2のキャッシュメモリが、受信したプリフェッチ要求に対応して前記記憶装置から入力するデータを保持するステップとを有することを特徴とする制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、お互いにラインサイズが異なる2階層以上のキャッシュメモリを備えるプロセッサに関し、更に詳しくはキャッシュメモリに対するプリフェッチ機能を有するプロセッサについての技術に関する。

【背景技術】

【0002】

従来からHPC(High Performance Computing)など、科学技術計算等で用いられるメモリへの連続アクセスを行うコンピュータでは、キャッシュレジスタに対してプリフェッチの技術が適用されている。

【0003】

プリフェッチは、近い将来必要とされている命令もしくはデータを予め予測してキャッシュメモリ等に読み込んでおく手法で、キャッシュメモリのキャッシュミスが減らすことが出来る。

【0004】

特許文献1は、プリフェッチ機能を備えたキャッシュシステムについて開示がある。特許文献1のシステムでは、メモリデータへの連続アクセスする場合において、キャッシュミスを起因として、連続アクセスにおいて次にアクセスするラインサイズ先の予測アドレスをキューに登録しておき、実際にアクセスアドレスがキューにヒットし、予測が当たった場合に、連続アクセスであると判断し、ラインサイズ先の次にアクセスするアドレスに対して、プリフェッチを発行している。

【0005】

複数階層のキャッシュメモリにおいて、上位階層のキャッシュメモリと下位階層のキャッシュメモリとでラインサイズが異なる場合、最下位階層のキャッシュミスによって、ムーブインされるデータサイズは、上位階層のラインサイズから下位階層のラインサイズまでのいずれかのサイズである。そして、ハードウェアプリフェッチが機能する連続アクセスのケースでは、もっともデータサイズの大きい下位階層のラインサイズであるケースがもっとも性能が高くなるため、上記ケースでは、ムーブインされるデータサイズは、下位階層のキャッシュのラインサイズである可能性が高い。

【0006】

たとえば、Columbus 2メモリシステムでは、キャッシュミスによりムーブインされるデータサイズは、メモリアクセスのケースでは下位階層のキャッシュのラインサイズであるが、コピーバックのケースでは上位階層のキャッシュのラインサイズである。

【0007】

連続アクセスの多いHPC系JOBでは、コピーバック率が低いいため、上記の連続アク

10

20

30

40

50

セスのケースでは、ムーブインされるデータサイズは、下位階層のキャッシュのラインサイズである可能性が高い。

【0008】

上位階層のキャッシュメモリと下位階層のキャッシュメモリとでラインサイズが異なるキャッシュメモリシステムにおいてプリフェッチを行うと、以下の問題点が生じる。

下位階層のキャッシュミスでムーブインされるデータサイズが、下位階層のキャッシュのラインサイズである場合、上位階層のキャッシュから下位階層のキャッシュに発行されるハードウェアプリフェッチの要求（下位階層のキャッシュへのムーブイン要求）は、下位階層のキャッシュのラインサイズにつき1回でよい。しかし、従来のキャッシュシステムでは上位階層のキャッシュのラインサイズごとに発行してしまい、無駄な下位階層キャッシュアクセスパイプラインを消費することとなる。

10

【0009】

下位階層のキャッシュミスでムーブインされるデータサイズが、下位階層のキャッシュのラインサイズである場合、下位階層のキャッシュに発行されるハードウェアプリフェッチの要求は、下位階層のキャッシュのラインサイズにつき1回でよい。しかし、ハードウェアプリフェッチは、実装上の制約により、プリフェッチ要求をロストしてしまうケースが時々あり、ハードウェアプリフェッチがロストした場合には、プリフェッチ要求の発行が1回のみだと、下位階層のキャッシュへのメモリデータのムーブイン要求が発行されなくなってしまう。

【0010】

20

下位階層のキャッシュレジスタで生じたキャッシュミスでムーブインされるデータサイズが、下位階層のキャッシュメモリのラインサイズである場合、下位階層のキャッシュメモリに対して発行されるプリフェッチ要求は、下位階層のキャッシュメモリのラインサイズにつき1回でよい。よって、上位階層のキャッシュレジスタでミスしたアドレスに対して、上位階層のキャッシュメモリのラインサイズ分先のアドレスをプリフェッチ要求のプリフェッチアドレスの初期値としてしまうと、下位階層のキャッシュにとって同一ラインである可能性があるため、無駄なプリフェッチ要求で下位階層キャッシュアクセスパイプラインを消費する。

【0011】

プリフェッチが機能するメモリへの連続アクセスの場合では、下位階層のキャッシュメモリへムーブインされるデータサイズは、下位階層のキャッシュメモリのラインサイズである可能性が高いが、場合によっては下位階層のキャッシュのラインサイズとは異なるデータサイズであることもある。

30

【特許文献1】特開2004-38345号公報

【発明の開示】

【0012】

本発明の課題は、上記問題点を解決したプリフェッチ機能を有するプロセッサを提供することである。

上記課題を解決するため、本発明によるプリフェッチ機能を有するプロセッサは、第1の階層のキャッシュメモリ、第2の階層のキャッシュメモリ、及びプリフェッチ制御部を備える。

40

【0013】

第1の階層のキャッシュメモリは、第1のラインサイズを持つ。

第2の階層のキャッシュメモリは、当該第1の階層のキャッシュメモリの下位階層で、前記第1のラインサイズとは異なる大きさの第2のラインサイズを持つ。

【0014】

プリフェッチ制御部は、前記第2のラインサイズ毎に、前記第1のラインサイズ分のブロックをプリフェッチするように、前記第1の階層のキャッシュメモリから前記第2の階層のキャッシュに対するプリフェッチ要求を発行する。

【0015】

50

この構成により、不必要なプリフェッチ要求が発行されるのを防ぐことが出来る。

前記プリフェッチ制御部は、前記第2のラインサイズ毎に1回乃至複数回前記プリフェッチ要求を発行する構成とすることも出来る。

【0016】

また前記プリフェッチ制御部は、前記第1のラインサイズの2倍以上のブロックをプリフェッチするように、前記プリフェッチ要求を発行する構成とすることも出来る。

この構成により、実装上の制約によって、プリフェッチ要求がロストする場合にも対処することが出来る。

【0017】

更に前記プリフェッチ制御部は、前記プリフェッチ要求を行うプリフェッチ先のアドレスを、前記第1の階層のキャッシュメモリでミスしたアドレスから前記第2のラインサイズ分先のアドレスとする構成とすることも出来る。

10

また前記プリフェッチ制御部は、前記第1のラインサイズ毎に前記プリフェッチ要求を発行するのと、前記第2のラインサイズ毎に前記プリフェッチ要求を発行するのとを、ムーブインしたデータの大きさに基づいて切り換える切換部を更に備える構成とすることも出来る。

【0018】

この構成により、コピーバック等第2のラインの大きさ以外のムーブインにも対処することが出来る。

20

本発明によれば、第1の階層のキャッシュメモリではなく、第2の階層のキャッシュメモリのラインサイズである第2のラインサイズ毎にプリフェッチ要求が発行されるので、無駄な発行によって、第2の階層のキャッシュメモリのアクセスパイプラインが消費されるのを抑制することが出来る。

【0019】

また、実装上の制約によって、プリフェッチ要求がロストしても、第2の階層のキャッシュメモリへのメモリデータのムーブイン要求が発行される可能性が高くなるようにすることで、性能向上が図れる。

