

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5059609号  
(P5059609)

(45) 発行日 平成24年10月24日(2012.10.24)

(24) 登録日 平成24年8月10日(2012.8.10)

(51) Int.Cl.		F I			
<b>G06F 12/08</b>	<b>(2006.01)</b>	G06F 12/08	505C		
<b>G06F 9/38</b>	<b>(2006.01)</b>	G06F 9/38	330B		
		G06F 9/38	330G		

請求項の数 13 (全 55 頁)

(21) 出願番号	特願2007-527950 (P2007-527950)	(73) 特許権者	501261300
(86) (22) 出願日	平成17年8月16日 (2005.8.16)		エヌヴィディア コーポレイション
(65) 公表番号	特表2008-510258 (P2008-510258A)		アメリカ合衆国, カリフォルニア 950
(43) 公表日	平成20年4月3日 (2008.4.3)		50, サンタ クララ, サン トーマス
(86) 国際出願番号	PCT/US2005/029135		エクスプレスウェイ 2701
(87) 国際公開番号	W02006/038991	(74) 代理人	100094318
(87) 国際公開日	平成18年4月13日 (2006.4.13)		弁理士 山田 行一
審査請求日	平成20年7月31日 (2008.7.31)	(74) 代理人	100123995
(31) 優先権主張番号	10/920, 610		弁理士 野田 雅一
(32) 優先日	平成16年8月17日 (2004.8.17)	(72) 発明者	ハクラ, ジヤド, エス.
(33) 優先権主張国	米国 (US)		アメリカ合衆国, カリフォルニア州,
(31) 優先権主張番号	10/920, 682		マウンテン ヴュー, ナンバー401,
(32) 優先日	平成16年8月17日 (2004.8.17)		ウェスト エル キャミノ リアル 2
(33) 優先権主張国	米国 (US)		400

最終頁に続く

(54) 【発明の名称】 メモリへの様々なタイプのアクセスを予測するため、およびキャッシュメモリに関連付けられた予測を管理するための、システム、装置、および方法

(57) 【特許請求の範囲】

【請求項1】

メモリへのアクセスを予測するためのプリフェッチャであって、  
ターゲットアドレスのグループをトリガアドレスに関連付けるように構成され、および

前記ターゲットアドレスのグループのうちの少なくとも1つのターゲットアドレスに基づき、予測アドレスのグループを発行するように構成された、

第1のアドレス予測器を備え、

前記ターゲットアドレスのグループのうちの少なくとも1つの前記ターゲットアドレスは、前記トリガアドレスに対して非順次であり、

前記第1のアドレス予測器は、前記トリガアドレスが検出された場合に、非順次予測として予測アドレスのグループを発行する非順次予測器を更に備え、

前記非順次予測器は、

前記ターゲットアドレスのグループと前記トリガアドレスとの関連を格納すると共に、前記ターゲットアドレスのグループ間の相対的優先順位を格納するためのリポジトリと

要求されたアドレスのストリーム内の前記トリガアドレスを検出するように構成され、さらに、検出された前記トリガアドレスとの関連と、最高優先順位であるその優先順位とに基づき、前記少なくとも1つのターゲットアドレスを、非順次予測として選択するよ

うに構成された、非順次予測エンジンと、をさらに備え、

当該プリフェッチャは、

前記トリガアドレスが、要求されたアドレスの順次ストリーム内にある場合、および前記非順次予測が、前記トリガアドレスとして、前記順次ストリーム内の次のアドレスの指定と比較して早く生成される場合、前記順次ストリームの第1のアドレスを前記トリガアドレスとして指定するための促進器と、

少なくとも1つの予測アドレスの生成を抑制するように構成された抑制器と、  
をさらに備える、プリフェッチャ。

【請求項2】

前記予測アドレスのグループは、順次アドレスのグループの基本アドレスとしての前記少なくとも1つのターゲットアドレスと共に順次アドレスのグループを含む、請求項1に記載のプリフェッチャ。

10

【請求項3】

前記最高優先順位が、プロセッサが前記少なくとも1つのターゲットアドレスを、前記ターゲットアドレスのグループのうちの他のアドレスに対して最も新しく要求したことを少なくとも示す、請求項1に記載のプリフェッチャ。

【請求項4】

前記第1のアドレス予測器がトリガアドレスからインデックスおよびタグを生成するように構成され、前記リポジットリガトリガアドレスとターゲットアドレスとの関連を格納するためのメモリロケーションをそれぞれ有するいくつかのwayを含み、第1のwayに格納されたトリガアドレスとターゲットアドレスとの関連が、第2のwayに格納された他のトリガアドレスとターゲットアドレスとの関連よりも高位の優先順位に関連付けられる、請求項1に記載のプリフェッチャ。

20

【請求項5】

第2のアドレスとマッチするターゲットアドレスを含む前記トリガアドレスとターゲットアドレスとの関連のうちの一つに関する優先順位を修正するように構成された優先順位調整器をさらに備える、請求項4に記載のプリフェッチャ。

【請求項6】

前記抑制器が、前記トリガアドレスが、データに関する要求またはプリフェッチ要求のいずれか、あるいはその両方に関する場合、前記ターゲットアドレスのグループに関するアドレスの数を削減するように構成され、それによって前記少なくとも1つの予測アドレスの生成が抑制される、請求項1に記載のプリフェッチャ。

30

【請求項7】

さらに前記抑制器が、前記トリガアドレスの検出から、前記非順次予測としての前記予測アドレスのグループの生成までの時間間隔が、予め設定された、時間に関するしきい値未満である場合、非順次予測としての前記予測アドレスのグループの生成を抑制するように構成される、請求項1に記載のプリフェッチャ。

【請求項8】

さらに前記抑制器が、複数のインタリーブされた、要求されたアドレスの順次ストリームの検出に基づいて、少なくとも前記予測アドレスの生成を抑制するように構成される、請求項1に記載のプリフェッチャ。

40

【請求項9】

前記複数のインタリーブされた順次ストリームのそれぞれが、スレッドに関連付けられている、請求項8に記載のプリフェッチャ。

【請求項10】

少なくとも1つの他のトリガアドレスに基づいて、順次予測を生成するための順次予測器を含む、第2のアドレス予測器をさらに備える、請求項1に記載のプリフェッチャ。

【請求項11】

50

前記順次予測が、  
 前記少なくとも1つの他のトリガアドレスから昇順に並べられた第1の数のアドレス、  
 または、  
 前記少なくとも1つの他のトリガアドレスから降順に並べられた第2の数のアドレス、  
 のいずれか、あるいは、  
 前記第1および第2の数のアドレスの両方  
 を含み、  
 さらに前記抑制器が、  
 前記少なくとも1つの他のアドレスが昇順の第1のアドレスストリームの一部である  
 ことを検出し、前記降順に並べられた前記第2の数のアドレスに基づく、前記数の追加の  
 予測アドレスを抑制するように、および  
 前記少なくとも1つの他のアドレスが降順の第2のアドレスストリームの一部である  
 ことを検出し、前記昇順に並べられた前記第1の数のアドレスに基づく、前記数の追加の  
 予測アドレスを抑制するように、  
 構成された、請求項10に記載のプリフェッチャ。

10

## 【請求項12】

前記順次予測が、  
 前記少なくとも1つの他のトリガアドレスから1つだけ降順のバックアドレス、または  
 前記少なくとも1つの他のトリガアドレスのバックセクタアドレス、のいずれか、ある  
 いはその両方、を含み、  
 さらに前記抑制器が、  
 前記順次予測が前記バックアドレス又は前記バックセクタアドレスのいずれかを含む場合  
 、前記順次予測の数を減少させるように構成された、請求項10に記載のプリフェッチャ  
 。

20

## 【請求項13】

予測を維持するようにそれぞれ構成された、複数のキューを備える、予測インベントリ  
 と、  
 フィルタリング済みのアドレスのサブセットを生成するためのインベントリフィルタで  
 あって、  
 前記予測インベントリ、または  
 前記予測アドレスのグループおよび前記順次予測の、いずれかにおいて、  
 冗長アドレスを除去するように構成された、前記インベントリフィルタと、  
 をさらに備え、  
 前記プリフェッチャが、前記フィルタリング済みアドレスのサブセットのうちの少なく  
 とも1つを提供するように構成された、請求項10に記載のプリフェッチャ。

30

## 【発明の詳細な説明】

## 【発明の簡単な説明】

## 【0001】

本発明は、一般にコンピューティングシステムに関し、より具体的には、たとえば、構  
 成可能な量の予測を生成すること、ならびに、たとえば予測インベントリおよび/または  
 マルチレベルキャッシュに対して、予測を抑制およびフィルタリングすることによる、メ  
 モリへの順次および非順次アクセスの予測に関する。

40

## 【発明の背景】

## 【0002】

プログラム命令およびプログラムデータをフェッチするためにプリフェッチャが使用さ  
 れるため、プロセッサは取り出された情報それ自体を必要に応じて容易に利用すること  
 ができる。プリフェッチャは、将来プロセッサがどの命令およびデータを使用するかを予測  
 するため、プロセッサは、典型的にはプロセッサよりも低速で動作するシステムメモリか  
 ら、命令またはデータにアクセスするのを待機する必要がない。プロセッサとシステムメ  
 モリとの間にプリフェッチャが実装されると、プロセッサがメモリからの要求データを待

50

ちながらアイドル状態を続ける可能性が低くなる。したがって、プリフェッチャは一般に、プロセッサの性能を向上させる。

【0003】

一般に、プリフェッチャによってより多くの予測が生成されるほど、プロセッサが利用できる必要な命令およびデータを有するようにプリフェッチャが調整できる可能性が高くなり、それによってプロセッサの待ち時間が減少する。しかし、従来のプリフェッチャは、典型的には予測プロセスを十分に管理することができない。このように管理されないと、これらのプリフェッチャは、予測されたアドレスの量がプリフェッチャの処理能力を超えた場合にメモリリソースに負荷をかけすぎる傾向がある。このようなりソースの過負荷を防止するために、従来のプリフェッチャは、プリフェッチャまたはメモリリソースのい  
10  
ずれかを過負荷にする可能性のある量の予測を生成しないように、予測の生成を控えめにする傾向がある。加えて、従来のプリフェッチャは、通常、こうした予測プロセスを実施するコストを考慮せずに予測を生成するため、予測プロセスおよびこれをサポートするために必要なリソースの量を合理化することの特典が実現できない。特に、従来タイプのプリフェッチャは、主に、本来は逐次的である予測を生成するための標準技術に依拠しており、計算に拠るまたは拠らないにかかわらず、リソースを節約するように予測を格納しない。また、従来のプリフェッチャは、通常、予測プロセスを十分に管理していないため、予測アドレスの量がプリフェッチャの処理能力を超える場合、計算およびメモリリソースに負荷をかけすぎる傾向がある。そこで、リソースの過負荷を防止するために、これらの  
20  
プリフェッチャは、プリフェッチャを過負荷にする可能性のある量の予測を生成しないように、予測の生成を控えめにする傾向がある。さらに、多くの従来のプリフェッチャは、予測が生成された後、およびプロセッサがそれらの予測を要求する前に、予測を管理するための機能が欠如している。通常、これらのプリフェッチャはプリフェッチデータを単一のキャッシュメモリに格納するが、このキャッシュメモリは、典型的にはすでにキャッシュに格納された予測に対して過剰な予測を制限するための機能が欠如している。従来のプリフェッチャで使用されるキャッシュメモリは、単にデータを格納するためのものであり、そこに格納された予測されたアドレスを効果的に管理するように十分な設計がなされていない。

【0004】

前述の内容に鑑み、メモリへのアクセスを効果的に予測するためのシステム、装置、および方法を提供することが望ましい。理論的に言えば、例示的なシステム、装置、または方法は、少なくとも前述の欠点を最小にするかまたは除去することになる。

【発明の概要】

【0005】

メモリへのアクセスを予測するためのシステム、装置、および方法が開示される。一実施形態では、例示的装置は、プログラム命令を実行しプログラムデータを処理するように構成されたプロセッサと、プログラム命令およびプログラムデータを含むメモリと、メモリプロセッサとを備える。メモリプロセッサは、プログラム命令またはプログラムデータを含むアドレスを受け取るように構成された、スペキュレータを含むことができる。こ  
40  
うしたスペキュレータは、構成可能な数の順次アドレスを生成するための順次予測器を備えることができる。スペキュレータは、アドレスのサブセットをアドレスに関連付けるように構成された、非順次予測器を含むこともできる。非順次予測器は、サブセットのうちの少なくとも1つのアドレスに基づいて、アドレスのグループを予測するように構成することも可能であり、サブセットのうちの少なくとも1つのアドレスは、アドレスに対してパターン化できない。一実施形態では、例示的な非順次予測器が、メモリへのアクセスを予想する。非順次予測器は、アドレスからインデックスおよびタグを生成するように構成された、予測生成器を含む。また、非順次予測器は、予測生成器に結合されたターゲットキャッシュも含む。ターゲットキャッシュは、それぞれがトリガ-ターゲット関連を格納するためのメモリロケーションを有する、メモリのいくつかの部分を含む。メモリの第1の部分に格納されたトリガ-ターゲット関連は、メモリの第2の部分に格納された他のトリ  
50

ガ - ターゲット関連よりも高い優先順位に関連付けられる。

【 0 0 0 6 】

本発明の一実施形態では、この装置は、それぞれが項目のグループを維持するように構成されたキューを含む、予測インベントリを含む。項目のグループは、典型的には、項目のグループに対応するトリガアドレスを含む。グループの各項目は、予測の1つのタイプである。またこの装置は、予測の数を、予測の数と同じ予測タイプを有するキューのうちの少なくとも1つと比較するように構成された、インベントリフィルタも含む。いくつかのケースでは、インベントリフィルタは、予測の数を、異なる予測タイプを有するキューのうちの他の少なくとも1つと比較するように構成される。たとえば、いくつかの前方順次予測は、バックキュー、または同様のものに対してフィルタリング可能である。少なくとも1つの実施形態では、装置が、メモリへの予測アクセスを管理するための戻りデータキャッシュメモリを含む。戻りデータキャッシュメモリは、たとえばしきい値より短い経過時間を有する予測を格納するように構成された短期キャッシュメモリと、たとえばしきい値より長いまたは等しい経過時間を有する予測を格納するように構成された長期キャッシュメモリとを、含むことができる。長期キャッシュメモリは、通常、短期キャッシュよりも多くのメモリ容量を有する。プリフェッチャは、複数の予測が短期キャッシュメモリまたは長期キャッシュメモリ、あるいはその両方の、いずれに格納されるかにかかわらず、1サイクルの操作中、または2サイクルにわたってなど、並行して検出するように構成されたインターフェースを含むことも可能であり、インターフェースは、短期キャッシュメモリおよび長期キャッシュメモリを検査する場合、複数の予測それぞれのうちの少なくとも2つの表現を使用する。

10

20

【 0 0 0 7 】

本発明は、添付の図面と共に説明される以下の詳細な記述に関連して、より完全に理解されよう。

【 0 0 0 8 】

同じ参照番号は、図面のいくつかの図全体を通じて対応する部分を表す。

【 例示的实施形態の詳細な説明 】

【 0 0 0 9 】

本発明は、プロセッサが必要とすると予測される可能性のあるプログラム命令およびプログラムデータを取り出すために、メモリへのアクセスを効果的に予測するためのシステム、装置、および方法を提供する。メモリへのアクセスを効果的に予測することにより、1つまたは複数のプロセッサに必要なデータを提供する待ち時間を最小限にすることができる。本発明の特定の実施形態によれば、装置は、メモリアccessを予測するように構成されたスペキュレータを含む。例示的スペキュレータは、予測生成レートを変化させるために構成可能な量の予測を生成するように構成することができる。他の実施形態では、スペキュレータは、そうでなければプリフェッチャが管理しなければならなくなる可能性のある、冗長予測などの不必要な予測の量を制限するために、一定の予測の生成を抑制することができる。特定の実施形態では、スペキュレータは、予測を含むキャッシュメモリまたはインベントリが、プロセッサに提示するためのより好適な予測を含むかどうかを探查することによって、不必要な予測をフィルタリングすることもできる。一実施形態では、キャッシュメモリは、短期キャッシュおよび長期キャッシュ内に予測を格納し、その両方が冗長予測をフィルタリング除去するために並行して検査される。

30

40

【 0 0 1 0 】

順次および非順次予測を生成するためのプリフェッチャおよびスペキュレータに関する例示的実施形態

図1は、本発明の特定の実施形態に従った、例示的スペキュレータを示すブロック図である。この例でスペキュレータ108は、プリフェッチャ106内に常駐するように示される。さらにプリフェッチャ106は、1つまたは複数のプロセッサによるメモリアccessを少なくとも制御するように設計された、メモリプロセッサ104内に常駐するように示される。プリフェッチャ106は、メモリ112からのプログラム命令およびプログラ

50

ムデータの両方を、要求される前に「フェッチ」し、その後フェッチされたプログラム命令およびプログラムデータを、プロセッサ102による要求に応じてプロセッサ102に提供するように動作する。使用に先立ってフェッチすること(すなわち「プリフェッチすること」)により、プロセッサのアイドル時間(たとえば、プロセッサ102がデータ不足である間の時間)が最小化される。プリフェッチャ106は、プリフェッチされたデータのプロセッサ102への提示を格納および管理するための、キャッシュメモリ110も含む。キャッシュメモリ110は、命令の実行およびデータ取り出しをスピードアップするためのデータストアとして働く。特に、キャッシュメモリ110はプリフェッチャ106内に常駐し、一般に、メモリコントローラ104とは別に何らかの待ち時間を減少させるために採用される、「L1」および「L2」キャッシュなどの他のメモリキャッシュを補足するように動作する。

10

#### 【0011】

動作時に、スペキュレータ108は、メモリ112にアクセスするためのプロセッサ102による要求(「読み取り要求」)について、システムバス103を監視する。特に、プロセッサ102がプログラム命令を実行する場合、スペキュレータ108は、プロセッサ102によってまだ使用されていないプログラム命令およびプログラムデータを含むアドレスに関する読み取り要求を検出する。考察のために、「アドレス」は、一般にメモリ112とキャッシュメモリ110との間で転送されるメモリのキャッシュラインまたは単位に関連付けられる。キャッシュラインの「アドレス」はメモリロケーションを表すことが可能であり、キャッシュラインはメモリ112の複数のアドレスからのデータを含むことができる。「データ」という用語は、プリフェッチ可能な情報の単位を表すのに対して、「プログラム命令」および「プログラムデータ」という用語は、それぞれ、プロセッサ102によってその処理中に使用される命令およびデータを表す。したがって、データ(たとえば任意のビット数)は、プログラム命令および/またはプログラムデータを構成する予測情報を表すことができる。また「予測」という用語は、「予測アドレス」という用語と同じ意味で使用することもできる。予測アドレスがメモリ112へのアクセスに使用される場合、典型的には、その予測アドレスならびに他の(予測されるかまたはされない)アドレスを含む1つまたは複数のキャッシュラインがフェッチされる。

20

#### 【0012】

検出された読み取り要求に基づいて、スペキュレータ108は、プロセッサ102によって次に要求される可能性のある、構成可能な数の予測アドレスを生成することができる。スペキュレータ108は、本発明の少なくとも1つの実施形態に従って、1つまたは複数の推測技法を使用することによってこれを実行する。スペキュレータ108は、これらの推測技法を予測器として実施するが、その実施については以下で説明する。さらにスペキュレータ108は、いくつかの予測の生成を抑制し、他の予測をフィルタリングする。一定の予測を抑制またはフィルタリングすることにより、あるいはそれらの両方を実行することにより、冗長予測の数が減少し、それによってリソースが保存される。保存されるリソースの例には、キャッシュメモリ110などのメモリリソース、およびメモリバス111などのバスリソース(たとえば帯域幅に関して)が含まれる。

30

#### 【0013】

スペキュレータ108の予測がオプションでフィルタリングされた後、メモリプロセッサ104は残った(すなわちフィルタリング除去されなかった)予測を、メモリバス111を介してメモリ112へ移送する。これに回答して、メモリ112はプリフェッチされたデータを予測アドレスと共に戻す。キャッシュメモリ110は、戻されたデータを、メモリプロセッサ104がそのデータをプロセッサ102に送るまでなど、一時的に格納する。メモリプロセッサ104は、適切な時点で、とりわけ待ち時間が最小になることを保証するために、プリフェッチされたデータをプロセッサ102へとシステムバス103を介して移送する。

40

#### 【0014】

図2は、本発明の一実施形態に従った、例示的スペキュレータを示す図である。スペキ

50

ュレータ108は、予測203の生成元である読み取り要求201を受け取るように構成される。図に示されるように、スペキュレータ108は、制御情報およびアドレス情報を順次予測器(「SEQ.予測器」)206および非順次予測器(「NONSEQ.予測器」)216へと提供するように構成された、予測コントローラ202を含み、それらはどちらも予測203を生成する。予測コントローラ202は、全体としてまたは部分的に、最適な量およびタイプの予測を提供するように予測生成プロセスを管理する働きをする。たとえば、予測コントローラ202は、読み取り要求201で指定された特定のキャッシュラインまたはキャッシュラインのグループに対して、生成される予測の数およびタイプを変えることができる。他の例として、予測コントローラ202は、ターゲットキャッシュ218内の使用可能なメモリなどのリソースを保存するように、または予測アドレスの重複によるメモリ112への不必要なアクセスを最小にするように、一定の予測の生成を抑制するための抑制器204を含む。予測コントローラ202は、非順次予測の生成を早めるために、オプションで促進器205を含むことができる。促進器208は、図8に示されるように、非順次予測が関係する非線形アドレスストリームの直前のアドレスの検出に先立って、非順次予測の生成をトリガするように動作する。予測コントローラ202については、順次予測器206および非順次予測器216の以下の説明の後で、より詳細に論じる。

10

#### 【0015】

順次予測器206は、ある程度の見込みを有する予測(すなわち予測アドレス)を生成するように構成される。すなわち、順次予測器206は、1つまたは複数のパターンの定期的な読み取り要求201に経時的に従うことが予測される予測を生成する。これらのパターンは、メモリ参照がそれらの間に空間的な局所性を有するという事実から生じる。たとえば、プロセッサ102がプログラム命令を実行する場合、読み取り要求201のストリームは、システムバス103をトラバースする際に事実上順次的とすることができる。順次パターンに従ってアドレスを予測するために、「前方順次予測」として以下で説明するあるタイプの推測技法が、順次アドレスを予測することができる。次に、このタイプの推測技法について説明する。

20

#### 【0016】

前方順次予測器208は、いくつかの順次アドレスを昇順で生成するように構成される。したがって、プロセッサ102が、昇順アドレスのストリームを含む一連の読み取り要求201をシステムバス103に伝送すると、前方順次予測器208は、追加の昇順アドレスをプリフェッチするためにいくつかの予測を生成することになる。前方順次予測器(「FSP」)208の一例が、図3Aに示される。図3Aに示されるように、FSP208は、アドレスA0などのアドレスを受け取り、A0アドレスから前方(すなわち昇順)順に1つまたは複数のアドレスを生成する。A0の表記法は、1つまたは複数の予測の形成元である基本アドレス(すなわちA+0)を識別する。したがって、表記法A1、A2、A3などはA+1、A+2、A+3などのアドレスを表し、表記法A(-1)、A(-2)、A(-3)などはA-1、A-2、A-3などのアドレスを表す。これらの表記法は1アドレスごとに昇順または降順の一連のアドレスを表すが、任意のパターン可能なアドレスセットを順次的と呼ぶことが可能である。全体を通じて使用される場合、順次アドレスを単一の文字によって表し、単一の文字とみなすことができる。たとえば、「A」はA0、A1、A2、A3などを表し、「B」はB0、B1、B2、B3などを表す。したがって、「A」および「B」はそれぞれ順次アドレスストリームを表すが、「B」のアドレスストリームは「A」のそれに対して非順次的である。

30

40

#### 【0017】

さらに図3Aでは、FSP208は少なくともイネーブル信号およびバッチ信号を受け取るように示され、その両方が予測コントローラ202によって提供される。イネーブル信号は、前方順次予測が生成されるかどうかを制御し、生成される場合、バッチ信号はFSP208が生成する順次アドレスの数を制御する。この例では、バッチ信号は、基本アドレスの他に「7つ」のアドレスが予測されることを示す。したがって、FSP2

50

08は、前方順のアドレスA1からA7を生成する。したがって、スペキュレータ108が読み取り要求201の一部としてA0などのアドレスを受け取る場合、順次予測器206は、予測203の一部としてアドレスA1、A2、A3、...、Abを提供することが可能であり、ここでbは「バッチ」の数である。

#### 【0018】

図2のブラインドバック順次予測器210は、1つの順次アドレスを生成するように構成されるが、基本アドレスからは降順である。ブラインドバック順次予測器（「ブラインドバック」）210の一例が図3Bに示され、ここでは、ブラインドバック順次予測器210がアドレスA0などの1つまたは複数のアドレスを受け取り、A0アドレスから後方（すなわち降順）順にアドレスA(-1)などの予測を1つだけ生成する。FSP

10

208の場合と同様に、ブラインドバック順次予測器210は、後方予測を生成するかどうかを制御するためのイネーブル信号も受け取る。

#### 【0019】

図2のバックセクタ順次予測器214は、システムバス103から他の特定のキャッシュラインを検出した後に、特定のキャッシュラインを予測として生成するように構成される。とりわけ、バックセクタ順次予測器214は、ある一定の読み取り要求201が高位キャッシュライン用であることを検出した場合、関連付けられた低位キャッシュラインが予測として生成される。高位キャッシュラインは、奇数アドレスを含む上位（「フロント」）セクタと呼ぶことが可能であり、低位キャッシュラインは、偶数アドレスを含む下位（「バック」）セクタと呼ぶことが可能である。例示のために、キャッシュラインは128

20

#### 【0020】

バックセクタ順次予測器214の一例が図3Cに示されており、ここでは、1つまたは複数のアドレスを受け取るバックセクタ順次予測器（「バックセクタ」）214が示されている。アドレスAUなどのキャッシュラインの上位またはフロントセクタに関する読み取り要求201を受け取ると、バックセクタ順次予測器214は1つだけの予測、アドレスALを生成する。このタイプの推測技法は、プロセッサ102が、典型的にはキャッシュラインの上位またはフロントセクタを要求した後に、下位またはバックセクタを要求する

30

#### 【0021】

図2の逆順次予測器212は、いくつかの順次アドレスを降順で生成するように構成される。したがって、プロセッサ102が、降順アドレスのストリームを含む一連の読み取り要求をシステムバス103に伝送すると、逆順次予測器212は、追加の降順アドレス用にいくつかの予測を生成することになる。逆順次予測器（「RSP」）212の一例が図3Dに示される。図3Dに示されるように、RSP 212は、アドレスA0、A(-1)、およびA(-2)などのアドレスのストリームを検出し、これに回答して、基本アドレスA0から逆（すなわち降順）順で1つまたは複数のアドレスを生成する。図3Dは、RSP 212が少なくともイネーブル信号、バッチ信号、および信頼レベル（「Conf.」）信号を受け取ることも示し、これら信号のすべてが予測コントローラ202によって提供される。イネーブル信号およびバッチ信号はFSP 208で使用される場合と同様に動作するが、信頼レベル（「Conf.」）信号は、逆順の予測の生成をトリガする時点を定義するしきい値を制御する。

40

#### 【0022】

さらに図3Dは、本発明の特定の実施形態に従った例示的RSP 212の挙動を示す図310を示す。ここでは、信頼レベル「2」がトリガレベル312を設定し、バッチ信号は、トリガアドレス以外に「5つ」のアドレスが予測されることを示す。トリガアドレスとは、予測器に予測を生成させるアドレスのことである。間隔I1中にA(0)を検出

50



した後、RSP 212が、続く間隔I2中にアドレスA(-1)も検出すると考えてみる。次に、間隔I3中にアドレスA(-2)を検出すると、検出されたストリームが一連の降順アドレスであるというある一定の信頼レベルに達する。トリガレベル312を上回った場合にこの信頼レベルに達し、これによって、RSP 212が逆順のアドレスA(-3)からA(-7)を生成する。したがって、スペキュレータ108が一連の読み取り要求201としてA0、A(-1)、およびA(-2)などの一定数のアドレスを受け取る場合、その後順次予測器206は、予測203の一部としてアドレスA(-3)、A(-4)、A(-5)、...、Abを提供することが可能であり、ここでbは「バッチ」の数である。いくつかの実施形態では、RSP 212は信頼レベルを採用せず、基本アドレス後に始まる予測を生成することに留意されたい。本発明の他の実施形態では、本明細書で説明される他の予測器において信頼レベルの概念が採用される。RSP 212および順次予測器206の他の構成予測器の制御については、以下でより詳細に論じるが、次に、図2の非順次予測器216について説明する。

#### 【0023】

非順次予測器216は、たとえアドレスが読み取り要求201の非線形ストリーム内にある場合でも、スペキュレータ108によって検出されたアドレスに続いて1つまたは複数の予測(すなわち予測アドレス)を生成するように構成される。典型的には、次のアドレスを予測する要求アドレスの顕著なパターンがない場合、先行アドレスのみに基づく予測は困難である。しかしながら本発明の実施形態によれば、非順次予測器216は、1つまたは複数の先行アドレスからパターン化不可能な予測アドレスを含む非順次予測を生成する。「パターン化不可能」な予測とは、先行アドレスによってパターン化できないか、または先行アドレスに対して不規則な、予測のことである。パターン化不可能な予測のタイプの1つが、非順次予測である。非順次予測が基づく先行アドレスは、即値アドレス、またはトリガアドレスとして構成された任意のアドレスの、いずれかとすることができる。とりわけ、読み取り要求201のストリーム内の2つまたはそれ以上のアドレスにわたって1つまたは複数のパターンが欠如していることは、プロセッサ102が、メモリロケーションの様々な空間的位置からの命令およびデータのフェッチに関してやや散漫な様式でプログラム命令を実行していることを示す。

#### 【0024】

非順次予測器216は、先行アドレスから非順次予測として分類可能な1つまたは複数の可能な非順次アドレスへの関連を格納するためのリポジトリとして、ターゲットキャッシュ218を含む。ターゲットキャッシュ218は、タイムリーな形で非順次予測を生成するために、そのコンテンツと着信する検出されたアドレスとを容易に比較するように設計される。非順次予測の生成元である検出されたアドレスは、「トリガ」アドレスと呼ばれ、結果として生じる予測は、この2つの間のパターン化不可能な関連の「ターゲット」である。次に、例示的非順次予測器216について説明する。