【0020】

更に、第1の階層のキャッシュメモリにミスしたアドレスに対して、第1のラインサイズ分先のアドレスではなく、第2の階層のキャッシュメモリのラインサイズ分先のアドレスをハードウェアプリフェッチのプリフェッチアドレスの初期値とすることによって、無駄な要求によって第2の階層のキャッシュアクセスパイプラインが消費されるのを抑制することが出来る。

30

【0021】

また最終的にムーブインしたデータサイズに応じて、ハードウェアプリフェッチ要求を発行することによって、第2の階層のキャッシュメモリへムーブインされるデータサイズが、第2の階層のキャッシュメモリのラインサイズと異なる場合でも、必要な要求がもれることなく、正しく要求が発行されるようになる。

【図面の簡単な説明】

40

【0022】

【図1】本実施形態におけるコンピュータシステムのプロセッサ及びその周辺構成の概略図である。

【図2】本実施形態におけるプロセッサのメモリ管理部分を中心に描いた図である。

【図3】第1の実施形態のプリフェッチキュー(PFQ)の構成例を示す図である。

【図4】加算器の出力アドレスと比較器の出力の関係を示す図である。

【図5】手順8、9、10における各状態を示した図である。

【図6】第1の実施形態のプリフェッチキュー(PFQ)の動作を示すフローチャートである。

【図7】第2の実施形態のプリフェッチキュー(PFQ)の構成例を示す図である。

50

【図 8】第 3 の実施形態のプリフェッチキュー（P F Q）の構成例を示す図である。

【図 9】第 4 の実施形態のプリフェッチキュー（P F Q）の構成例を示す図である。

【発明を実施するための最良の形態】

【0023】

以下に本発明の一実施形態を図面を参照しながら説明する。

図 1 は本実施形態におけるコンピュータシステムのプロセッサ及びその周辺構成の概略図である。

【0024】

図 1 の構成では、プロセッサユニット 1、プリフェッチ制御装置 2、1 次キャッシュ 3、2 次キャッシュ 4 及び主記憶装置 5 を有している。

プロセッサユニット 1 は、A L U、レジスタ等を含み、実際の計算やデータ処理を司る部分である。また同図の構成では、分岐予測等もプロセッサユニット 1 内で行われ、予測結果に基いたリクエストを 1 次キャッシュ 3 に行く。プリフェッチ制御装置 2 は、プリフェッチ処理の制御全般を受け持つ装置で、プロセッサユニット 1 から 1 次キャッシュ 3 へのリクエストアドレスを監視しながら、2 次キャッシュ 4 にプリフェッチを要求する。1 次キャッシュ 3 は、1 次キャッシュシステムで、アクセス速度の早いメモリと 1 次キャッシュ制御装置から構成されている。2 次キャッシュ 4 は、2 次キャッシュシステムで、主記憶装置 5 よりアクセス速度が早く 1 次キャッシュ 3 より容量の大きなメモリと 2 次キャッシュ制御装置から構成されている。また本実施形態では、プリフェッチされたデータはこの 2 次キャッシュ 4 に保持される。主記憶装置 5 は、D R A M 等によって構成されるメモリである。

【0025】

プロセッサユニット 1 が、主記憶装置 5 上のデータにアクセスする際は、要求アドレスをリクエストアドレス 6 から指定し、読み出し時にはフェッチデータ 7 を読み出し、書き込み時にはストアデータ 8 として 1 次キャッシュ 3 に出力する。

【0026】

1 次キャッシュ 3 は、プロセッサユニット 1 からの読み出し要求に対して、要求アドレスのデータを自己が保持していれば、そのデータをフェッチデータ 7 としてプロセッサユニット 1 に出力し、保持していない場合には、リクエストバス 1 1 からそのデータを含む 1 ライン分のデータを 2 次キャッシュ 4 に対して要求すると共にキャッシュミス 9 としてプリフェッチ制御装置 2 に通知する。そして、フェッチデータ 1 2 を受け取ると、プロセッサユニット 1 に要求されたデータをフェッチデータ 7 として出力する。また 1 次キャッシュ 3 は、自己が保持しているキャッシュデータが更新された場合、適当なタイミングでデータバス 1 3 からそのデータを 2 次キャッシュ 4 にライトバックする。

【0027】

2 次キャッシュ 4 は、1 次キャッシュ 3 からのデータの要求に対して、そのデータを保持していれば、そのデータを含む 1 ライン分のデータをフェッチデータ 7 として 1 次キャッシュ 3 に出力し、保持していない場合には、リクエストバス 1 4 からそのデータを含む 1 ライン分のデータを主記憶装置 5 に対して要求する。そして、フェッチデータ 1 5 を受け取ると、1 ライン分のデータを 1 次キャッシュ 3 に出力する。また 2 次キャッシュ 4 は、1 次キャッシュ 3 と同様、自己が保持しているキャッシュデータが更新されると、適当なタイミングでデータバス 1 6 からそのデータを主記憶装置 5 にライトバックする。

【0028】

プロセッサユニット 1 が 1 次キャッシュ 3 に対してデータを要求する時、アドレスバス 6 でアドレスを指定するが、このアドレス値をプリフェッチ制御装置 2 は監視し、自己が備えているプリフェッチアドレスキューをこのアドレス値によって検索する。そしてこのアドレスが、プリフェッチアドレスキューに存在するアドレスを先頭とする 1 ブロック中にある（以下、ヒットするという）場合、プリフェッチアドレスバス 1 0 から 2 次キャッシュ 4 にプリフェッチ要求アドレスを出力してプリフェッチ要求を行うと共にアドレスをプリフェッチアドレスキュー 2 5 に登録し、またプリフェッチアドレスキュー内に存在し

10

20

30

40

50

ない場合はプリフェッチを要求しない。

【0029】

なお本実施形態では、1次キャッシュ3と2次キャッシュ4はお互いに異なるラインサイズを持つキャッシュメモリであり、以下の説明では、1次キャッシュ3のラインサイズは64バイト(以下Bと記す)、2次キャッシュ4のラインサイズは256Bであるとする。

【0030】

図2は本実施形態におけるプロセッサのメモリ管理部分を中心に描いた図である。

同図において、プロセッサは、メモリ管理用の構成要素として、フェッチポート(FP)21、ストアポート(SP)22、1次キャッシュアクセスパイプライン23及び1次キャッシュムーブインバッファ(L1\$MIB)24を1次キャッシュ3内に備え、プリフェッチキュー25をプリフェッチ制御装置2内に備え、2次キャッシュムーブインポート(L2\$MIP)26、2次キャッシュプリフェッチポート(L2\$PFP)、2次キャッシュアクセスパイプライン28及び2次キャッシュムーブインバッファ(L2\$MIP)29を2次キャッシュ4内に備え、システムコントローラムーブインポート(SCMIP)30を主記憶装置5内に備えている。

【0031】

フェッチポート(FP)21は、プロセッサユニット1からのload命令やstore命令等を受け付けるポートである。またストアポート(SP)22は、ストアコミットしたstore命令が、キャッシュにデータを書き込むためのポートである。また2次キャッシュムーブインバッファ(L2\$MIP)29及びシステムコントローラムーブインポート(SCMIP)30は、それぞれ2次キャッシュ4及び主記憶装置5に対するムーブイン要求を受け付けるポートである。