#### 【0025】

図4は、本発明の一実施形態に従った、例示的非順次予測器216を示す図である。非順次予測器216は、ターゲットキャッシュ422であるリポジトリに動作可能に結合された、非順次予測エンジン(「NonSeq.予測エンジン」)420を含む。ターゲットキャッシュ422は、各トリガアドレスと1つまたは複数の対応するターゲットアドレスとの間の関連を維持する。図4が、非順次アドレスを関連付ける多くの方法のうちの1つを示すことに留意されたい。ここでは、ツリー構造が特定のトリガアドレスをその対応するターゲットアドレスに関係付けている。この例では、ターゲットキャッシュ422は、アドレス「B」、「X」、および「L」などの可能な非順次予測のアドレスへの関連を形成する元となる、トリガアドレスとしてのアドレス「A」を含む。これら3つのターゲットアドレスは、それぞれアドレス「C」および「G」、「Y」、ならびに「M」に対するトリガアドレスでもある。ターゲットキャッシュ422の形成および動作について、以下でより詳細に考察する。アドレス「A」が、図4に示されていないトリガアドレスに対するターゲットアドレスともなり得ることに留意されたい。さらに、図示されていないア

10

20

30

40

50

ドレス間の多くの他の関連も可能である。

【0026】

非順次予測エンジン420は、少なくとも4つの信号および任意数のアドレス402を受け取るように構成される。非順次予測エンジン420の動作を制御するために、予測コントローラ202は「バッチ」信号および「イネーブル」信号を提供し、これらはどちらも前述の信号と事実上同様である。予測コントローラ202は、2つの他の信号、幅（「W」）信号および深さ（「D」）信号も提供する。これらの信号はターゲットキャッシュ422の形成を制御するものであり、幅信号Wはトリガアドレスの予測元となり得る可能なターゲットの数を設定し、深さ信号Dはトリガアドレスに関連付けられるレベルの数を設定する。後者の例は、Dが深さ「4」を示す場合である。これは、アドレスAが第1のレベルであり、アドレスBが第2のレベルであり、アドレスCおよびGが第3のレベルであり、アドレスDが第4のレベルであることを意味する。前者の例は、Wが「2」に設定された場合である。これは、3つのアドレス「B」、「X」、および「L」のうちのみが非順次予測に使用されることを意味する。

10

【0027】

図4は、それぞれが以前に検出されたアドレスにパターン化不可能なアドレスを含む、概念上、非順次アドレスストリーム404、406、408、410、および412で示されたアドレスなどの、例示的地址402を、予測コントローラ202から受け取るように構成された、非順次予測エンジン420も示す。たとえば、ストリーム404はアドレス「A」を含み、その後アドレス「B」が続き、さらにその後アドレス「C」が続く。非順次アドレスの場合と同様に、「A」から「B」を予測するため、および「B」から「C」を予測するためのパターンを検出することは、プロセッサ102からの読み取り要求201を監視するだけでは困難である。このため、非順次予測器216は、特定のトリガアドレスとそのターゲットアドレスとの間のパターン化不可能な関連の予測を実行可能にするために、ターゲットキャッシュ422を形成する。非順次予測エンジン420は、非順次予測を形成すると、関連付けられたターゲットアドレスから予測のグループを生成する。したがって、トリガアドレス「A」がアドレス「B」（すなわち、基本アドレスとしてのB0）の非順次予測につながる場合、予測アドレスはB0、B1、B2、...、Bbを含むことになり、ここでbはバッチ信号によって設定される数である。

20

【0028】

本発明の一実施形態では、非順次予測エンジン420は、アドレス402のそれぞれから後続のアドレスへの関連を格納しながら、ターゲットキャッシュ422を形成する。たとえば、ストリーム404のアドレスAを検出すると、非順次予測エンジン420は、AからBへの関連、BからCへの関連、CからDへの関連などの関連を、ターゲットキャッシュ422に追加する。非順次予測エンジン420は、他のストリーム406、408などのアドレスを検出した場合も、同様に実行する。

30

【0029】

特定の実施形態によれば、ターゲットキャッシュ420は、これらの関連を表430、440、および450などの表形式で格納する。これらの表は、それぞれトリガアドレスとターゲットアドレスとの間の関連を格納するための、トリガ列426およびターゲット列428を含む。次に、すべてのストリームのアドレス402が、ターゲットキャッシュ422の表430、440、および450に格納されるものと考えてみる。表430に示されるように、トリガ-ターゲット関連432、434、および436は、それぞれAからB、BからC、およびGからQへの関連を記述する。他のトリガ-ターゲット関連438は、CからDなどの関連を含む。同様に、表440はAからXへの関連を記述するためのトリガ-ターゲット関連442を含み、表450はAからLへの関連を記述するためのトリガ-ターゲット関連452を含む。

40

【0030】

図4は、表430、440、および450がそれぞれ、同じトリガアドレスに関する複数のトリガ-ターゲット関連の相対的優先順位を記述する、「Way 0」、「Way

50

1」、および「Way 2」として識別されることを示す。この場合、Way 0は最高優先順位に、Way 1は第2位優先順位に、という具合に関連付けられる。この例では、表430のトリガ-ターゲット関連432は、AからBへの関連が、表440のトリガ-ターゲット関連442であるAからXへの関連よりも優先順位が高いことを示す。したがって、ターゲットキャッシュ422がこれらの関連を含んだ後、次に非順次予測エンジン420がアドレスAを検出した場合（予測コントローラ202が非順次予測エンジン420を動作させられる限り）、表の相対的優先順位により、アドレスBは最高優先順位、続いてアドレスXは第2位優先順位、として予測されることになる。

#### 【0031】

本発明の一実施形態によれば、相対的優先順位は少なくとも2つの方法で決定される。第1に、トリガ-ターゲット関連が最初に検出され、ターゲットキャッシュ422内に配置された場合、これには最高優先順位が関連付けられる。第2に、非順次予測エンジン420が、トリガ-ターゲット関連が成功である（たとえば、その特定の関連に基づく非順次予測の結果として生じた、最新のキャッシュヒットが存在する）と判定した場合、そのトリガ-ターゲット関連には最高優先順位が関連付けられる。「最新の」キャッシュヒットとは、特定のトリガアドレスに関連付けられたターゲットアドレスのうちの少なくとも1つの新しいキャッシュヒットのことである。さらに、以前の「最高優先順位」（leg 0としても指定される）は、対応する関連をWay 1の表に移動させることによって、第2位優先順位（leg 1としても指定される）にシャッフルされる。一例として、AからXへの関連が第1のトリガ-ターゲット関連としてターゲットキャッシュ422に導入される場合の、第1の時点について考えてみる。結果として、表430内（すなわちway 0）に配置されることにより、最高優先順位（すなわち初期にはleg 0）が関連付けられることになる。何らかのその後の時点で、ターゲットキャッシュ422はAからBへの関連を表430に挿入する（最高優先順位、leg 0）。また、AからXへの関連も表440に移動される（第2位優先順位、leg 1）。本発明の特定の実施形態では、トリガ-ターゲット関連が格納される表は、インデックスを構成するアドレスビットの一部に依存する。

#### 【0032】

再度図2を参照すると、予測コントローラ202は、順次予測器206および非順次予測器216の両方を制御するように構成される。予測コントローラ202は、順次予測器206または非順次予測器216のいずれか、あるいはその両方によって生成される予測の量ならびにタイプを制御する。また、予測コントローラ202は、冗長または重複予測などの、そうでなければ不要な予測203の生成も抑制する。予測器208、210、212、214、および216のそれぞれが同時に動作可能な場合、予測203の数はプリフェッチャリソースに負荷をかけすぎないように管理しなければならない。予測コントローラ202は、この操作および他の同様の操作を実行するために、抑制器204を採用する。

#### 【0033】

本発明の一実施形態では、抑制器204が生成される予測の量を制御する。これは、第1に読み取り要求201の一定の属性を確認することによって実行される。具体的に言えば、抑制器204は、読み取り要求201がプログラム命令（すなわち「コード」）またはプログラムデータ（すなわち「非コード」）のどちらに関連するかを判定する。通常、プログラムデータ以外のコードを取り出すための読み取り要求201は、事実上より順次的であるか、または少なくともパターン化可能である傾向がある。これは、プロセッサ102が一般に、プログラムデータに対する要求よりも線形的な方法で命令を実行するためである。したがって、抑制器204は、順次予測器206および非順次予測器216に対して、読み取り要求201がプログラムデータに関係する場合は予測生成を抑制するように指示することができる。これにより、擬似予測の生成防止に役立つ。

#### 【0034】

抑制器204は、読み取り要求201が非プリフェッチ「デマンド」またはプリフェッ

10

20

30

40

50

子のいずれであるかを確認することによって、順次予測器 206 および非順次予測器 216 が生成する予測の量を調整することもできる。プロセッサ 102 は、通常、確実に必要な何らかのケースにおいて、プログラム命令またはプログラムデータをメモリ 112 から取り出すように要求する（非プリフェッチデマンドとして）が、プロセッサ 102 は、後で必要となることを見込んで、プログラム命令またはプログラムデータをプリフェッチするように要求するだけでも可能である。確実に必要である方が必要となる見込みであるよりも重要な処理である可能性が高いため、抑制器 204 は特定の予測器に対し、デマンド読み取り要求 201 に基づく予測を優先して、プリフェッチ読み取り要求 201 に基づく予測を抑制するように指示することができる。

【0035】

表 1 は、生成される予測の数を抑制するための例示的技法を示す。すなわち、読み取り要求 201 がコードおよびデマンドの両方に関連する場合、抑制器 204 は最も抑制的でなくなる。すなわち、予測コントローラ 202 は「バッチ」を表 1 でバッチサイズ (4) と示されるような大規模サイズに設定することになる。特定の例では、バッチサイズ (4) は 7 に設定することができる。しかしながら、前述の理由では、読み取り要求 201 がプログラムデータ（すなわち非コード）およびプロセッサ生成プリフェッチの両方に関係する場合、抑制器 204 は最も抑制的となる。したがって、予測コントローラ 202 は「バッチ」を表 1 でバッチサイズ (1) と示されるような小規模サイズに設定することになる。一例として、バッチサイズ (1) は 1 に設定することができる。他のケースでは、予測コントローラ 202 は、バッチサイズ (2) およびバッチサイズ (3) などの他のバッチサイズを使用して、予測抑制のレベルを変更することができる。本発明の一実施形態に従った抑制器は、プロセッサ要求がデータまたはプリフェッチ要求、あるいはその両方に対するものである場合、「バッチ」量を減少させることによって、少なくとも 1 つの予測アドレスの生成を抑制するように構成されるが、表 I はこれに限定されない。たとえば、コードまたは命令に関するプロセッサ要求は、「バッチ」サイズを増加させるのではなく減少させることができる。他の例として、デマンドに関する要求も、「バッチ」サイズを増加させるのではなく減少させることができる。当業者であれば、多くの変形形態が本発明の範囲内であることを理解されよう。

【表 1】

コードまたはデータ	デマンドまたはプリフェッチ	バッチ
非コード (すなわちデータ)	プリフェッチ	バッチサイズ (1)
非コード (すなわちデータ)	デマンド	バッチサイズ (2)
コード	プリフェッチ	バッチサイズ (3)
コード	デマンド	バッチサイズ (4)

表1. 読み取り要求のタイプ

【0036】

抑制器 204 は、順次予測器 206 および非順次予測器 216 が生成する予測のタイプも調整することができる。第 1 に、予測コントローラ 202 が前方順次予測器 208 および逆順次予測器 212 の両方を同時に実行可能にできるものと考えてみる。したがって、抑制器 204 は、プロセッサ 102 が降順でアドレス読み取りを要求している場合、昇順でのアドレスの予測を最小限にするために、逆順次予測器 212 がトリガする（すなわち信頼レベルを超える）場合、少なくとも前方順次予測器 208 を実行不可にするよう、予

測コントローラ 202 に指示する。

【0037】

第2に、予測コントローラ 202 が順次予測（すなわち、前方順次予測器 208 または逆順次予測器 212 のいずれか）を動作可能にする場合、あるアドレスが、バック予測（すなわち、ブラインドバック順次予測器 210 またはバックセクタ順次予測器 214 のいずれか）をトリガすると考えてみる。この場合、抑制器 204 は、前方順次予測器 208 または逆順次予測器 212 のいずれかについて、その初期量から 1 ずつバッチを抑制する。すなわち、「バッチ」が初期に 7 に設定されていた場合、ブラインドバック順次予測器 210 またはバックセクタ順次予測器 214 のいずれかがトリガまたは活動化されると同時に、「バッチ」は 1 つだけ減少することになる。たとえば、前方順次予測器 208 に関するバッチがアドレス A0、A1、A2、...、A7 を生成するように設定されている場合、およびブラインドバック順次予測器 210 が 1 つまたは複数の読み取り要求 201 に対して実行可能である場合、前方順次予測器 208 は予測 A1、A2、...、A6 のみを生成する。最終結果は、それらの読み取り要求 201 に対して、予測 A(-1)、A(0)、A1、A2、...、A6 のセットとなり、バック予測は予測 A(-1) を提供する。

10

【0038】

第3に、予測コントローラ 202 はオプションで、プロセッサからのアドレス 201 の順次ストリームで最初に予測が生成された後に、ブラインドバック順次予測器 210 またはバックセクタ順次予測器 214 のいずれかに対して、それらの予測を抑制不可にすることができる。これは、シーケンスの基本アドレスが確立された後、後続の前方または逆順次予測も後方タイプの推測（たとえ 1 アドレス後方であっても）を予測するためである。たとえば、前方順次予測 A2、A3、および A4 は、すべてがすでに予測されている（基本アドレスが A0 の場合）後方タイプの予測 A1、A2、および A3 もカバーしている。抑制器 204 は、他のタイプの予測を抑制するように構成することが可能であり、次にその例について説明する。

20

【0039】

図5は、本発明の一実施形態に従った、非順次予測を抑制する例示的技法を示す図である。この技法によれば、抑制器 204 は、ターゲットキャッシュ 422 へのトリガ - ターゲット関連の格納を必要とする、そうでなければ非順次とみなされる可能性のあるインタリーブされた順次ストリームを検出する。リソース、特にターゲットキャッシュ 422 内の使用可能なメモリを保存するために、抑制器 204 はストリーム 502 内などの非順次アドレスを解析し、それらの非順次アドレスをインタリーブされた順次ストリームとしてモデル化する。図に示されるように、ストリーム 502 は、それぞれが各間隔 I1、I2、I3、I4、I5、I6、I8、および I9 中に検出された、アドレス A0、B0、C0、A1、B1、C1、A2、B2、および C2 からなる。抑制器 204 は、非順次アドレスを順序どおりにモデル化するための、表 504 などのデータ構造を含む。表 504 は、ストリーム 502 を分解するための任意数のストリームトラックを含むことができる。具体的に言えば、ストリームトラック 520、522、および 524 は、それぞれ順次ストリーム B0、B1、および B2 と、A0、A1、および A2 と、C0 および C1 とをモデル化するように設計される。A7（図示せず）などの、ストリーム 502 から後に検出される読み取りアドレスがこれらのストリームと比較され、追跡されているストリームに対して非順次予測が依然として抑制可能であるかどうか調べられる。