【0032】

1次キャッシュアクセスパイプライン23及び2次キャッシュアクセスパイプライン28は、1次キャッシュ3及び2次キャッシュ4に対するアクセス要求を受け付けるパイプラインである。1次キャッシュアクセスパイプライン23は、P、T、M、B及びRの5つのステージを持ち、Pステージではアドレスを選択してそのアドレスを転送し、Tステージでは転送されたアドレスで1次キャッシュのタグとTLB(トランスレーションルックアップテーブル)を参照し、MステージではTステージの参照結果として得られたデータの比較(マッチング)を行ない、Bステージでは比較結果に基づいて、1次キャッシュのデータを選択して転送し、Rステージでは1次キャッシュミスやTLBミスなどに対して転送したデータが有効かあるいは無効かを示すフラグを計算して送る。2次キャッシュアクセスパイプライン28は、PR1、XP0-14のステージを持ち、各ステージでは、ポートの選択、L2\$タグ検索、アドレス比較、L2\$ミス時にL2\$MIBに登録、L2\$ヒット時にL2\$データの読み出し、L2\$データのL1\$MIBへの転送などを行っている。

【0033】

1次キャッシュムーブインバッファ(L1\$MIB)24及び2次キャッシュムーブインバッファ(L2\$MIB)29は、1次キャッシュ3及び2次キャッシュ4に対して生じたムーブイン命令をバッファリングするものである。

【0034】

プリフェッチキュー(PFQ)25は、以前プリフェッチを行ったアドレスの1ライン分先のアドレスを登録しており、1次キャッシュ3でキャッシュミスが生じると、キャッシュミスが生じたアドレスとプリフェッチキュー(PFQ)25内に登録されているアドレスをマッチングし、プリフェッチキュー(PFQ)25に一致するアドレスが登録されていれば、2次キャッシュプリフェッチポート(L2\$PFP)27に、プリフェッチ要求を発行する。2次キャッシュプリフェッチポート(L2\$PFP)27は、プリフェッチキュー(PFQ)25からのプリフェッチ要求を受け付けるものである。

【0035】

同図における動作を以下に説明する。

10

20

30

40

50

プロセッサユニット1で、load命令等をデコードし、メモリの読み出し要求が発行されると、この要求は、フェッチポート(FP)25から1次キャッシュアクセスパイプライン23に入力される。読み出し要求に対して1次キャッシュ2がヒットすればそのままデータをフェッチポート(FP)25から要求を発行したプロセッサユニット1に返して、データをレジスタ31に書きこむ。

【0036】

1次キャッシュ2がミスしたときは、2次キャッシュ3からデータを持ってこなければいけないので、1次キャッシュムーブインバッファ(L1\$MIB)24に要求を入れる。1次キャッシュムーブインバッファ(L1\$MIB)24は、2次キャッシュ3に対して読み出し要求を出す。これは2次キャッシュ3のリクエストを受け取る2次キャッシュムーブインポート(L2\$MIP)26を介して2次キャッシュアクセスパイプライン28に入る。

10

【0037】

そしてこの読み出し要求が2次キャッシュ3でヒットすれば、そのデータを1次キャッシュムーブインバッファ(L1\$MIB)24に入れ、1次キャッシュムーブインバッファ(L1\$MIB)24は1次キャッシュラインアクセスパイプラインを獲得して1次キャッシュ2にデータを書きこむ(1次キャッシュミス2次キャッシュヒットの場合)。

【0038】

次にハードウェアプリフェッチを行う場合について説明する。

20

1次キャッシュ2でミスして、プリフェッチキュー(PFQ)25にハードウェアプリフェッチとして動作すべきアドレスが登録されていない場合、そのアドレスを一旦プリフェッチキュー(PFQ)27に登録する。このとき特許文献1に示してあるように、64B分先のアドレスを登録する、次に64バイト先のアクセスしに行ったときには1次キャッシュがミスすると同時にプリフェッチキュー(PFQ)25はヒットする。このときプリフェッチキュー(PFQ)25は、更に64Bを足して+128Bのアドレスのプリフェッチのリクエストをプリフェッチポート(L2\$PFP)27に出す。

【0039】

1次キャッシュミスは、2次キャッシュムーブインポート(L2\$MIP)26と2次キャッシュプリフェッチポート(L2\$PFP)27に登録され、2次キャッシュにアクセスしてヒットすればデータを返す。また、ミスすれば2次キャッシュムーブインバッファ(L2\$MIB)29に登録してシステムコントロールムーブインポート(SCMIP)30に出力して主記憶装置5にリクエストを出す。そして主記憶装置5データが帰ってきたら、それを2次キャッシュアクセスパイプライン28を介して2次キャッシュ3に書きこみ、同時にバイパスで1次キャッシュアクセスパイプライン23に返しこれを1次キャッシュ2に書き込む。

30

【0040】

図3は、第1の実施形態のプリフェッチキュー(PFQ)25の構成例を示す図である。

同図のプリフェッチ3、選択回路44、加算器45、選択回路46及び47、及び加算器48を備え、各エントリ41-1~41-nはそれぞれ、エントリ41に登録されるアドレス値等がセットされるレジスタ49、レジスタ49内のアドレスとリクエストアドレスを比較する比較器50及び比較器50の比較結果と後述するレジスタ49内の有効ビットとのANDを求めるAND回路51を有している。

40

【0041】

レジスタ49は、アドレス値の他に状態フラグとして働く有効ビット、待機ビット及びL2\$PFP登録許可フラグを記録している。

レジスタ49内の有効ビットは、レジスタ48にセットされているアドレス値が有効かどうかを示すもので、アドレス値がレジスタ48に登録される時セットされ、このエントリ41からアドレス値が読み出された時にリセットされる。待機ビットは、有効ビットが

50

セットされているエントリ 4 1 において、1 次キャッシュアクセスパイプライン 2 3 からのリクエストアドレスがレジスタ 4 8 内に登録されているアドレス値にマッチした場合セットされる。プリフェッチアドレスキュー (P F Q) 2 5 は、この待機ビットの状態から読み出しを行うエントリ 4 1 - 1 ~ 4 1 - n を決定する。L 2 \$ P F P 登録許可フラグは、1 次キャッシュアクセスパイプライン 2 3 からのリクエストアドレスと、このエントリ 4 1 に登録されているアドレスがマッチ (ヒット) したときに次の 2 5 6 B の連続アドレスを 2 次キャッシュプリフェッチポート (L 2 \$ P F P) 2 7 に登録するかどうかの判断に用いられるもので、L 2 \$ P F P 登録許可フラグに 1 がセットされていれば 2 次キャッシュプリフェッチポート (L 2 \$ P F P) に登録を行ない、0 がセットされていれば 2 次キャッシュプリフェッチポート (L 2 \$ P F P) に登録を行わない。

10

【 0 0 4 2 】

1 次キャッシュアクセスパイプライン 2 3 からリクエストアドレスが入力され、これが新規登録される場合、有効ビットには 1、待機ビットには 0、L 2 \$ P F P 登録許可フラグには 1、及びアドレスにはリクエストアドレスに加算器 4 5 で 6 4 を加えた値が入力される。

【 0 0 4 3 】

またレジスタ 4 9 の登録アドレスが更新される場合には、L 2 \$ P F P 登録許可フラグには、比較器 4 2 による加算器 4 5 の出力のビット [7 : 6] が 0 かどうかの比較結果と新規登録かどうかの O R を O R 回路 4 3 で求めた結果が入力される。この L 2 \$ P F P 登録許可フラグの内容は、選択回路 4 6 及び 4 7 の選択信号となっており、L 2 \$ P F P 登録許可フラグが 1 のとき P F P リクエスト信号として 1 が出力され、また P F P リクエストアドレスとしてそのエントリ 4 1 のレジスタ 4 9 に登録されているレジスタ値に加算器 4 8 によって 2 5 6 B 加算した値が出力される。そしてこれらの出力によって、2 次キャッシュ P F P (L 2 \$ P F P) 2 7 に P F P リクエストアドレスが登録される。

20

【 0 0 4 4 】

図 4 は、加算器 4 5 の出力アドレスと比較器 4 2 の出力の関係を示す図である。

加算器 4 5 によって、レジスタ 4 9 にセットされているアドレス値を 6 4 B インクリメントしてレジスタ 4 9 に登録すると、その出力アドレスのビット [7 : 6] は、4 回に一回 0 となり、よって比較器 4 2 からは 4 回に 1 回 1 が出力され、これが O R 回路 4 3 を介して L 2 \$ P F P 登録許可フラグにセットされる。なおレジスタ 4 9 への登録が新規登録の場合、L 2 \$ P F P 登録許可フラグには 1 がセットされる。よって、アドレスの新規登録からアドレスが 4 回更新されるごとに 1 回 L 2 \$ P F P 登録許可フラグに 1 がセットされ、2 次キャッシュ P F P (L 2 \$ P F P) 2 7 に P F P リクエストアドレスが登録される。

30

【 0 0 4 5 】

この L 2 \$ P F P 登録許可フラグは、プリフェッチキュー (P F Q) 2 5 に新規登録を行うとき及び 2 5 6 B 境界の先頭 6 4 B アドレスをプリフェッチキュー (P F Q) 2 5 に更新登録するときにセットされる。また 2 5 6 B 境界の先頭 6 4 B 以外のアドレスを P F Q に更新登録するときリセットされる。

【 0 0 4 6 】

このプリフェッチキュー (P F Q) 2 5 は、1 次キャッシュアクセスパイプライン 2 3 からプリフェッチのリクエストアドレスが入力されると、このアドレス値はレジスタ 4 9 内のアドレス値と比較器 5 0 によって比較され、この比較結果と有効ビットとの A N D を A N D 回路 5 1 で取り、結果を P F Q にヒットしたかどうかを示す P F Q ヒット信号として 1 次キャッシュアクセスパイプライン 2 3 に出力する。したがって、リクエストアドレスとレジスタ 4 9 が一致し、且つ有効ビットが 1 のとき P F Q ヒット信号が 1 になる。

40

【 0 0 4 7 】

また 1 次キャッシュ 2 のラインの大きさが 6 4 B に対して 2 次キャッシュ 3 のラインの大きさが 2 5 6 B と、上位と下位のキャッシュメモリのラインサイズが異なるときであっても、2 次キャッシュプリフェッチポート (L 2 \$ P F P) 2 7 へのアドレス値の登録は

50

4 回に 1 回 (2 5 6 B / 6 4 B) にすることが出来、プリフェッチ要求は 2 次キャッシュ 3 のラインサイズに付き 1 回となる。よって無駄な下位階層キャッシュアクセスパイプラインの消費を抑止し、性能向上を図ることが出来る。

【 0 0 4 8 】

次に、図 2、図 3 を用い、プリフェッチ動作を含む、プロセッサのメモリアクセス命令に対する処理の詳細手順を以下に示す。

以下の説明では、アドレス A、A + 8、A + 1 6、. . .、A + 5 6 に対する l o a d 命令がプロセッサユニットでデコードされた場合を例として示す。

1 : l o a d 命令が、フェッチポート (F P) 2 1 を介して 1 次キャッシュアクセスパイプライン 2 3 を獲得。

2 : 1 次キャッシュアクセスパイプライン 2 3 において、1 次キャッシュに対してアドレス A でアクセス。

3 : 1 次キャッシュアクセスパイプライン 2 3 において、2 の結果、1 次キャッシュミスを検出。

4 : 1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 にミスアドレスを登録。

【 0 0 4 9 】

4 . 1 : 1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 は、2 次キャッシュ 3 から 1 次キャッシュ 2 へのムーブイン要求を、2 次キャッシュムーブインポート (L 2 \$ M I P) 2 6 に対して発行。

【 0 0 5 0 】

4 . 2 : 2 次キャッシュムーブインポート (L 2 \$ M I P) 2 6 は、2 次キャッシュアクセスパイプライン 2 8 を獲得し、2 次キャッシュに対してアドレス A でアクセス。

4 . 3 : 2 次キャッシュアクセスパイプライン 2 8 において、手順 4 . 2 の結果、2 次キャッシュミスを検出。

【 0 0 5 1 】

4 . 4 : 2 次キャッシュムーブインバッファ (L 2 \$ M I B) 2 9 にミスアドレスを登録。

4 . 5 : 2 次キャッシュムーブインバッファ (L 2 \$ M I B) 2 9 は、主記憶装置 5 から 2 次キャッシュ 3 へのムーブイン要求を、システムコントローラムーブインポート (S C M I P) 3 0 に対して発行。

【 0 0 5 2 】

4 . 6 : システムコントローラムーブインポート (S C M I P) 3 0 は、主記憶装置 5 からミスアドレス A から 2 5 6 B 分のデータを取り出し、2 次キャッシュムーブインバッファ (L 2 \$ M I B) 2 9 にムーブイン。

【 0 0 5 3 】

4 . 7 : 2 次キャッシュムーブインバッファ (L 2 \$ M I B) 2 9 は、2 次キャッシュアクセスパイプライン 2 8 を獲得し、2 次キャッシュ 4 に 2 5 6 B のムーブインデータを書き込む。

【 0 0 5 4 】

4 . 8 : 2 次キャッシュムーブインバッファ (L 2 \$ M I B) 2 9 は、1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 に 6 4 B のムーブインデータをバイパス転送。

4 . 9 : アドレス A で 1 次キャッシュミスした l o a d 命令は、1 次キャッシュアクセスパイプライン 2 3 を獲得し、1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 に転送されたムーブインデータをバイパスして読み出し、プロセッサユニット 1 内のレジスタ 3 1 にデータを書き込む。

【 0 0 5 5 】

4 . 1 0 : 1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 は、1 次キャッシュアクセスパイプライン 2 3 を獲得し、1 次キャッシュ 2 に 6 4 B のムーブインデータを書き込む。

5 : プリフェッチキュー (P F Q) 2 5 がミスを検出。

10

20

30

40

50

6：次の連続アドレス（ $A + 64$ ）をプリフェッチキュー（ PFQ ）25に新規登録。レジスタ49内の $L2\$PFP$ 登録許可フラグをセット。

7：同様に、連続アドレス（ $A + 8$ 、 $A + 16$ 、...、 $A + 56$ ）にアクセスする $load$ 命令が、1次キャッシュアクセスパイプライン23を獲得。

8：このとき2次キャッシュ3からのムーブインデータが到着していなければ、1次キャッシュ MIB ヒット、データミスを検出し、1次キャッシュアクセスパイプライン23をアポート。アポートされた要求は、フェッチポート（ FP ）21に戻る。

9：また2次キャッシュ3からのムーブインデータが到着していて、1次キャッシュ2にデータが書き込まれていなければ、1次キャッシュ MIB ヒット、データヒットを検出し、1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24のデータをバイパスして読み出し、レジスタ31にデータを書き込む。

10：また2次キャッシュ3からのムーブインデータが到着していて、1次キャッシュにデータが書き込まれていなければ、1次キャッシュヒットを検出し、1次キャッシュ2からデータを読み出し、レジスタ31にデータを書き込む。