30

40

【0040】

抑制器 204 は動作時に、シーケンスの第1のアドレスなどの基本アドレス 510 を格納することによって、順次ストリームを追跡する。その後抑制器 204 は、最新検出アドレス 514 を維持する。各新しい最新検出アドレスについて（たとえば、ストリームトラック 520 の B2）、前の最新検出アドレス（たとえばストリームトラック 520 の B1）は、オプション列である列 512 内に配置されることによって無効（「無効」）とされる。抑制器 204 は、この例示的技法を使用して、他のタイプの予測が使用可能な場合に

50

不必要な非順次予測の生成を抑制する。したがって図5に示された例では、前方順次予測器208がストリーム502に関する予測を適切に生成することができる。

【0041】

図6は、本発明の一実施形態に従った、非順次予測を抑制する他の例示的技法を示す図である。この技法によれば、抑制器204は、図5に記載されたプロセスと同様のインタリーブされた順次ストリームとして非順次アドレスをモデル化する。しかしながら、図6の技法は、任意数のスレッドにわたる順次ストリームを検出するためにそれぞれが使用される、複数のデータ構造を実施する。この例では、表604、606、および608は、それぞれスレッド(0)('T')、スレッド(1)('T')、およびスレッド(2)('T')用のストリームトラックを含む。したがってストリーム602の非順次アドレスは、この技法を使用して、非順次予測を抑制するように、複数のスレッドにわたる複数の順次ストリームとしてモデル化することができる。この技法は、逆順次ストリームまたは他のタイプの予測に適用可能であることに留意されたい。

10

【0042】

図7は、本発明の特定の実施形態に従った、非順次予測を抑制するための他の技法を示す図である。アドレスストリーム702の場合、アドレスA4とB0との間に非順次性が存在する。しかしながらいくつかのケースでは、これらの要求された読み取りアドレス間の時間差が非常に短い場合、非順次予測を採用する十分な時間がないことになる。抑制器204のマッチャ706は、アドレスA4とB0との間の時間差dを比較するように動作する。dがしきい値THに等しいかまたはこれより大きい場合、マッチャ706は非順次予測器216を動作可能にする(すなわち「抑制しない」)信号を送る。しかしながら、dがTHよりも小さい場合、マッチャ706は非順次予測器216を使用不可にする信号を送り、それによって予測を抑制する。

20

【0043】

抑制器204によって実施可能な他の抑制メカニズムは、以下のとおりである。一般に、プロセッサ102がフロントセクタアドレスを要求した後、バックセクタアドレスに関して要求を出すまでに経過する時間量は有限である。この時間量が十分に長い場合、バックセクタアドレスの読み取り要求は不規則(すなわちフロントセクタに対してパターン化不可能)に見える可能性がある。これを防止するために、抑制器204は、プロセッサ102によるフロントセクタ読み取りのリストを維持するように構成される。フロントセクタアドレスの検出に続いて、そのフロントセクタアドレスとアドレスが比較される。対応するバックセクタに達すると、そのように認識されることになる。したがって、その他の形の非順次性ならびにその予測を抑制することができる。

30

【0044】

図8は、本発明の特定の実施形態に従った、予測の生成を促進するための例示的技法を示す図である。具体的に言えば、促進器205(図2)は、この技法に従って非順次予測の生成を早めるように動作する。この例では、ストリーム802は2つの隣接する順次ストリームA0からA4およびB0からB3を含む。非順次予測器216は、通常、アドレスA4をトリガアドレス808として指定し、アドレスB0をターゲットアドレス810として指定する。しかしながら、非順次予測を生成するための時間を減少させるために、トリガアドレス808を新しいトリガアドレス804(すなわちA0)に変更することができる。したがって、ターゲットアドレスに新しいトリガアドレスを指定することで、次にプロセッサ102がストリーム802のアドレスを要求した場合、非順次予測器216は最近のアドレスよりも前のアドレスを検出すると同時に直ちにその予測を生成することができる(すなわち、A4ではなくA0が「新しい」トリガアドレスとして検出された場合、予測を生成する)。これにより、最も適切な時点での非順次予測の生成が保証される。

40

【0045】

図9は、本発明の一実施形態に従った、他の例示的スペキュレータを示す図である。この例では、プリフェッチャ900は、不必要な予測生成を最小限に維持するように冗長ア

50

ドレスをフィルタリングするためのフィルタ 914 を備えたスペキュレータ 908 を含む。図 9 のプリフェッチャ 900 は、マルチレベルキャッシュ 920 および予測インベントリ 916 も含む。ここで、マルチレベルキャッシュ 920 は、第 1 レベルの戻りデータキャッシュ (「DRC1」) 922 および第 2 レベルの戻りデータキャッシュ (「DRC2」) 924 からなる。第 1 レベルの戻りデータキャッシュ 922 は、一般に、短期データストアとして説明可能であり、第 2 レベルの戻りデータキャッシュ 924 は、一般に、長期データストアとして説明可能である。マルチレベルキャッシュ 920 は、メモリ 112 からプリフェッチされたプログラム命令およびプログラムデータを、プロセッサ 102 が必要とするまで格納する。同様に、予測インベントリ 916 は、メモリ 112 にアクセスするためにアービタ 918 によって選択されるまで、生成された予測に一時ストレージを提供する。アービタ 918 は、アービトレーション規則に従って、メモリ 112 にアクセスして命令およびデータをプリフェッチするために、生成されたどの予測が発行されるかを決定するように構成される。

10

#### 【0046】

フィルタ 914 は、キャッシュフィルタ 910 およびインベントリフィルタ 912 の、少なくとも 2 つのフィルタを含む。キャッシュフィルタ 910 は、新しく生成された予測と、マルチレベルキャッシュ 920 内にすでに格納された命令およびデータをプリフェッチした以前の予測とを、比較するように構成される。したがって、マルチレベルキャッシュ 920 に関して、1 つまたは複数の新しく生成された予測が、任意の以前に生成された予測と重複する場合、予測の数を最小にするようにそれらの冗長予測は無効とされる。さらにインベントリフィルタ 912 は、新しく生成された予測と、すでに生成され、予測インベントリ 916 に格納された予測とを、比較するように構成される。したがって、1 つまたは複数の新しく生成された予測が予測インベントリ 916 に格納された予測と重複する場合、予測の数を最小にするように任意の冗長予測は無効とされ、それによってプリフェッチャリソースが解放される。

20

#### 【0047】

非順次予測器に関する例示的实施形態

図 10 は、本発明の特定の实施形態に従った、例示的非順次 (「NONSEQ」) 予測器 1010 を示すブロック図である。この例では、非順次予測器 1010 は、順次予測を生成するための順次予測器 1012 も含むスペキュレータ 1008 内に常駐するように示される。プリフェッチャ 1006 はスペキュレータ 1008 を含み、要求される前に (図示せず) メモリからプログラム命令およびプログラムデータの両方を「フェッチ」し、その後プロセッサ (図示せず) によって要求されると、フェッチされたプログラム命令およびプログラムデータをそのプロセッサに提供するように動作する。それらを使用する前にフェッチする (すなわち「プリフェッチ」) ことによって、プロセッサのアイドル時間 (たとえば、プロセッサがデータ不足である間の時間) が最小となる。非順次予測器 1010 は、予測を生成するための非順次予測エンジン (「予測エンジン」) 1020 と、予測を格納および優先順位付けするためのターゲットキャッシュ 1030 とを含む。

30

#### 【0048】

プリフェッチャ 1006 は、フィルタ 1014、オプションの予測インベントリ 1016、オプションのアービタ 1018、およびマルチレベルキャッシュ 1040 も含む。ここで、フィルタ 1014 は、新しく生成された予測と、プログラム命令およびプログラムデータをすでにマルチレベルキャッシュ 1040 にプリフェッチされた状態にする以前の予測とを比較するように構成された、キャッシュフィルタ (図示せず) を含む。したがって、任意の新しく生成された予測が、マルチレベルキャッシュ 1040 に格納された任意の以前に生成された予測と重複する場合、予測の数を最小にするようにその冗長予測は無効とされ、それによってプリフェッチャリソースが解放される。予測インベントリ 1016 は、メモリにアクセスするためにアービタ 1018 によって選択されるまで、生成された予測を格納するための一時ストレージを提供する。アービタ 1018 は、メモリにアクセスして命令およびデータをプリフェッチするために、生成されたどの予測が発行される

40

50

かを決定するように構成される。

【0049】

マルチレベルキャッシュ1040は、第1レベルの戻りデータキャッシュ(「DRC1」)1042および第2レベルの戻りデータキャッシュ(「DRC2」)1044からなる。第1レベルの戻りデータキャッシュ1042は、一般に、短期データストアとして説明可能であり、第2レベルの戻りデータキャッシュ1044は、一般に、長期データストアとして説明可能である。本発明の実施形態によれば、第1レベルの戻りデータキャッシュ1042または第2レベルの戻りデータキャッシュ1044、あるいはその両方が、予測アドレス(すなわちターゲットアドレス)に基づいてプリフェッチされた、プロフェッチ済みプログラム命令およびプログラムデータを格納することができる。図に示されるように、マルチレベルキャッシュ1040に格納されたプリフェッチ済み予測情報は、データ(TRT1)およびデータ(TRT2)として表される。この表記法は、ターゲットアドレスTRT1およびTRT2が、予測情報を表すデータのプリフェッチに寄与したことを意味する。図に示されるように、ならびに以下で論じるように、データ(TRT1)およびデータ(TRT2)は、それぞれ、予測識別子(「PID」)1および2と共にマルチレベルキャッシュ1040内に格納される。データ(TRT1)またはデータ(TRT2)のいずれかがプロセッサによって要求された場合、対応するターゲットアドレス(たとえばTRT1)および予測識別子が非順次予測器1010に送られることになる。

10

【0050】

スペキュレータ1008は、動作時に、プロセッサがメモリへのアクセスを要求する(「読み取り要求」)場合、システムバスを監視する。プロセッサがプログラム命令を実行する場合、スペキュレータ1008は、まだプロセッサによって使用されていないプログラム命令およびプログラムデータを含むアドレスに関する読み取り要求を検出する。考察のために、「アドレス」は、一般にメモリと、マルチレベルキャッシュ1040などのキャッシュメモリとの間で転送される、メモリのキャッシュラインまたは単位に関連付けられる。キャッシュメモリは、ターゲットキャッシュ1030の外部リポジトリの一例であることに留意されたい。

20

【0051】

検出された読み取り要求に基づいて、非順次予測器1010は、プロセッサによって次に要求される可能性のある構成可能な数の予測アドレスを生成することができる。具体的に言えば、非順次予測器1010は、たとえアドレスが読み取り要求の非線形ストリーム内にある場合でも、そのアドレスの検出に続いて、1つまたは複数の予測(すなわち予測アドレス)を生成するように構成される。典型的には、次のアドレスを予測する要求アドレスの顕著なパターンがない場合、先行アドレスのみに基づく予測は困難である。しかしながら本発明の実施形態によれば、非順次予測エンジン1020は、1つまたは複数の先行アドレスからパターン化不可能な予測アドレスを含む、非順次予測を生成する。「パターン化不可能」な予測とは、先行アドレスによってパターン化できないか、または先行アドレスに対して不規則な、予測のことである。パターン化不可能な予測のタイプの1つが、非順次予測である。非順次予測が基づく先行アドレスは、即値アドレス、またはトリガアドレスとして構成された任意のアドレスの、いずれかとすることができる。とりわけ、読み取り要求のストリーム内の2つまたはそれ以上のアドレスにわたって1つまたは複数のパターンが欠如していることは、プロセッサが、メモリロケーションの様々な空間的位置からの命令およびデータのフェッチに関してやや散漫な様式でプログラム命令を実行していることを示す。

30

40

【0052】

非順次予測器1010は、先行アドレスから非順次予測としてそれぞれ分類可能な1つまたは複数の潜在的な非順次アドレスへの関連を格納するためのリポジトリとして、ターゲットキャッシュ1030を含む。ターゲットキャッシュ1030は、迅速な形で非順次予測を生成するために、そのコンテンツと着信する検出されたアドレスとを比較するように設計される。さらにターゲットキャッシュ1030は、たとえばキャッシュメモリ内の

50



ヒットにตอบสนองして、それらの非順次予測を優先順位付けするように構成される。あるいは、非順次予測器 1010 は、新しい非順次予測と特定のトリガアドレスとの間の関連を確立する第 1 のインスタンスを優先順位付けすることができる。「トリガ」アドレスとは、非順次予測器 1010 が非順次予測を生成する元となる検出されたアドレスのことであり、2 つの間のパターン化不可能な関連の「ターゲット」と呼ばれる、結果として生じる予測を伴う。本発明の少なくとも 1 つの実施形態によれば、ターゲットキャッシュ 1030 は、それ以外の場合には複数ポートメモリによって使用される、リソースを節約するための単一ポートメモリとすることができることに留意されたい。

#### 【0053】

プリフェッチャ 1006 が非順次予測器 1010 から予測を発行した後に、非順次予測を使用してメモリにアクセスする。これにตอบสนองして、メモリは、予測アドレスに関する参照と共にプリフェッチされたデータを戻し、ここで参照は、予測識別子（「PID」）および対応するターゲットアドレスを含むことができる。その後、マルチレベルキャッシュメモリ 1040 は、プロセッサが要求する時点などまで、戻されたデータを一時的に格納する。以下で説明するように、プロセッサがプリフェッチ済みデータ（すなわち予測情報）を要求する場合、必要であれば非順次予測の優先順位を再調整するために参照が非順次予測器 1010 に送信される。

#### 【0054】

図 11 は、本発明の一実施形態に従った、例示的非順次予測器 1010 を示す図である。非順次予測器 1010 は、ターゲットキャッシュ 1130 によって例示されたりポジトリに動作可能に結合された、非順次予測エンジン（「Non Seq. 予測エンジン」）1120 を含む。さらに非順次予測エンジン 1120 は、予測生成器 1122 および優先順位調整器 1124 を含む。予測生成器 1122 は予測を生成し、ターゲットキャッシュ 1130 に格納されたトリガ - ターゲット関連を管理する。優先順位調整器 1324 は、たとえば、最も新しい正常なターゲットアドレスから、最も古いかまたは正常でないターゲットアドレスへと、トリガ - ターゲット関連を優先順位付けするように動作する。予測生成器 1122 および優先順位調整器 1124 については、それぞれ図 12 および 13 でより詳細に説明する。

#### 【0055】

ターゲットキャッシュ 1130 は、各トリガアドレス（「TGR」）と 1 つまたは複数の対応するターゲットアドレス（「TRT」）との間の関連を維持する。図 11 は、非順次アドレスを関連付けるために使用する多くの方法のうちの 1 つを示すことに留意されたい。ここでは、ツリー構造が特定のトリガアドレスをその対応するターゲットアドレスに關係付けている。この例では、ターゲットキャッシュ 1130 は、アドレス「B」、「X」、および「L」などの可能な非順次予測のアドレスへの関連を形成する元となる、トリガアドレスとしてのアドレス「A」を含む。これら 3 つのターゲットアドレスは、それぞれアドレス「C」および「G」、「Y」、ならびに「M」に対するトリガアドレスでもある。特に、予測生成器 1122 が新しいトリガ - ターゲット関連を発見し、その関連をターゲットキャッシュ 1130 に挿入する場合の、ターゲットキャッシュ 1130 の形成および動作について、以下でより詳細に考察する。アドレス「A」が、図 11 に示されていないトリガアドレスに対するターゲットアドレスともなり得ることに留意されたい。さらに、図示されていないアドレス間の多くの他の関連も可能である。