【0056】

図5は、この手順8、9、10における各状態を示したものである。

手順8の状態では、1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24からアドレスは到着しているがデータは到着しておらず、また1次キャッシュ2にはデータは書き込まれていないので、1次キャッシュアクセスパイプライン23をアポートする。

【0057】

また手順9の状態では、1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24からアドレス及びデータが到着しているが、1次キャッシュ2にはデータは書き込まれていないので、1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24からデータを読み出してレジスタ31にデータを書き込む。

【0058】

また手順10の状態では、1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24からアドレス及びデータが到着しており、また1次キャッシュ2にはデータが書き込まれているので、1次キャッシュ2からデータを読み出してレジスタ31にデータを書き込む。

【0059】

以下に続けて連続アドレス（ $A + 64$ ）にアクセスする $load$ 命令についての処理を説明する。

11：手順1と同様に、連続アドレス（ $A + 64$ ）にアクセスする $load$ 命令が、1次キャッシュアクセスパイプライン23を獲得。

12：手順11の結果、1次キャッシュミスを検出。

【0060】

12.1：1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24にミスアドレスを登録し、2次キャッシュ3にアクセス。

12.2：2次キャッシュヒットを検出し、2次キャッシュから64Bのデータを読み出し、1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24に転送。

【0061】

12.3：アドレス（ $A + 64$ ）で、1次キャッシュミスした $load$ 命令は、1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24のデータをバイパスして読み出し、レジスタ31にデータを書き込む。

12.4：1次キャッシュムーブインバッファ（ $L1\$MIB$ ）24は、1次キャッシュ2に64Bデータを書き込む。

13：プリフェッチキュー（ PFQ ）25のヒットを検出。レジスタ49内の待機ビットをセット。

14：次の連続アドレス（ $A + 128$ ）をプリフェッチキュー（ PFQ ）25に登録。レジスタ49内の $L2\$PFP$ 登録許可フラグをリセット。

15：手順14でリセットされるまで $L2\$PFP$ 登録許可フラグがセットされていたの

10

20

30

40

50

で、次の256B連続アドレス(A+64+256)を2次キャッシュプリフェッチポート(L2\$PFP)27に登録。

【0062】

15.1: 2次キャッシュプリフェッチポート(L2\$PFP)27は、2次キャッシュアクセスパイプライン28を獲得し、2次キャッシュ3に対して、アドレス(A+64+256)でアクセス。

【0063】

15.2: 手順15.1の結果、2次キャッシュミスを検出。

15.3: 2次キャッシュムーブインバッファ(L2\$MIB)29にキャッシュミスアドレスに登録。

10

【0064】

15.4: 2次キャッシュムーブインバッファ(L2\$MIB)29は、主記憶装置5から2次キャッシュ3へのムーブイン要求を、システムコントローラムーブインポート(SCMIP)30に発行。

【0065】

15.5: システムコントローラムーブインポート(SCMIP)30は、主記憶装置5からミスアドレス(A+64+256)から256B分のデータを取り出し、それを2次キャッシュムーブインバッファ(L2\$MIB)29にムーブイン。

【0066】

15.6: 2次キャッシュムーブインバッファ(L2\$MIB)29は、2次キャッシュアクセスパイプライン28を獲得し、2次キャッシュ3に256Bのムーブインデータを書き込む。

20

16: 同様に、連続アドレス(A+64+8、A+64+16、...、A+64+56)にアクセスするload命令が、1次キャッシュアクセスパイプライン23を獲得。

17: 2次キャッシュ4からのムーブインデータが到着していなければ、1次キャッシュMIBヒット、データミスを検出し1次キャッシュアクセスパイプライン23をアポート。アポートした要求は、フェッチポート(FP)21に戻る。

18: 2次キャッシュ3からのムーブインデータが到着していて、1次キャッシュ2にデータが書き込まれていなければ、1次キャッシュMIBヒット、データヒットを検出し、1次キャッシュムーブインバッファ(L1\$MIB)24のデータをバイパスして読み出し、それをレジスタ31に書き込む。

30

19: 2次キャッシュ3からのムーブインデータが到着していて、1次キャッシュ2にデータが書き込まれていければ、1次キャッシュヒットを検出し、1次キャッシュ2からデータを読み出し、それをレジスタ41に書き込む。

【0067】

以下に続けて連続アドレス(A+128)にアクセスするload命令についての処理を説明する。

20: 手順1、11と同様に、連続アドレス(A+128)にアクセスするload命令が、1次キャッシュアクセスパイプライン23を獲得。

21: 手順20の結果、1次キャッシュミスを検出。

40

【0068】

21.1: 1次キャッシュムーブインバッファ(L1\$MIB)24にキャッシュミスアドレスに登録し、2次キャッシュ3にアクセス。

21.2: 手順21.1の結果、2次キャッシュヒットを検出し、2次キャッシュ3から64Bのデータを読み出し、1次キャッシュムーブインバッファ(L1\$MIB)24に転送。

【0069】

21.3: アドレス(A+128)で1次キャッシュミスしたload命令は、1次キャッシュムーブインバッファ(L1\$MIB)24からデータをバイパスして読み出し、これをレジスタ31に書き込む。

50

【 0 0 7 0 】

2 1 . 4 : 1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 は、1 次キャッシュ 2 に 6 4 B データを書き込む。

2 2 : プリフェッチキュー (P F Q) 2 5 のヒットを検出。レジスタ 4 9 の待機ビットをセット。

2 3 : 次の連続アドレス (A + 1 9 2) をプリフェッチキュー (P F Q) 2 5 に登録。L 2 \$ P F P 登録許可フラグをリセット。

2 4 : (手順 2 3 でリセットされるまで L 2 \$ P F P 登録許可フラグがセットされていたので、次の 2 5 6 B 連続アドレス (A + 1 2 8 + 2 5 6) は、2 次キャッシュプリフェッチポート (L 2 \$ P F P) 2 7 に登録しない。)

2 5 : 同様に、連続アドレス (A + 1 2 8 + 8、A + 1 2 8 + 1 6、. . .、A + 1 2 8 + 5 6) にアクセスする l o a d 命令が、1 次キャッシュアクセスパイプライン 2 3 を獲得。

2 6 : 2 次キャッシュ 3 からのムーブインデータが到着していなければ、1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4、データミスを検出し 1 次キャッシュアクセスパイプライン 2 3 をアボート。アボートした要求は、フェッチポート (F P) 2 1 に戻る。

2 7 : 2 次キャッシュ 3 からのムーブインデータが到着していて、1 次キャッシュにデータが書き込まれていなければ、1 次キャッシュ M I B ヒット、データヒットを検出し、1 次キャッシュ M I B のデータをバイパスして読み出し、レジスタ 3 1 にデータを書き込む。

2 8 : 2 次キャッシュ 3 からのムーブインデータが到着していて、1 次キャッシュ 2 にデータが書き込まれていければ、1 次キャッシュヒットを検出し、1 次キャッシュ 2 からデータを読み出し、レジスタ 3 1 にデータを書き込む。

【 0 0 7 1 】

以下に続けて連続アドレス (A + 1 9 2) にアクセスする l o a d 命令についての処理を説明する

2 9 : 手順 1、1 1、2 0 と同様に、連続アドレス (A + 1 9 2) にアクセスする l o a d 命令が、1 次キャッシュアクセスパイプラインを獲得。