#### 【0056】

図に示されるように、ターゲットキャッシュは、本発明の一実施形態により、幅（「w」）、深さ（「d」）、および高さ（「h」）の少なくとも 3 つの変数に従って、たとえば非順次予測エンジン 1120 によって構築することができる。幅 w は、トリガアドレスの予測元とすることができる可能なターゲットの数を設定し、深さ d は、トリガアドレスに関連付けられるレベルの数を設定する。高さ h は、非順次予測を生成するために使用される連続するトリガアドレスの数を設定する。一例として、d が深さ「4」を示すものと考えてみる。これは、アドレス A が第 1 のレベルであり、アドレス B が第 2 のレベルであ

10

20

30

40

50

り、アドレスCおよびGが第3のレベルであり、アドレスDが第4のレベルであることを意味する。他の例として、wが「2」に設定されるものと考えてみる。これは、3つのアドレス「B」、「X」、および「L」のうちの2つのみが  $1 \leq g \leq 0$  および  $1 \leq g$

1として非順次予測に使用され、3つのアドレスすべてが第2のレベルにあることを意味する。特定の実施形態では、変数hは、マルチレベルの予測生成を達成するために第1のレベルを超えるレベル数を設定する。

#### 【0057】

図11に示されるように、hが2に設定されるものと考えてみる。これは、第1のレベルのトリガアドレス(たとえばアドレスA)および次に続く第2のレベルのトリガアドレス(たとえばアドレスB)という、2つのレベルのトリガアドレスが存在することを意味する。したがって、hを2に設定すると、トリガアドレスAにตอบสนองして第1の予測グループが形成される。すなわち、第2のレベルのターゲットアドレスのうちのいずれかが非順次アドレスの1つまたは複数のグループを生成することができる。たとえば、アドレス「B」、「X」、および「L」のうちのいずれかが非順次予測を生成するための基準となることが可能であり、これらのアドレスの数は、非順次予測エンジン1120によって定義されたアクティブ  $1 \leq g$  の数(たとえば  $1 \leq g \leq 0$  から  $1 \leq g \leq 2$  まで)によって選択される。しかしながら、マルチレベル予測生成(およびhを2に設定すること)に従って、アドレス「B」、「X」、および「L」のそれぞれを、次の下位レベルのターゲットアドレスに基づいて予測の第2グループを生成するための連続するトリガアドレスとすることが可能である。したがって、第3レベルのターゲットアドレスCおよびGを使用して、連続するトリガアドレスBに基づいて追加の非順次予測を生成することができる。同様に、ターゲットアドレスYおよびMを使用して、それぞれ連続するトリガアドレスXおよびLに基づいて非順次予測を生成することもできる。当業者であれば、前述の3つの変数のうちの1つまたは複数を変更することによって、多くの実施が可能であることを理解されよう。

#### 【0058】

非順次予測エンジン1120は、読み取り要求の例示的地址1101を受け取るように構成される。図11は、それぞれが以前に検出されたアドレスにパターン化不可能なアドレスを含む、非順次アドレスストリーム1102、1104、1106、1108、および1110を概念的に示す。たとえば、ストリーム1102は、アドレス「A」を含み、その後アドレス「B」が続き、さらにその後アドレス「C」が続く。非順次アドレスの場合と同様に、「A」から「B」を予測するため、および「B」から「C」を予測するためのパターンを検出することは、読み取り要求1101を監視するだけでは困難である。このため、予測生成器1122は、特定のトリガアドレスとそのターゲットアドレスとの間のパターン化不可能な関連の予測を実行可能にするために、ターゲットキャッシュ1130のコンテンツを確立する。たとえば、ストリーム1102のアドレスA(ならびに後続のアドレス)を検出すると、予測生成器1122は、AからBへの関連、BからCへの関連、CからDへの関連などの関連を、ターゲットキャッシュ1130に追加する。非順次予測エンジン1120は、他のストリーム1104、1106などのアドレスを検出する場合も、同様に実行する。

#### 【0059】

特定の実施形態によれば、ターゲットキャッシュ1130はこれらの関連を、表1140、1150、および1160などの表形式で格納する。これらの表は、それぞれトリガアドレスとターゲットアドレスとを格納するための、トリガ列(「TGR」)およびターゲット列(「TGT」)を含む。次に、すべてのストリームのアドレス1101が、表1140、1150、および1160に格納されるものと考えてみる。表1140に示されるように、トリガ-ターゲット関連1142、1144、および1146は、それぞれAからB、BからC、およびGからQへの関連を記述する。他のトリガ-ターゲット関連1148は、CからDなどの関連を含む。同様に、表1150はAからXへの関連を記述するためのトリガ-ターゲット関連1152を含み、表1160はAからLへの関連を記述

10

20

30

40

50

するためのトリガ - ターゲット関連 1 1 6 2 を含む。

【 0 0 6 0 】

図 1 1 は、表 1 1 4 0、1 1 5 0、および 1 1 6 0 がそれぞれ、同じトリガアドレスに関するターゲットキャッシュ 1 1 3 0 内の複数のトリガ - ターゲット関連の相対的位置を記述する、「way 0」、「way 1」、および「way 2」として識別されることを示す。優先順位調整器 1 1 2 4 は、通常はメモリロケーションに優先順位を割り当てることによって、トリガ - ターゲット関連に優先順位、すなわち予測を割り当てる。この場合、way 0 は最高優先順位に、way 1 は第 2 位優先順位に、という具合に関連付けられる。この例では、表 1 1 4 0 のトリガ - ターゲット関連 1 1 4 2 は、A から B への関連が、表 1 1 5 0 のトリガ - ターゲット関連 1 1 5 2 である A から X への関連よりも優先順位が高いことを示す。したがって、ターゲットキャッシュ 1 1 3 0 がこれらの関連を含んだ後、次に非順次予測エンジン 1 1 2 0 がアドレス A を検出した場合、非順次予測エンジン 1 1 2 0 が 1 つまたは複数の予測を提供することができる。通常、非順次予測エンジン 1 1 2 0 は、優先順に生成される非順次予測を生成する。具体的に言えば、非順次予測エンジン 1 1 2 0 は、優先順位が下位の予測を生成する前に最高優先順位を有する予測を生成する。したがって、非順次予測エンジン 1 1 2 0 は、優先順位に基づいて構成可能な数の予測を生成することができる。たとえば、非順次予測エンジン 1 1 2 0 は、予測の数を、leg 0 および leg 1 (すなわち、トリガ - ターゲット関連の上位 2 つ) の 2 つに制限することができる。これは、何らかのケースで、非順次予測エンジン 1 1 2 0 が、表の相対的優先順位により、アドレス X ではなくアドレス B を提供する傾向がより高くなることを意味する。トリガ - ターゲット関連間での相対的優先順位がまさにそれ、すなわち相対的であることに留意されたい。これは、ターゲットキャッシュ 1 1 3 0 が、特定のトリガアドレスの最高優先順位関連をたとえば way 4 に位置付け、第 2 位の優先順位関連を way 9 に位置付けることができることを意味する。しかしながら、ターゲットキャッシュ 1 1 3 0 は、1 つのアドレスから、単なる leg 0 および leg 1 を超える任意の量の「leg」を含むことができることに留意されたい。

【 0 0 6 1 】

図 1 2 は、本発明の実施形態に従った、例示的予測生成器 1 2 2 2 を示す図である。この例では、予測生成器 1 2 2 2 は予測を生成するため、ならびにその中に格納されたトリガ - ターゲット関連を管理するために、ターゲットキャッシュ 1 2 3 0 に結合される。予測生成器 1 2 2 2 は、インデックス生成器 1 2 0 4、タグ生成器 1 2 0 6、ターゲット特定器 1 2 0 8、およびコンパイナ 1 2 1 0 を含む。また、予測生成器 1 2 2 2 は、発見されたトリガ - ターゲット関連をターゲットキャッシュ 1 2 3 0 に挿入するための挿入器 1 2 0 2 も含む。

【 0 0 6 2 】

予測を生成する場合、インデックス生成器 1 2 0 4 およびタグ生成器 1 2 0 6 は、他のアドレスに先行するアドレスとすることが可能な第 1 のアドレス「addr\_1」を表すために、それぞれインデックスおよびタグを作成するように動作する。インデックス生成器 1 2 0 4 は、addr\_1 からターゲットキャッシュ 1 2 3 0 内のメモリロケーションのサブセットにアクセスするための、インデックス「index(addr\_1)」を形成する。典型的には、index(addr\_1) の値は、選択される各 way の対応する各メモリロケーションを選択する。さらに、タグ生成器 1 2 0 6 はタグ「tag(addr\_1)」を形成するため、予測生成器 1 2 2 2 は、addr\_1 に関連付けられたターゲットキャッシュ 1 2 3 0 内の特定のトリガ - ターゲット関連にアクセスすることができる。

【 0 0 6 3 】

一例として、addr\_1 が「G」として考えてみる。このアドレスの場合、予測生成器 1 2 2 2 は、そのインデックスに関連付けられたメモリロケーションを選択するために index(G) を生成する。このインスタンスでは、index(G) は値(I)を有し、これは 3 である(すなわち I = 3)。これは、index(G) を使用して、wa

10

20

30

40

50

y ( 「 w a y 0 」 ) 1 2 4 0、 w a y ( 「 w a y 1 」 ) 1 2 5 0 から、 w a y ( 「 w a y N 」 ) 1 2 6 0 について、 I = 3 によって識別された各メモリロケーションを選択できることを意味し、ここでNは、ターゲットキャッシュ1230で使用可能なwayの数を表す構成可能な数である。同じアドレスGについて、タグ生成器1206は、Gに関連付けられた特定のメモリロケーションを識別するために、アドレスGのタグをtag ( G ) として作成することになる。したがって、index ( G ) のインデックスおよびtag ( G ) のタグが与えられた場合、図12に示されるように、ターゲットアドレスQおよびP (またはその代替表現) は、way 1240およびway 1250内のそれぞれのメモリロケーションから取り出すか、そこに格納することができる。特定の実施形態では、各アドレスは36ビットからなる。ビット28 : 18はアドレスのタグを表すことが可能であり、ビット19 : 9、18 : 8、17 : 7、またはビット16 : 6の任意のグループは、そのアドレスの構成可能なインデックスを表すことが可能である。一実施形態では、アドレスの一部がターゲットアドレスを交互に表す。たとえば、36ビットターゲットアドレスのビット30 : 6はターゲットキャッシュ1230のTRT列内で維持される。ターゲットおよびトリガアドレスの両方の表現が減少すると、必要なハードウェアが減少し、それによって材料、リソースなどに関するコストが削減される。

10

## 【0064】

ターゲット特定器1208は、トリガ - ターゲット関連が特定のトリガに対して存在するかどうかを特定し、存在する場合、そのトリガの各ターゲットアドレスを特定する。引き続き前の例を見ると、ターゲット特定器1208は、tag ( G ) が他のトリガアドレスを表すindex ( G ) のタグとマッチするのに応答して、ターゲットアドレスQおよびPを取り出す。当業者であれば、よく知られた比較回路 ( 図示せず ) が、マッチングタグを識別するために、予測生成器1222またはターゲットキャッシュ1230いずれかでの実施に好適であることを理解されよう。1つまたは複数のターゲットアドレスが見つかった場合、それらのアドレスがコンバイナ1210に渡される。コンバイナ1210は各ターゲットアドレス1214を、トリガアドレスのインデックスおよびタグからなる予測識別子 ( 「 P I D 」 ) 1212に関連付ける。PID 1212は、ターゲットアドレスQおよびPを予測させるトリガアドレスを識別する。したがって、PID 1212を [ index ( G ) , tag ( G ) ] として表すことができる場合、予測生成器1222によって生成される非順次予測は、基準として [ [ index ( G ) , tag ( G ) ] , Q ] の形を有することになる。予測としてのQは、 [ index ( G ) , tag ( G ) ] が関連付けられる場合、「基準予測」とみなされることに留意されたい。したがって、キャッシュメモリにプリフェッチされた予測情報は、data ( Q ) + [ [ index ( G ) , tag ( G ) ] , Q ] として表すことができる。

20

30

## 【0065】

コンバイナ1210は、トリガアドレスに対して非順次的ないくつかの追加予測を生成するための「バッチ」信号1226を受け取るように構成することができる。たとえば、バッチ信号1226がコンバイナ1210に対して、マッチしたターゲットアドレスを含む領域を有する予測のグループとして、「n」個の予測を生成するように指示すると考えてみる。そこで、トリガアドレス「G」がアドレス「Q」(すなわち、基本アドレスとしてのQ0)の非順次予測を生成した場合、予測アドレスはQ0、Q1、Q2、...、Qbを含むことが可能であり、ここでbはバッチ信号によって設定される数である。バックセクタまたはブラインドバック順次予測が同時に生成されるいくつかのケースでは、バッチbがb - 1に設定できることに留意されたい。したがって、予測アドレスのグループは、Q(-1)、Q0、Q1、Q2、...、Q(b-1)を含むことになる。予測アドレスのグループ内のそれぞれをPID 1212に関連付けることもできることに留意されたい。特定の実施形態では、ターゲットアドレス1214はトリガアドレスの属性を継承し、ここでこうした属性は、トリガアドレスがコードまたはプログラムデータに関連付けられているかどうか、およびトリガアドレスがプロセッサデマンドアドレスであるか否かを示す。他の特定の実施形態では、グループ内の予測アドレスよりも少ない数をPID

40

50

1212に関連付けることもできる。一例では、ターゲットアドレスQ0のみがPID1212に関連付けられ、グループ内の他の1つまたは複数(たとえば、Q(-1)、Q2、Q3など)をPID1212に関連付ける必要がない。したがって、トリガアドレスGに遭遇し、その後ターゲットアドレスQ0が続く場合、PID1212は非順次予測器に報告される。その後、Q2またはグループのうちの任意の他のアドレスに遭遇した場合、PID1212は報告されない。これにより、ターゲットキャッシュ内の冗長エントリの数が削減される。したがって、関連「G->Q0」のみが格納され、その予測のヒットの結果として再度優先順位付けされる。アドレスストリーム内でアドレスQ1が検出された場合、非順次予測器は関連「G->Q1」を挿入する必要がなくなる。

#### 【0066】

次に、ターゲット特定器1208が、addr\_1に関するターゲットアドレスを検出しないと考える。次にターゲット特定器1208は、addr\_1に関するトリガ-ターゲット関連が存在しない旨を挿入器1202に伝える。これに回答して挿入器1202は、addr\_1に関するトリガ-ターゲット関連を形成し、その関連をターゲットキャッシュ1230に挿入する。このように実行するために、挿入器1202は第1に、tag(addr\_1)を格納するために使用されるindex(addr\_1)を使用してメモリロケーションを識別する。挿入器1202は、トリガアドレスaddr\_1へのターゲットアドレスとして格納するために後続のアドレス「addr\_2」を受け取るようにも構成される。新しく形成されたトリガ-ターゲット関連に先立って存在するトリガ-ターゲット関連がない場合、挿入器1202はtag(addr\_1)およびaddr\_2を、最高優先順位のway(すなわちway\_0)であるway\_1240のTRG列およびTGT列にそれぞれ格納する。たとえば、図11のアドレスストリーム1104について考えてみると、このストリームは「Y」の後に「Z」が続く第1のインスタンスを示す。「tag(Y) to Z」のトリガ-ターゲット関連が存在しないことが特定された後、図12の挿入器1202はindex(Y)で新しいトリガ-ターゲット関連を格納する。したがって、「tag(Y) to Z」がトリガ-ターゲット関連1242としてway\_1240に格納される。特定の実施形態では、挿入器1202は優先順位調整器1324から挿入信号(「INS」)1224を受け取るが、これについて次に説明する。