3 0 : 手順 2 9 の結果、1 次キャッシュミスを検出。

【 0 0 7 2 】

3 0 . 1 : 1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 にキャッシュミスアドレスを登録し、2 次キャッシュ 3 にアクセス。

3 0 . 2 : 2 次キャッシュヒットを検出し、2 次キャッシュから 6 4 B のデータを読み出し、1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 に転送。

【 0 0 7 3 】

3 0 . 3 : アドレス (A + 1 9 2) で、1 次キャッシュミスした l o a d 命令は、1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 からデータをバイパスして読み出し、これをレジスタ 3 1 を書き込む。

【 0 0 7 4 】

3 0 . 4 : 1 次キャッシュムーブインバッファ (L 1 \$ M I B) 2 4 は、1 次キャッシュ 2 に 6 4 B データを書き込む。

3 1 : プリフェッチキュー (P F Q) 2 5 のヒットを検出。レジスタ 4 9 の待機ビットをセット。

3 2 : 次の連続アドレス (A + 2 5 6) をプリフェッチキュー (P F Q) 2 5 に登録。レジスタ 4 9 の L 2 \$ P F P 登録許可フラグをセット。

3 3 : (手順 3 2 でセットされるまで L 2 \$ P F P 登録許可フラグがリセットされていたので、次の 2 5 6 B 連続アドレス (A + 1 9 2 + 2 5 6) は、2 次キャッシュプリフェッチポート (L 2 \$ P F P) 2 7 に登録しない。)

3 4 : 同様に、連続アドレス (A + 1 9 2 + 8、A + 1 9 2 + 1 6、. . .、A + 1 9 2

10

20

30

40

50

+ 5 6) にアクセスする `load` 命令が、1 次キャッシュアクセスパイプライン 2 3 を獲得。

3 5 : 2 次キャッシュ 3 からのムーブインデータが到着していなければ、1 次キャッシュムーブインバッファ (`L 1 $ M I B`) 2 4 ヒット、データミスを検出し 1 次キャッシュアクセスパイプライン 2 3 をアボート。アボートされた要求はフェッチポート (`F P`) 2 1 に戻る。

3 6 : 2 次キャッシュ 3 からのムーブインデータが到着していて、1 次キャッシュ 2 にデータが書き込まれていなければ、1 次キャッシュムーブインバッファ (`L 1 $ M I B`) 2 4 ヒット、データヒットを検出し、1 次キャッシュムーブインバッファ (`L 1 $ M I B`) 2 4 からデータをバイパスして読み出し、これをレジスタ 3 1 に書き込む。

3 7 : 2 次キャッシュ 3 からのムーブインデータが到着していて、1 次キャッシュ 2 にデータが書き込まれていれば、1 次キャッシュヒットを検出し、1 次キャッシュ 2 からデータを読み出し、これをレジスタ 3 1 に書き込む。

【 0 0 7 5 】

以下に続けて連続アドレス (`A + 2 5 6`) にアクセスする `load` 命令についての処理を説明する。

3 8 : 同様に、連続アドレス (`A + 2 5 6`) にアクセスする `load` 命令が、1 次キャッシュアクセスパイプライン 2 3 を獲得。

3 9 : 手順 3 8 の結果 1 次キャッシュ 2 ミスを検出。

【 0 0 7 6 】

3 9 . 1 : 1 次キャッシュムーブインバッファ (`L 1 $ M I B`) 2 4 にミスアドレスを登録し、2 次キャッシュ 3 にアクセス。

3 9 . 2 : 2 次キャッシュヒットを検出し、2 次キャッシュ 3 から 6 4 B のデータを読み出し、これを 1 次キャッシュムーブインバッファ (`L 1 $ M I B`) 2 4 に転送。

【 0 0 7 7 】

3 9 . 3 : アドレス (`A + 2 5 6`) で、1 次キャッシュミスした `load` 命令は、1 次キャッシュムーブインバッファ (`L 1 $ M I B`) 2 4 からデータをバイパスして読み出し、これをレジスタ 3 1 に書き込む。

【 0 0 7 8 】

3 9 . 4 : 1 次キャッシュ `M I B` は、1 次キャッシュに 6 4 B データを書き込む。

4 0 : プリフェッチキュー (`P F Q`) 2 5 ヒットを検出。レジスタ 4 9 内の待機ビットをセット。

4 1 : 次の連続アドレス (`A + 3 2 0`) をプリフェッチキュー (`P F Q`) 2 5 に登録。レジスタ 4 9 内の `L 2 $ P F P` 登録許可フラグをリセット。

4 2 : 手順 4 1 でリセットされるまで `L 2 $ P F P` 登録許可フラグがセットされていたので、次の 2 5 6 B 連続アドレス (`A + 2 5 6 + 2 5 6`) を 2 次キャッシュプリフェッチポート (`L 2 $ P F P`) 2 7 に登録。

【 0 0 7 9 】

4 2 . 1 : 2 次キャッシュプリフェッチポート (`L 2 $ P F P`) 2 7 は、2 次キャッシュアクセスパイプライン 2 8 を獲得し、2 次キャッシュ 3 に対してアドレス (`A + 2 5 6 + 2 5 6`) でアクセス。

【 0 0 8 0 】

4 2 . 2 : 手順 4 2 . 1 の結果、2 次キャッシュミスを検出。

4 2 . 3 : 1 次キャッシュムーブインバッファ (`L 1 $ M I B`) 2 4 にミスアドレスを登録。

【 0 0 8 1 】

4 2 . 4 : 1 次キャッシュムーブインバッファ (`L 1 $ M I B`) 2 4 は、主記憶装置 5 から 2 次キャッシュ 3 へのムーブイン要求を、システムコントローラムーブインポート (`S C M I P`) 3 0 に対して発行。

【 0 0 8 2 】

10

20

30

40

50

42.5: システムコントローラムーブインポート (SCMI P) 30 は、主記憶装置 5 からキャッシュミスアドレス ($A + 256 + 256$) から 256 B 分のデータを取り出し、このデータを 2 次キャッシュムーブインバッファ (L2 \$ MIB) 29 にムーブイン

【0083】

42.6: 2 次キャッシュムーブインバッファ (L2 \$ MIB) 29 は、2 次キャッシュアクセスパイプライン 28 を獲得し、2 次キャッシュ 3 に 256 B のムーブインデータを書き込む。

【0084】

以下、連続アドレス ($A + 320$)、($A + 384$)、・・・にアクセスする load 命令について同様の処理を繰り返す。

図 6 は、図 3 に示した第 1 の実施形態のプリフェッチキュー (PFQ) 25 の動作を示すフローチャートである。

【0085】

ステップ S1 において、アドレス A にて 1 次キャッシュ 2 にアクセスし、1 次キャッシュ 2 がキャッシュミスし (ステップ S2、Y)、且つプリフェッチキュー (PFQ) 25 もミスしたら (ステップ S3、Y)、ステップ S4 としてプリフェッチキュー (PFQ) 25 に 1 次キャッシュ 2 の 1 ライン分先のアドレス ($A + 64$) を登録し、またプリフェッチキュー (PFQ) 25 内のレジスタ 49 の L2 \$ 登録許可フラグをセットして、処理をステップ S1 に戻す。

【0086】

またステップ S2 において、1 次キャッシュ 2 がヒットし (ステップ S2、N)、プリフェッチキューがミスしたとき (ステップ S3、Y)、処理をステップ S1 に戻す。

ステップ S2 において、1 次キャッシュ 2 がミスし (ステップ S2、Y)、プリフェッチキュー (PFQ) 25 はヒットしたとき (ステップ S3、N)、及び 1 次キャッシュ 2 がヒットし (ステップ S2、N)、プリフェッチキュー (PFQ) 25 はヒットしたとき (ステップ S3、N)、処理をステップ S6 に移し、プリフェッチを行う。

【0087】

ステップ S6 では、プリフェッチキュー (PFQ) 25 のレジスタ 49 内の待機ビットをセットする。そしてレジスタ 49 内の L2 \$ 登録許可フラグがセットされていたら (ステップ S7、Y)、2 次キャッシュプリフェッチポート (L2 \$ PFP) 27 に PFP リクエストアドレス ($A + 64$) でプリフェッチ要求を登録し、L2 \$ 登録許可フラグがセットされていなければ (ステップ S7、N)、2 次キャッシュプリフェッチポート (L2 \$ PFP) 27 にリクエストを登録しない。

【0088】

次にステップ S9 としてプリフェッチキュー (PFQ) 25 の登録アドレスを $A + 64$ に更新した後、($A / 64 + 1$) を 4 で割った余りが 0 ならばステップ S11 としてレジスタ 49 内の L2 \$ 登録許可フラグをセットして、処理をステップ S1 に戻す。また $A / 64 + 1$ を 4 で割った余りが 0 でないのならばステップ S12 として L2 \$ 登録許可フラグをリセットして、処理をステップ S1 に戻す。

【0089】

このように第 1 の実施形態では、下位層のキャッシュレジスタのラインの大きさ毎にプリフェッチ要求を発行することが出来るので、無駄なプリフェッチ要求がアクセスパイプラインを占めることが無く、性能向上を図ることが出来る。

【0090】

次に、プリフェッチキュー (PFQ) 25 の第 2 の構成例について説明する。

第 1 の実施形態のプリフェッチキュー (PFQ) 25 が、下位階層のキャッシュの 1 ラインの大きさが上位階層のキャッシュの 1 ラインの大きさの n 倍であったとき、n 回の連続アクセスに対して 1 回 2 次キャッシュプリフェッチポート (L2 \$ PFP) 27 に対してプリフェッチ要求を登録していたが、第 2 の実施形態のプリフェッチキュー (PFQ)

10

20

30

40

50

25は、n回に2回以上プリフェッチ要求を登録する。

【0091】

下位階層のキャッシュのミスでムーブインされるデータサイズが、下位階層のキャッシュのラインサイズである場合、下位階層のキャッシュに発行されるハードウェアプリフェッチの要求は、第1の実施形態のプリフェッチキュー(PFQ)25のように下位階層のキャッシュのラインサイズにつき1回でよい。

【0092】

しかし、プリフェッチ要求は、ハードウェア実装上の制約によりロストしてしまうケースが時々あるため、プリフェッチ要求の発行が1回のみだと、ハードウェアプリフェッチがロストした場合に、下位階層のキャッシュへのメモリデータのムーブイン要求が発行されなくなってしまう。この実装上の制約とは、例えば2次キャッシュ3がミスすると2次キャッシュムーブインバッファ(L2\$MIB)29に登録されるが、2次キャッシュムーブインバッファ(L2\$MIB)29が一杯のときは、再登録を行わずにプリフェッチ要求がロストしてしまうことがある。

【0093】

この点に対処し、第2の実施形態のプリフェッチキュー(PFQ)25は、下位階層のキャッシュのラインサイズにつき複数回プリフェッチ要求を2次キャッシュプリフェッチポート(L2\$PFP)27に発行する。

【0094】

図7は、第2の実施形態のプリフェッチキュー(PFQ)25の構成例を示す図である。なお同図は、図3の第1の実施形態のプリフェッチキュー(PFQ)25と対比する形で記載されている。

【0095】

図7を図3の第1の実施形態のプリフェッチキュー(PFQ)25と比較すると、比較器42aの入力が、加算器45aから出力されるアドレスのうちビット[6]のみになっている。よって第1の実施形態ではレジスタ49のアドレスが4回更新されると1回L2\$PFP登録許可フラグが1にセットされているが、第2の実施形態のプリフェッチキュー(PFQ)25では、アドレスが4回更新されると2回L2\$PFP登録許可フラグに1がセットされ、2次キャッシュプリフェッチポート(L2\$PFP)27にリクエストが登録される。

【0096】

これにより第2の実施形態では、1つのプリフェッチ要求が、ハードウェアの実装上の問題でロストしても、対処することが出来る。

次に、プリフェッチキュー(PFQ)25の第3の実施形態について説明する。

【0097】

第3の実施形態のプリフェッチキュー(PFQ)25も、第2の実施形態と同様、ハードウェアの実装上の問題で、プリフェッチ要求がロストする場合に対処したものである。

第3の実施形態では、上位階層のキャッシュのラインサイズの2倍以上のブロックをプリフェッチするように2次キャッシュプリフェッチポート(L2\$PFP)27にプリフェッチ要求を発行する。これによりプリフェッチ要求は、2次キャッシュプリフェッチポート(L2\$PFP)27で2倍以上に展開され、複数のプリフェッチ要求が発行される。

【0098】

図8は、第3の実施形態のプリフェッチキュー(PFQ)25の構成例を示す図である。なお同図も、図3の第1の実施形態のプリフェッチキュー(PFQ)25と対比する形で記載されている。

【0099】

図8の第3の実施形態のプリフェッチキュー(PFQ)25の構成を図3の第1の実施形態の構成と比較すると、2次キャッシュプリフェッチポート(L2\$PFP)27に出力されるPFPリクエストブロックサイズ61が128Bと、1次キャッシュ2のライン

10

20

30

40

50

の2倍になっている。なお不図示であるが、図3の第1の実施形態ではこのPFPリクエストブロックサイズは、1次キャッシュ2のラインサイズと同じ64Bとなっている。

【0100】

この構成により第3の実施形態のプリフェッチキュー(PFQ)25では、1次キャッシュ2のラインサイズの2倍のサイズのブロックサイズを指定して、プリフェッチ要求を2次キャッシュプリフェッチポート(L2\$PFP)27に対して発行するので、2次キャッシュプリフェッチポート(L2\$PFP)27では2回のプリフェッチ要求が発行されることになる。

【0101】

これにより第3の実施形態でも、1つのプリフェッチ要求が、ハードウェアの実装上の問題でロストしても、対処することが出来る。

次に、プリフェッチキュー(PFQ)25の第4の実施形態について説明する。

【0102】

第4の実施形態のプリフェッチキュー(PFQ)25は、本実施形態で行われている下位層キャッシュのラインサイズ(256B)毎のムーブインと、従来のプロセッサで行われている上位層キャッシュのラインサイズ(64B)毎のムーブインを切り換えることが出来るようにしたものである。

【0103】

これにより、コピーバックの際に実行される上位層キャッシュのラインサイズ(64B)毎のムーブインにも対処することが出来る。

図9は、第4の実施形態のプリフェッチキュー(PFQ)25の構成例を示す図である。なお同図も、図3の第1の実施形態のプリフェッチキュー(PFQ)25と対比する形で記載されている。

【0104】

同図の第4の実施形態のプリフェッチキュー(PFQ)25と図3の第1の実施形態を比較すると、図9の構成では、レジスタ49bにムーブイン(MI)データサイズが記憶されている。このMIデータサイズは、0がセットされるとプリフェッチキュー(PFQ)25は256B毎のムーブインを行ない、1がセットされると64B毎のムーブインを行う。