#### 【0067】

図13は、本発明の実施形態に従った、例示的優先順位調整器1324を示す図である。一般に、優先順位調整器1324は、最も新しい正常なターゲットアドレスから最も古いまたは正常でないターゲットアドレスへと、トリガ-ターゲット関連を優先順位付けするように動作する。たとえば、項目に以前のターゲットが存在しない場合、トリガ-ターゲット関連には最高優先順位が割り当てられる(すなわち、way\_0に格納される)ことになる。さらに、予測されたターゲットアドレスが正常に提供された場合(たとえば、プロセッサによりデータが読み取られ、そのデータが非順次予測に基づいてプリフェッチされた場合)、トリガ-ターゲット関連に最高優先順位を割り当てることができる。この例では、優先順位調整器1324は、とりわけその中に格納されたトリガ-ターゲット関連を優先順位付けするために、ターゲットキャッシュ1230に結合される。優先順位調整器1324は、レジスタ1302、インデックス復号器1308、タグ復号器1310、ターゲット特定器1318、マッチャ1314、および優先順位変更器1316を含む。

#### 【0068】

一般に、優先順位調整器1324は、特定アドレスがプロセッサの要求したデータの提供に成功したことを示す、非順次予測器1010外部の情報を受け取る。こうした情報は、図10で説明したマルチレベルキャッシュ1040などのキャッシュメモリによって生成することができる。優先順位調整器1324は、この情報を「Hit Info」としてレジスタ1302内で受け取る。Hit Infoは、少なくともデータ(たとえば、プロセッサによって実際に要求されたプログラム命令および/またはプログラムデータ)

10

20

30

40

50

のアドレス1304を含む参照である。アドレス1304はaddr\_2としてラベル付けされている。この参照は、アドレス1304に関連付けられたPID 1306も含む。

#### 【0069】

インデックス復号器1308およびタグ復号器1310は、それぞれ、addr\_2が適切な優先順位レベルを有するかどうかを特定するために、PID 1306からindex(addr\_1)およびtag(addr\_1)を抽出する。このように実行するために、優先順位調整器1324は、addr\_2が、ターゲットキャッシュ1230内の既存のトリガ-ターゲット関連のターゲットアドレスであるかどうかを識別する。優先順位調整器1324がターゲットキャッシュ1230にtag(addr\_1)およびindex(addr\_1)を適用した後、ターゲットキャッシュ1230のTRG列にある任意のマッチングトリガアドレスがターゲット特定器1318によって受け取られることになる。addr\_1に関連付けられた1つまたは複数のターゲットアドレスを検出すると、ターゲット特定器1318はそれらのターゲットアドレスをマッチャ1314に提供する。

10

#### 【0070】

しかしながら、ターゲット特定器1318が、トリガ-ターゲット関連内にターゲットアドレスが存在しない(すなわち、アドレスaddr\_1に関連付けられたいかなるaddr\_2も存在しない)ことを特定した場合、新しいトリガ-ターゲット関連を挿入するために、図12の挿入器1202に挿入信号(「INS」)1224を送ることになる。挿入信号1224は、通常、addr\_1およびaddr\_2などのアドレス情報を含む。典型的には、Hit InfoのPID 1306に関するマッチングターゲットアドレスが存在しない状況は、プロセッサが以前に発行された非順次予測でヒットしたことを意味する。しかしながら、その後ターゲットキャッシュ1230は、その以前に発行された非順次予測に関する基準を形成したトリガ-ターゲット関連をパージした。したがって、非順次予測エンジン1010は、プロセッサによって正常に使用された非順次アドレスを予測するために再度使用可能なトリガ-ターゲット関連を挿入、または再挿入することになる。

20

#### 【0071】

ターゲット特定器1318は、1つまたは複数のターゲットアドレスを検出した場合、検出されたターゲットアドレスをマッチャ1314に提供する。マッチャ1314は、検出された各ターゲットアドレスを、addr\_2(すなわちアドレス1304)と比較して、addr\_1に関していくつの関連付けられたターゲットアドレスが存在するか、および、既存の各ターゲットアドレスに関して、対応するトリガ-ターゲット関連が常駐するwayを特定する。マッチャ1314は、その比較の結果を優先順位変更器1316に提供し、必要であれば優先順位を修正する。

30

#### 【0072】

第1に、1つまたは複数のターゲットアドレスが、トリガアドレスとしてaddr\_1を表すPID 1306(すなわちaddr\_1)に関連付けられているとして検出されているが、addr\_2を含むトリガ-ターゲット関連は存在しない、というインスタンスを考えてみる。したがって、優先順位変更器1316は、最高優先順位を表す位置(たとえばway 0)に新しいトリガ-ターゲット関連を挿入し、同じトリガの既存のトリガ-ターゲット関連の優先順位を降格することになる。たとえば、図12に示されるように、「tag(A)-to-X」トリガ-ターゲット関連が最高優先順位を表すメモリロケーションにあり、「tag(A)-to-L」関連がより低位の優先順位を有すると考えてみる。次に、PID 1306がアドレスAをaddr\_1として表し、addr\_2はアドレスBであると想定する。優先順位変更器1316は図13に示されるように、「tag(A)-to-B」関連をway 0に格納し、他の以前の関連はより優先順位の低い他のwayに格納されるように動作することになる。

40

#### 【0073】

50

第2に、2つのターゲットアドレスがPID 1306(すなわちaddr\_1)に関連付けられているとして検出されているが、2つのトリガ-ターゲット関連がそれらの優先順位を不適切にスワップされている、というインスタンスを考えてみる。このケースでは、優先順位変更器1316は、最高優先順位を表す位置(たとえばway 0)に最高優先順位トリガ-ターゲット関連を挿入し、第2位の優先順位を表す他の位置(たとえばway 1)に以前の最高優先順位トリガ-ターゲット関連を挿入することになる。たとえば、図12に示されるように、「tag(B)-to-G」トリガ-ターゲット関連が最高優先順位を表すメモリロケーションにあり、「tag(B)-to-C」関連がより低位の優先順位を有すると考えてみる。次に、PID 1306がアドレスBをaddr\_1として表し、アドレスCはaddr\_2であると想定する。優先順位変更器1316は図13に示されるように、「tag(B)-to-C」関連をway 0に格納し、他の関連はより優先順位の低いway 1に格納されるように動作することになる。この優先順位付けの技法は、少なくとも上位2つの優先順位が、それぞれ、最高優先順位および第2位優先順位としての「leg 0」および「leg 1」として維持される場合に有用であることに留意されたい。

10

## 【0074】

次に、2つのターゲットアドレスがPID 1306(すなわちaddr\_1)に関連付けられているとして検出されており、2つのトリガ-ターゲット関連にはそれらの優先順位が適切に割り当てられている、というインスタンスを考えてみる。このケースでは、優先順位変更器1316は、対応するトリガ-ターゲット関連が正しいとして、何のアクションも実行しない。

20

## 【0075】

図14は、本発明の特定の実施形態に従った、非順次予測を形成するために予測器生成器を動作させるための例示的パイプライン1400を示す図である。図14では、実線の四角形がステージ中またはステージ間のストレージを表し、破線の四角形が非順次予測器によって実行されるアクションを表す。ステージ0では、図13のインデックス復号器1308およびタグ復号器1310の融合とすることが可能なタグ/インデックス組み合わせ生成器1402によって、読み取り要求のaddr\_1が復号される。一実施形態では、タグ/インデックス組み合わせ生成器1402は、addr\_1をアドレスの第1部分およびアドレスの第2部分に分離するように構成されたマルチプレクサである。第1の部分はtag(addr\_1)として1406で保持され、第2の部分はindex(addr\_1)として1408で保持される。また、このステージでは、トリガ-ターゲット関連を記述するデータを取り出すために、index(addr\_1)が1410のターゲットキャッシュに印加される。オプションで、ターゲットキャッシュが書き込まれる間、読み取り要求のaddr\_1を一時的にバッファ1404に格納することができる。

30

## 【0076】

ステージ1では、tag(addr\_1)およびindex(addr\_1)がそれぞれ1412および1414で保持されたままである。1416で、ターゲットアドレスがターゲットキャッシュから読み取られる。ステージ2では、第1に1418でtag(addr\_1)とindex(addr\_1)に関連付けられたタグとをマッチングすることによって、非順次予測エンジンが好適な非順次予測を選択する。1420で、非順次予測エンジンは、たとえば最高優先順位のターゲットアドレスを(すなわち最高優先順位のトリガ-ターゲット関連を格納しているwayから)1422でleg 0予測キューへ転送するように、および、第2位優先順位のターゲットアドレスを(すなわち第2位優先順位のトリガ-ターゲット関連を格納しているwayから)1424でleg 1予測キューへ転送するように、マルチプレクサを構成する。ステージ3では、1430で、これら2つの非順次予測がたとえばコンパイナへと出力される。図14では非順次予測を4つのステージで生成するが、他の実施形態の他の非順次予測パイプラインは、これよりも多いかまたは少ないステージを有することができることに留意されたい。

40

## 【0077】

50

図15は、本発明の特定の実施形態に従った、非順次予測を優先順位付けするように優先順位調整器を動作させるための例示的パイプライン1500を示す図である。実線の四角形がステージ中またはステージ間のストレージを表し、破線の四角形が優先順位調整器によって実行可能なアクションを表す。パイプライン1500は、トリガ-ターゲット関連をターゲットキャッシュに挿入し、ターゲットキャッシュ関連の優先順位を変更する、例示的方法を示す。ステージ1は、優先順位調整器が挿入するか優先順位付けするかを決定する。優先順位調整器が挿入を実行しようとする場合、1502の読み取り要求のアドレス`addr_1`が、このステージ中に1506で格納される。このアドレスは、ターゲットアドレス用のトリガアドレスとなる可能性を有する。優先順位調整器が優先順位付けを実行しようとする場合、1504で、優先順位調整器は外部ソース(たとえばキャッシュメモリ)から`addr_1`アドレスを表すPID 1508を受け取り、このステージ中に1510でアドレス`addr_2`も受け取る。

【0078】

図14および15は、1つのレベルの予測を使用して非順次予測を例示することに留意されたい。マルチレベルの予測生成を達成するためには、それぞれのパイプライン1400および1500の終わりに、生成された予測を入力アドレスとしてパイプライン1400および1500にフィードバックするように、例示的パイプライン1400および1500を修正することができる。その後、これらの予測は、他のレベルの予測生成用にキューに入れられる。たとえばAが検出された場合、ターゲットキャッシュ1130はターゲットアドレスBおよびXを(たとえば、2つの最高優先順位`way`として)生成する。その後、連続するトリガアドレスとして、アドレスBがパイプラインのトップに再入力され、これによってターゲットキャッシュ1130はアドレスCおよびGを生成する。端的に言えば、複数レベルの予測を実施するために、例示的パイプライン1400および1500にフィードバックループを追加することができる。

【0079】

第1に、ステージ0で、優先順位調整器がトリガ-ターゲット関連挿入を実行していると考えてみる。このインスタンスでは、`addr_1`がタグ/インデックス組み合わせ生成器1514によって復号され、`addr_2`が1512からマルチプレクサ1516を介して選択される。タグ/インデックス組み合わせ生成器1514は、インデックス生成器およびタグ生成器の集合機能を実行する。一実施形態では、タグ/インデックス組み合わせ生成器1514は、1506または1508のいずれかからアドレスを選択するように構成されたマルチプレクサである。このケースでは、タグ/インデックス組み合わせ生成器1514は、1520で`tag(addr_1)`として保持される第1のアドレス部分を形成し、1522で`index(addr_1)`として保持される第2の部分形成する。また、このステージでは、`index(addr_1)`が、トリガ-ターゲット関連を記述するデータを取り出すために、マルチプレクサ1518を介して1524でターゲットキャッシュに印加される。次に、ステージ0で、優先順位調整器がターゲットキャッシュの優先順位付けを実行していると考えてみる。このインスタンスでは、`addr_1`(またはその代替表現)が1508から受け取られ、`addr_2`がマルチプレクサ1516を介して1510から選択される。その後、タグ/インデックス組み合わせ生成器1514は、PID 1508から第1および第2の部分形成する。その後、PID 1508から形成された`Index(addr_1)`が、トリガ-ターゲット関連を記述するデータを取り出すために、マルチプレクサ1518を介して、1524でターゲットキャッシュに印加される。ステージ1からステージ3では、優先順位調整器が挿入または優先順位付けを実行しているかどうかにかかわらず、パイプライン1500は同様に挙動する。

【0080】

ステージ1では、`tag(addr_1)`および`index(addr_1)`は、1530および1532でそれぞれ保持されたままである。1534では、ターゲットアドレスがターゲットキャッシュから読み取られる。ステージ2では、優先順位調整器が第1に

10

20

30

40

50



tag(addr\_1)をタグとマッチングさせる。1540でマッチするタグがない場合、1542でマルチプレクサがトリガ-ターゲット関連を挿入する準備をするように構成される。しかしながら、ターゲットキャッシュのwayからの少なくとも1つのタグが1544でマッチした場合、および最高優先順位のトリガ-ターゲット関連が最高優先順位に対応するwayに常駐していない場合、1554でトリガ-ターゲット関連の優先順位が変更される。これを実行するために、1552で、新しいトリガ-ターゲット関連を優先順位変更または挿入するためのマルチプレクサが選択される。ステージ3では、完全に接続された優先順位変更マルチプレクサが、1556からのaddr\_2を格納するように構成される。このアドレスは、1550で保持されたindex(addr\_1)によって特定された場合、ステージ0でway 0のターゲットアドレスとして書き込まれることになる。図に示されるように、完全に接続された優先順位変更マルチプレクサによって1560で特定された他のトリガ-ターゲット関連も、1550で保持されるindex(addr\_1)を使用して、キャッシュ書き込みデータとして1524でターゲットキャッシュに書き込まれる。パイプライン1500がステージ0に戻った後、優先順位調整器は適宜動作を続行する。

10

#### 【0081】

インベントリから予測を発行するための例示的实施形態

図16は、本発明の特定の实施形態に従った、例示的予測インベントリ1620を示すブロック図である。この例で、予測インベントリ1620は、プリフェッチャ1606内に常駐するように示される。さらにプリフェッチャ1606は、1つまたは複数のプロセッサによって少なくともメモリアクセスを制御するように設計された、メモリプロセッサ1604内で動作するように示される。プリフェッチャ1606は、必要とされる前にプログラム命令およびプログラムデータの両方をメモリ1612から「フェッチ」し、その後、プロセッサの要求時に、フェッチしたプログラム命令およびプログラムデータをそのプロセッサ1602に提供するように動作する。使用に先立ってフェッチすること(すなわち「プリフェッチ」)によって、プロセッサのアイドル時間(たとえばプロセッサ1602がデータ不足である間の時間)が最小化される。プリフェッチャ1606は、予測を生成するためのスペキュレータ1608および不必要な予測を除去するためのフィルタ1622も含む。

20

#### 【0082】

フィルタ1622は、インベントリフィルタまたはインベントリ後フィルタ、あるいはその両方の表現である。不必要な予測を除去することにより、プリフェッチャ1606は、そうでなければ重複する予測を不必要に管理するために使用されることになる計算リソースおよびメモリリソースを保存することができる。インベントリフィルタ(インベントリ前フィルタとしての)が、予測インベントリ1620への挿入に先立って不必要な予測を除去するように動作するのに対して、インベントリ後フィルタは、メモリ1612への発行に先立って不必要な予測を除去する。インベントリ後フィルタの一例は、図20に示されている。次に、プリフェッチャ1606の動作およびその構成要素について説明する。

30

#### 【0083】

動作時に、スペキュレータ1608は、メモリ1612にアクセスするためのプロセッサ1602による要求(「読み取り要求」)について、システムバス1603を監視する。プロセッサ1602がプログラム命令を実行する場合、スペキュレータ1608は、プロセッサ1602によってまだ使用されていないプログラム命令およびプログラムデータを含むアドレスに関する読み取り要求を検出する。考察のために、「アドレス」は、一般にメモリ1612とキャッシュメモリ(図示せず)との間で転送されるメモリのキャッシュラインまたは単位に関連付けられる。キャッシュメモリとは、予測インベントリの外部にある予測のリポジトリの一例である。キャッシュラインの「アドレス」はメモリロケーションを表すことが可能であり、キャッシュラインはメモリ1612の複数のアドレスからのデータを含むことができる。「データ」という用語は、プリフェッチ可能な情報の単

40

50

位を表すのに対して、「プログラム命令」および「プログラムデータ」という用語は、それぞれ、プロセッサ1602によってその処理中に使用される命令およびデータを表す。したがって、データ(たとえば任意のビット数)は、プログラム命令および/またはプログラムデータを構成する予測情報を表すことができる。

#### 【0084】

検出された読み取り要求に基づいて、スペキュレータ1608は、プロセッサ1602によるメモリ1612へのアクセスを正確に予測する機会を改善するために多数の予測を生成することが可能であり、それらの多数の予測は冗長予測を含む可能性がある。こうした予測の例には、前方順次予測、逆順次予測、バックブラインド順次予測、バックセクタ順次予測、非順次予測などが含まれる。こうした冗長を除去するために、インベントリフィルタ1622は重複予測をフィルタリング除去して、存続している予測を生成し、その後それらが予測インベントリ1620に格納される。冗長を除去するために、インベントリフィルタ1622は生成された予測をキャッシュ(図示せず)のコンテンツと比較した後、それら予測を予測インベントリ1620に挿入する。予測と予測インベントリ1620内に常駐する予測との間にマッチが見つかった場合、インベントリフィルタ1622はその予測を無効にする。しかしながらマッチが見つからない場合、インベントリフィルタ1622は存続している予測を予測インベントリ1620に挿入する。新しい予測グループ(すなわち、1つのイベントまたは同じトリガアドレスによって生成された予測)内のいくつかの予測はキャッシュのコンテンツとマッチするが、他の予測はマッチしないケースがあることに留意されたい。この場合、インベントリフィルタ1622は、キャッシュ内の予測とマッチする個々の予測を無効にし、マッチしなかった(たとえば「無効」とマーク付けされていない)予測を予測インベントリ1620に挿入する。

#### 【0085】

いったん予測インベントリ1620に入ると、予測はインベントリの「アイテム」として維持される。「アイテム」という用語は、予測インベントリ1620内に格納される「予測」または「トリガアドレス」(予測を生成する)のいずれかを表す。これらのアイテムは、フィルタリングの目的で後に生成される予測と比較することができる。プリフェッチャ1606は、インベントリ内のこれらのアイテムを様々なレートでメモリ1612に発行しながら、管理する。発行のレートは予測のタイプ(たとえば、前方順次予測、非順次予測など)、各予測のタイプの優先順位、および以下で説明する他の要因に依存する。

#### 【0086】

予測が冗長になる可能性のある場合の1つが、プロセッサ1602が特定のアドレスに関する実際の読み取り要求を発行し、そのアドレスに関する予測が予測インベントリ1620内にすでに存在する場合である。このケースでは、予測はフィルタリング除去(すなわち無効化)され、プロセッサ1602の実際の読み取り要求が維持される。これは、順次タイプおよびバックタイプの予測などの場合、特にあてはまる。また、いくつかの予測は、予測インベントリ1620がそれらの予測を受け取り、プリフェッチャ1606がそれらをメモリ1612に発行するまでの時間に冗長となり、プリフェッチャ1606はアイテムを発行するのに先立って予測をフィルタリング除去することもできる。これにより、重複時間中に発生する冗長予測の数が再度削減されるが、後に生成される予測は予測インベントリ1620内に挿入される。また、冗長予測の数が減少すればするほど、保存されるリソースは多くなる。

#### 【0087】

プリフェッチャ1606が予測インベントリ1620から予測を発行した後、メモリプロセッサ1604は、メモリバス1611を介して残りの(少なくともインベントリ後フィルタによってフィルタリング除去されなかった)予測をメモリ1612に移送する。これに回答して、メモリ1612は、予測アドレスを参照しながらプリフェッチされたデータを戻す。プリフェッチャ1606内に常駐するかまたは常駐しないキャッシュメモリ(図示せず)は、メモリプロセッサ1604がそのデータをプロセッサ1602に送信するまでなどの間、戻されたデータを一時的に格納する。適切な時点で、メモリプロセッサ1

604は、プリフェッチされたデータを、とりわけ待ち時間が最小になるのを保証するために、システムバス1603を介してプロセッサ1602に移送する。

【0088】

図17は、本発明の一実施形態に従った、例示的予測インベントリ1620を示す図である。予測インベントリ1620は、予測を格納するためのいくつかのキュー1710、1712、1714、および1716を含み、キューは、それぞれの予測が発行またはフィルタリング除去されるまで格納するためのバッファまたは任意の同様の構成要素とすることができる。予測インベントリ1620は、インベントリマネージャ1704および1つまたは複数のキュー属性1706も含み、これによってインベントリマネージャ1704は、対応するキュー属性1706に従ってキューそれぞれの構造および/または動作を構成する。

10

【0089】

個々のキューは、予測をアイテムとして維持し、そのすべてが一般に、前方順次予測などの同じ特定タイプの予測のものである。図に示されるように、予測インベントリ1620は4つのキュー、すなわち順次キュー(「Sキュー」)1710、バックキュー(「Bキュー」)1712、非順次ゼロキュー(「NS0キュー」)1714、および非順次1キュー(「NS1キュー」)1716を含む。順次キュー1710は、前方順次予測または逆順次予測のいずれかを含むように構成可能であり、バックキュー1712はブラインドバック順次予測またはバックセクタ順次予測のいずれかを含むことができる。考察のために、前方順次予測、逆順次予測などは、まとめて「シリーズタイプ」予測と呼ぶことが可能であり、ブラインドバック順次予測、バックセクタ順次予測などは、まとめて「バックタイプ」予測と呼ぶことが可能であることに留意されたい。

20

【0090】

予測インベントリ1620は、「0番目」の非順次キューおよび「1番目」の非順次キューを含む。非順次(「0」)キュー1714および非順次(「1」)キュー1716は、それぞれ「最高」および「第2位」の優先順位を有する非順次予測を含む。特に、非順次0キュー1714は、対応するトリガアドレスによって生成可能な(任意数のターゲットアドレスのうちの)最高優先順位のターゲットアドレスを含む、非順次予測を維持する。「トリガ」アドレスとは、スペキュレータ1608の予測生成元となる検出されたアドレスのことである。こうした予測(すなわち予測アドレス)は、ターゲットを生成するトリガでパターン化不可能(たとえば非順次)な、「ターゲット」アドレスである。同様に、非順次1キュー1716は非順次予測を維持しないが、代わりに対応するトリガアドレスによって生成可能な第2位の優先順位のターゲットアドレスを含む。

30

【0091】

各キューは、グループ0、1、2、および3などの任意数のグループ1720からなるものとしてすることができる。各グループ1720は、トリガアドレス、およびトリガアドレスが生成する対応する予測などの、構成可能な数のアイテムを含む。たとえば、順次キュー1710のグループ1720は、それぞれトリガアドレスおよび7つの順次予測を含むことが可能であり、バックキュー1712のグループ1720は、それぞれトリガアドレスおよび1つのバックタイプ予測を含む(または、場合によってはこれらのキューはアイテムとしての予測のみを含む)ことが可能である。さらに、非順次0キュー1714または非順次1キュー1716、あるいはその両方の、グループ1720は、トリガアドレスと、4つの非順次予測のグループとを含む(または、場合によってはアイテムとしての予測のみを含む)ことが可能である。特定の実施形態では、スペキュレータ1608は、特定数の予測を生成するためにその「バッチ」数を設定することによって、予測インベントリ1620に格納されたグループ1720あたりのアイテム数を決定する。予測をグループ化されたアイテムとして予測インベントリ1620に格納することによって、グループ1720は、通常は各予測を個別に管理するために使用される情報の量を削減し、これによって予測を発行する場合のアービトレーションが容易になる。

40

【0092】

50

インベントリマネージャ 1704 は、各キュー内のアイテムのインベントリを管理するように、ならびにキューの構造および/または動作を制御するように、構成される。予測インベントリ 1620 を管理するために、インベントリマネージャ 1704 は全体としてまたは部分的に、1つまたは複数のキュー属性 1706 を使用してこれを実行する。キュー属性の第 1 の例は、キューのタイプである。たとえば、キュー 1710 から 1716 のうちのいずれかを、先入れ先出し(「FIFO」)バッファ、後入れ先出し(「LIFO」)バッファ、または任意の他のタイプのバッファとなるように構成可能である。FIFO または LIFO などのキューのタイプは、キューに関してアイテムを挿入および除去する方法に影響を与える。一実施形態では、順次キュー 1710 は LIFO として構成され、非順次 0 キュー 1714 および非順次 1 キュー 1716 はそれぞれ FIFO として構成される。

10

#### 【0093】

キュー属性の第 2 の例は、キュー、グループ、またはアイテムに割り当て可能な満了時間または存続時間である。この属性は、予測に関する陳腐化の程度を制御する。任意のグループ 1720 またはキュー内の予測が古くなる、または陳腐になるにつれて、正確な予測を反映する可能性が次第に少なくなる。したがって、古くなったアイテムを最小にするために、インベントリマネージャ 1704 は、グループがその現在のインベントリをある満了時間まで維持できるようにし、その時間が過ぎると、インベントリマネージャ 1704 は古くなったグループ全体またはまだ発行されていない任意の残りのアイテムのいずれかをパージする。本発明の一実施形態では、キュー、グループ、またはアイテムの存続時間を、それらが無制限に保持するように構成することができる。すなわち、それらを「不滅」として設定することが可能であり、これは、不滅の予測が発行されるまでまたは不滅性が撤回されるまで、キュー内に常駐することを意味する。特定の实施形態では、グループがキューに挿入される場合、満了時間はそのグループに関連付けられる。その後、ゼロに達した場合、そのグループの残りのアイテムが無効化されるように、タイマが満了時間からカウントダウンする。他の実施形態では、非順次予測が発行され、その結果としてデータキャッシュ内でヒットすることになる確率を上げるために、非順次 0 キュー 1714 または非順次 1 キュー 1716 のいずれかのグループ 1720 に関する満了時間が、順次キュー 1710 のグループ 1720 よりも長く設定される。

20

#### 【0094】

キュー属性の第 3 の例は、キューが満杯の場合、インベントリマネージャ 1704 が予測をどのようにキューに挿入するかを示すために、キューに関連付けられた挿入インジケータである。1つのインスタンスでは、挿入インジケータは、インベントリマネージャ 1704 が新しく生成された予測を挿入されないようにするか、または特定のキューに常駐する古いアイテムを上書きするかを示す。挿入インジケータが「ドロップ」された場合、インベントリマネージャ 1704 は、そうでなければ挿入されることになる任意の新しい予測を廃棄する。挿入インジケータが「上書き」された場合、インベントリマネージャ 1704 は特定のキューが対応するキューのタイプに応じて、2通りのアクションのうちの 1つを実行する。キューが LIFO として構成された場合、インベントリマネージャ 1704 は事実上新しい予測をスタックとして LIFO に押し入れ、これが最も古いアイテム

30

40

#### 【0095】

キュー属性の第 4 の例は、次のアイテムの発行元となる特定のキューを決定するために、それぞれのキューに関連付けられた優先順位である。一実施形態では、優先順位は、次の予測を選択するためにキュー間でアービトレーションするためのキュー 1710、1712、1714、および 1716 のそれぞれに関して設定される。シリーズタイプの予測がより多量に生成される応用例では、順次キュー 1710 を処理することが重要である。したがって、このキューは通常比較的高い優先順位に関連付けられる。たとえば、これは非順次 0 キュー(「NS0 キュー」) 1714 および非順次 1 キュー(「NS1 キュー」

50

） 1716 が、順次キュー 1710 に比べて低い優先順位に設定される確率が最も高いことを意味する。キュー属性の他の例は、どれだけの予測を一時的に格納できるかを特定するために各キューに関連付けられるキューサイズである。たとえば、順次キューは 2 つのグループのサイズまたは深さを有することが可能であり、バックキューは 1 つのグループの深さを有することが可能であり、非順次キューは 4 つのグループの深さを有することが可能である。キューサイズは、異なるタイプの予測にどれだけのインベントリメモリが割り当てられるかを制御することによって、プリフェッチャ 1606 によって発行される予測の数を制御することができることに留意されたい。

#### 【0096】

本発明の一実施形態によれば、バックキュー 1712 の優先順位は、順次キュー 1710 のそれよりも高くなるように、動的に促進または修正することが可能である。この特徴は、スペキュレータ 1608 が上位または「フロント」セクタを検出した後に、メモリ 1612 から予測情報を取り出す際のものである。これは、プロセッサ 1602 が、キャッシュラインの上位またはフロントセクタを要求した直後に、より下位または「バック」セクタを要求する可能性が高いためである。したがって、バックキュー 1712 の優先順位を上げることによって、特に、バックセクタ順次予測を維持している場合、プリフェッチャ 1606 が適切なバックセクタ順次予測をメモリ 1612 に発行することになる確率が増加する。特定の実施形態では、バックキューカウンタ（図示せず）が、バックキュー 1712 以外のキューから発行されるアイテムの数をカウントする。このカウンタがしきい値に達すると、バックキュー 1712 は少なくとも順次キュー 1710 よりも高い優先順位へと促進される。その後、アイテム（たとえばバックセクタアイテム）をバックキュー 1712 から発行することができる。少なくとも 1 つのバックタイプアイテムを発行するかまたはバックキュー 1712 が（たとえば、古くなることまたはすべてのアイテムを発行することによって）空になった後、バックキュー 1712 の優先順位はその初期の優先順位に戻り（または逆戻りし）、バックキューカウンタがリセットされる。