【0105】

このMIデータサイズには、初期値として0がセットされ、キャッシュミスしたムーブインアドレスとレジスタ49bにセットされているアドレスを比較器71で比較した結果、両者が一致し、且つムーブインデータサイズが64Bであったとき、AND回路72の出力によって1がセットされる。またこのMIデータサイズは、OR回路73によってL2\$PFP登録許可フラグとのORが取られた結果が選択回路46bに入力される。よって、MIデータサイズに0がセットされているときは256Bのムーブインを行ない、1がセットされているときは64Bのムーブインを行なう、というようにムーブインの大きさを切り換えることができる。

【0106】

以上のように本実施形態によれば、上位層のキャッシュレジスタと下位層のキャッシュレジスタのラインの大きさが異なっても、下位層のキャッシュレジスタのラインの大きさ毎にプリフェッチ要求を発行することが出来るので、無駄なプリフェッチ要求でキャッシュアクセスパイプラインが消費されるのを抑止し、性能向上を図ることが出来る。

【0107】

また実装上の制約により、プリフェッチ要求がロストしてしまう点にも対処することが出来る。

更には、コピーバックの際に実行される上位層キャッシュのラインサイズのムーブインにも対処することが出来る。

【0108】

なお上記例では、本発明を2次キャッシュメモリと主記憶装置間のプリフェッチに適用

10

20

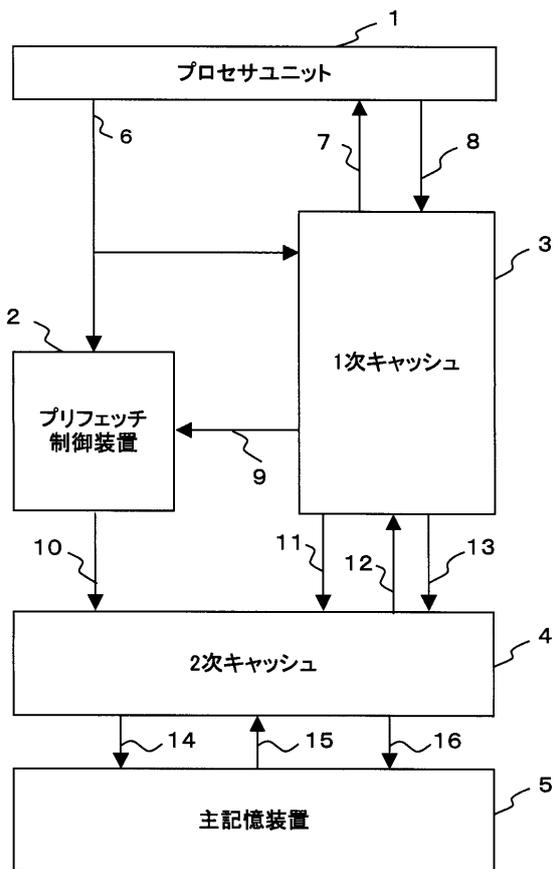
30

40

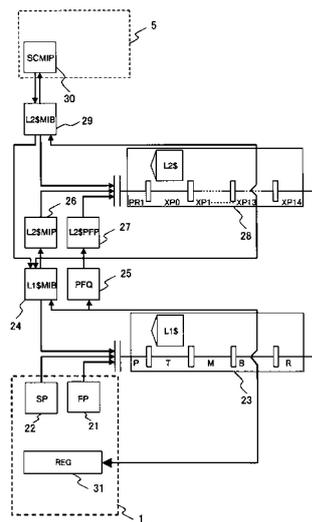
50

した場合を例として示したが、本発明はこれに限定されるものではなく、システムが3次キャッシュ以上のキャッシュメモリを備えている場合、2次キャッシュと3次キャッシュの間、3次キャッシュと主記憶装置の間等にも適用することが出来る。
また、上記例では、プリフェッチの連続アクセス方向が、昇順であるケースについて適用した場合を例として示したが、本発明はこれに限定されるものではなく、プリフェッチの連続アクセス方向が、降順であるケースについても、適用することができる。

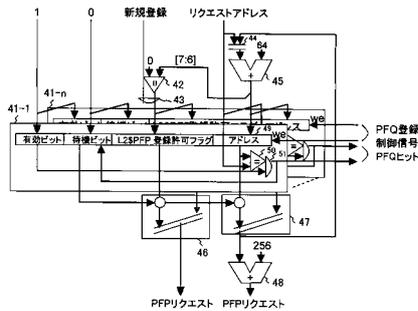
【図1】



【図2】



【図3】



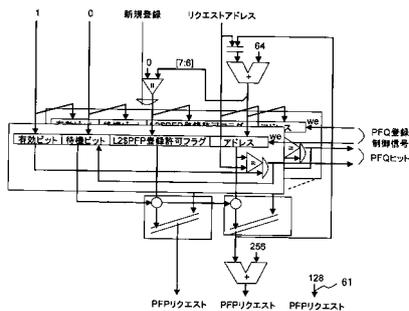
【図4】

比較器	42	7	6	5	4	3	2	1	0	
	0	0	1	0	0	0	0	0	0	+64
	0	1	0	0	0	0	0	0	0	+128
	0	1	1	0	0	0	0	0	0	+192
	1	0	0	0	0	0	0	0	0	+256

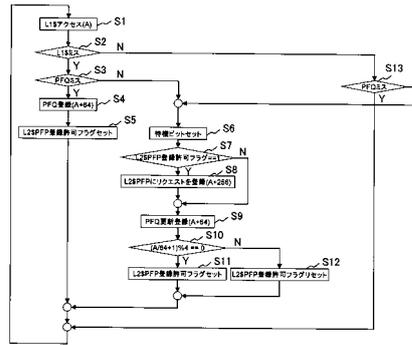
【図5】

	L1\$ MIB Address	L1\$ MIB Data	L1\$ data	処理
手順 8	○	×	×	ア bort
手順 9	○	○	×	L1\$MIBから
手順 10	○	○	○	L1\$ 直接

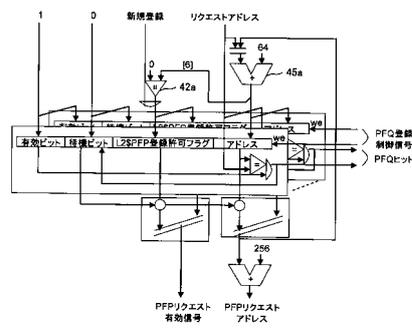
【図8】



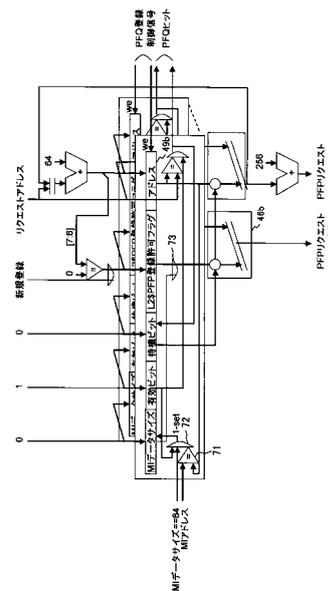
【図6】



【図7】



【図9】



フロントページの続き

- (56)参考文献 特開平02-301843(JP,A)
特開平09-128293(JP,A)
特開2004-038345(JP,A)
特開2006-040141(JP,A)
特開昭61-169949(JP,A)
特開平06-149669(JP,A)
特開平07-129464(JP,A)
特開2000-339157(JP,A)
特開2006-048181(JP,A)
特開2006-048182(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/08-12/12