#### 【0097】

一般に、予測の非順次グループのうちの任意のグループ 1720 について、非順次予測に関するターゲットアドレスとして、シリーズタイプおよびバックタイプの予測の混合が存在する可能性がある。特に、非順次アドレスのグループは、シリーズタイプ（すなわち、前方または逆のいずれか）の予測のみを含むことができる。しかしながらこれらのグループは、バックタイプと混合されたいくつかのシリーズタイプの予測を含むこともできる。前者の一例として、スペキュレータ 1608 が、トリガアドレス「A」がターゲットアドレス「B」および他のターゲットアドレス「C」に関連付けられていることを特定すると考えてみる。ターゲットアドレス B が C より高い優先順位の場合、B は非順次 0 キュー 1714 内に維持されると共に、予測のグループはトリガアドレス A に対して非順次的である。その後、グループは予測 B0（すなわちアドレス B）、B1、B2、および B3 を含むことが可能であり、そのすべてがアドレス A に対して非順次的であるが、すべて前方シリーズタイプである。後者の一例として、グループ 1720 は非順次予測 B（-1）（すなわちアドレス B-1）、B0、B1、および B2 を含むことが可能であり、ここで予測 B（-1）は他のシリーズタイプ予測と混合されたバックタイプ予測である。あるいは、グループ 1720 は、本明細書では具体的に説明しない予測の任意の他の配置構成を含むことができる。C は B よりも 2 番目に高位の優先順位を有するため、C は非順次予測の同様のグループと共に非順次 1 キュー 1716 内で維持される。したがって、予測 B0、B1、B2、および B3 を非順次 0 キュー 1714 のグループ 3 として挿入することが可能であり、予測 C0、C1、C2、および C3 を非順次 1 キュー 1716 のグループ 3 として挿入することができる。

#### 【0098】

図 17 は、一実施形態で、予測インベントリ 1620 が予測パスを存続させるインベントリフィルタ 1702 を介して予測 1701 を受け取るように構成されることも示す。その後、存続する予測が適切なキューに挿入され、前述のようにインベントリマネージャ 1

10

20

30

40

50

704によって管理される。次に、例示的なインベントリフィルタ1702について説明する。

【0099】

図18は、本発明の特定の実施形態に従った、インベントリフィルタ1702の例を示す図である。この例は、図17の順次キュー1710などの、順次キューに対する前方順次予測のフィルタリングに適用されるが、任意のタイプの予測をフィルタリングするために任意のキューと協働してインベントリフィルタ1702を使用することができる。すなわち、インベントリフィルタ1702は、任意の予測タイプの任意数の予測を、異なる予測タイプの予測を含む少なくとも1つの他のキューと比較するように構成することができる。たとえば、いくつかの前方順次予測をバックキューなどに対してフィルタリングすることができる。インベントリフィルタ1702は、グループ内のアイテム1806およびいくつかの予測1802をマッチングするために少なくともマッチャ1804を含む。グループ1806は、アイテムA1からA7を含み、そのそれぞれがアイテムA0に関連付けられる。A0は、以前にアイテムA1からA7として識別された予測を生成したトリガアドレスである。また、グループ1806は、順次キュー1710内の任意のグループ1720として常駐することができる。予測1802の数に関しては、これらはトリガアドレスとしての「TA」および予測B1からB7を含み、そのすべてがアドレスTAの検出時にスペキュレータ1608によって生成されたものである。図18は1つのグループ(すなわちグループ1806)のみを示しているが、同じキューの他のグループ1720が同じ方法および同時にフィルタリング可能であることに留意されたい。

【0100】

特定の実施形態では、マッチャ1804は、CMP0、CMP1、CMP2、...、CMPM(図示せず)として識別されるいくつかのコンパレータからなる。コンパレータCMP0は、TAとグループ1806内のN個のアイテムとを比較し、コンパレータCMP1、CMP2、...、CMPMはそれぞれ予測1802からの予測とグループ1806内のN個のアイテムのうちいくつかとを比較するように構成され、ここでMは生成された最大数の予測を収容するように設定される。一例として、Mは7であり、それによって7つのコンパレータが必要であり、Nは3であり、その結果各コンパレータは1802内の1つの要素と1806内の3つのアイテムとを比較すると考える。さらに、予測1802の各要素が、同じ位置を有する対応するアイテムとマッチングされる(たとえば、1番目と1番目、2番目と2番目など)と考える。したがって、CMP0はTAとA0、アイテムA1、およびアイテムA2とを比較し、CMP1は予測B1とアイテムA1、A2、およびA3とを比較するという具合である。数字Nは、コンパレータハードウェアの量を最小にするように、ただし、連続するストリームと、システムバス1603上で検出されたアドレスのストリーム内の小さな(すなわちNより大きくない)ジャンプから生じる可能性のあるそれらの予測とを十分にフィルタリング除去するように、設定可能である。

【0101】

一実施形態では、キューは、A0を表すためのページアドレスを格納し、アイテムA1、アイテムA2などを表すそれぞれをオフセットする。このケースでマッチが存在するかどうかを判別するために、アドレスTAのページアドレスおよび予測1802からの特定の予測のオフセットは、それぞれA0のページアドレスおよび対応するオフセットと比較される。本発明の特定の実施形態では、インベントリフィルタ1702は非順次予測に対して順次予測をフィルタリングしないため、非順次0キュー1714または非順次1キュー1716とは協働しない。これは、順次予測に存在するほど多くの冗長が非順次スペキュレーションにはない可能性が高いためである。

【0102】

図19Aおよび19Bは、本発明の特定の実施形態に従って冗長をフィルタリング除去する例示的技法を示す図である。マッチャ1804がマッチを特定すると、新しく生成された予測(すなわち、新しいアイテムK)または以前に生成されたアイテム(すなわち古いアイテムK)のいずれかが無効化される。図19Aは、新しいアイテムKまたは古いA

アイテムKのいずれかがフィルタリング除去または無効化されるのかを示す。このケースでは、キュー1902はFIFOである。したがって、新しいアイテムKが無効化されることになり、これによって古いアイテムKは維持される。これに対して、図19Bは、キュー1904がLIFOの場合、古いアイテムKが無効化されることになり、これによって新しいアイテムKが維持されることを示す。一般に、新しいアイテムKまたは古いアイテムKのいずれかのうち、最も早く発行したものが維持され、最も新しく発行したものは無効化されることになる。当業者であれば、インベントリフィルタ1702は、本発明の範囲および精神を逸脱しない他の技法が採用できることを理解されたい。

#### 【0103】

図20は、本発明の一実施形態に従った、プリフェッチャ内に配置される他の例示的予測インベントリを示す図である。この例では、プリフェッチャ2000はスペキュレータ1608およびフィルタ2014を含む。図20のプリフェッチャ2000は、マルチレベルキャッシュ2020および予測インベントリ1620も含む。ここでマルチレベルキャッシュ2020は、第1レベルの戻りデータキャッシュ(「DRC1」)2022および第2レベルの戻りデータキャッシュ(「DRC2」)2024からなる。第1レベルの戻りデータキャッシュ2022は一般に短期データストアとして説明し、第2レベルの戻りデータキャッシュ2024は一般に長期データストアとして説明することができる。マルチレベルキャッシュ2020は、メモリ1612からプリフェッチされたプログラム命令およびプログラムデータを、プロセッサ1602が必要とするまで格納する。マルチレベルキャッシュ2020のキャッシュは、新しく生成された予測がマルチレベルキャッシュ2020に対してフィルタリング可能なように、プリフェッチされた予測情報を生成した予測への参照も格納する。たとえば、DRC1 2022およびDRC2 2024は、キャッシュラインまたはメモリの単位に関するデータに加えて、(1)新しい予測に対するフィルタリングに使用される格納済みキャッシュラインに関するアドレス、および(2)キャッシュラインが予測の結果としてキャッシュに入れられた場合のトリガアドレス、という2つのタイプの情報を、参照として格納する。特に、トリガアドレスは、スペキュレータ1608内の非順次予測の優先順位をシャッフルするために使用される。

#### 【0104】

予測インベントリ1620は、生成された予測に対して、アービタ2018によって選択されるまでの一時ストレージを提供する。予測インベントリ1620内の格納済み予測は、そうでなければ発行されることになる冗長をフィルタリング除去するために使用される。アービタ2018は、アービトレーション規則に従って、命令およびデータをプリフェッチするためにどの生成された予測を発行するかを決定するように構成される。一般に、こうしたアービトレーション規則は、予測を発行するために特定のキューを選択する際の基準を提供する。たとえばアービタ2018は、部分的または全体として、キューおよび/またはグループ間の相対的優先順位に基づいて、予測を選択および発行する。

#### 【0105】

フィルタ2014は、キャッシュフィルタ2010およびインベントリフィルタ1702の、少なくとも2つのフィルタを含む。キャッシュフィルタ2010は、新しく生成された予測と、プリフェッチされた命令およびデータをすでにマルチレベルキャッシュ2020内に格納させた以前の予測とを、比較するように構成される。したがって、1つまたは複数の新しく生成された予測が、マルチレベルキャッシュ2020に関して任意の以前に生成された予測と重複する場合、処理を必要とする予測の数を最小にするために冗長予測は無効とされる。冗長予測(すなわち、余分な不必要な予測)は新しく生成された予測の可能性もあることに留意されたい。インベントリフィルタ1702は、新しく生成された予測と、すでに生成されて予測インベントリ1620内に格納された予測とを比較するように構成される。一実施形態では、インベントリフィルタ1702は、図18に示された構造および/または機能と同様である。ここでも、1つまたは複数の新しく生成された予測が、予測インベントリ1620内に以前に格納された予測と重複する場合、プリフェッチャリソースを解放するために任意の冗長予測が無効化することができる。

10

20

30

40

50

## 【0106】

冗長予測の数をさらに減少させるために、プリフェッチャ2000内にインベントリ後フィルタ2016が含まれる。プリフェッチャ1606が予測インベントリ1620から予測を発行した後、または発行する直前に、インベントリ後フィルタ2016は、予測インベントリ1620がそれらの予測を最初に受け取った時間からアービタ2018が発行する予測を選択する時間までの間に発生した、冗長予測をフィルタリング除去する。これらの冗長は、通常、予測インベントリ内のアイテムの同じ予測アドレスを表す予測が、予測インベントリ1620からメモリに発行された可能性があるが、いずれの予測情報もまだキャッシュ2020に戻されていない(すなわち、フィルタリングの対象となる参照がキャッシュ2020内にはない)可能性があることによって発生する。一実施形態では、インベントリ後フィルタ2016は、図18に示されたインベントリフィルタ1702またはキャッシュフィルタ2002のいずれかと同じ構造および/または機能とすることができる。

10

## 【0107】

一実施形態では、インベントリ後フィルタ2016は、予測インベントリ1620内の各グループ1720の各アイテムに関する発行情報を維持する。特にこの発行情報は、特定グループのうちどのアイテムが発行されるかを示す。しかしながら、インベントリ後フィルタ2016は、予測インベントリ1620から発行されたアイテムを除去しない。むしろ、入ってくる冗長予測をフィルタリング除去する場合に比較の対象とすることができるため、それらを残しておく。その特定グループ内の各アイテムを発行する場合、発行情報はこれを反映するように更新される。すべてのアイテムが発行されるとグループはページされ、追加のアイテムを受け入れるためにキューが解放される。

20

## 【0108】

一実施形態では、アービタ2018は、予測インベントリ1620の予測の発行に関する何らかの側面を制御することができる。特にアービタ2018は、最も有利な予測を発行するために、キュー、グループ、またはアイテム間の相対的優先順位を修正することができる。特定の実施形態では、アービタ2018は、メモリ1612、キャッシュメモリ2020、またはメモリサブシステムの他の構成要素などの、メモリに過度に負荷をかける(すなわち、メモリの過剰利用)多数の予測の生成を抑制するために、相対的優先順位を効果的に修正するように構成される。たとえばアービタ2018は、構成可能な負荷しきい値を各キューに割り当てることができる。このしきい値は、特定のキューが予測を発行できる最高レートを示す。この負荷しきい値と、メモリ1612から要求された累積作業単位を維持する作業負荷アキュムレータ(図示せず)のコンテンツとが比較される。作業単位とは読み取り、書き込みなどの、メモリ1612の任意の要求アクションのことである。メモリ1612の追加の作業単位が要求されると、作業負荷アキュムレータ内の値が増加する。しかしながら、時間が経過するにつれて(たとえば、あらゆる一定数のクロックサイクルについて)、その値は減少する。動作時に、アービタ2018は各キューの負荷しきい値と作業負荷アキュムレータの値とを比較する。作業負荷値が負荷しきい値を超えると、アービタ2018は2つの例示的アクションのうちの一つを実行する。アービタ2018は、その特定のキューに関する予測の入手を停止するように、予測インベントリ1620に指示することが可能であり、その結果、その中のアイテムが発行されるかまたは古くなる。あるいはアービタ2018は、アイテムを上書きすることによって、キューのアイテムを取り除くことができる。作業負荷値が負荷しきい値よりも低くなったことをアービタ2018が検出すると、予測の発行にキューを再度使用できるようになる。

30

40

## 【0109】

キャッシュメモリ内の予測情報に関して先読みルックアップを実行するための例示的実施形態

図21は、本発明の特定の実施形態に従った、例示的マルチレベルキャッシュ2120を含むプリフェッチャ2100を示すブロック図である。この例で、マルチレベルキャッシュ2120は、キャッシュフィルタ2110、第1レベルの戻りデータキャッシュ(「

50



DRC1) 2122、および第2レベルの戻りデータキャッシュ(DRC2) 2124を含む。キャッシュフィルタ2110は、第1レベルのDRC 2122および第2レベルのDRC 2124の両方について、迅速に検査するかまたは「先読みルックアップ」を実行し、どちらのキャッシュにおいても予測アドレスなどの入力アドレスの有無を検出するように構成される。先読みルックアップとは、たとえばマルチレベルキャッシュ2120内にいくつかの予測がすでに存在するかどうかを並行して特定するための、キャッシュメモリの検査である。

#### 【0110】

予測の有無に応じて、マルチレベルキャッシュ2120は、以下で説明するキャッシュポリシー、例に従って、第1レベルのDRC 2122および第2レベルのDRC 2124の両方のコンテンツを管理する。第1レベルのDRC 2122は一般に短期データストアとして説明し、第2レベルのDRC 2124は一般に長期データストアとして説明することが可能であり、これによって第1レベルのDRC 2122内の予測は、プロセッサがそれらの予測を要求しない場合、最終的に第2レベルのDRC 2124に移行される。本発明の実施形態によれば、第1レベルのDRC 2122または第2レベルのDRC 2124のいずれか、あるいは両方が、予測アドレスならびにプロセッサ要求アドレスに基づいて、プリフェッチされたプログラム命令およびプログラムデータを格納することができる。また、キャッシュフィルタ2110、第1レベルのDRC 2122、および第2レベルのDRC 2124は、冗長予測を減少させることによって、ならびに、たとえば予測情報のプリフェッチ速度を上げること(たとえば、ページオープン動作を予想すること)によって、プリフェッチされたプログラム命令およびプログラムデータを提供する待ち時間を削減するように協働する。以下の考察はマルチレベルキャッシュメモリ(すなわち複数キャッシュ)に関するが、以下のいずれの例示の実施形態も単一のキャッシュメモリを含むことができることに留意されたい。

#### 【0111】

キャッシュフィルタ2110は、入力アドレスの範囲と、複数キャッシュが本来階層状である、いくつかの複数キャッシュのそれぞれとを、並行して比較するように構成される。たとえば、第1キャッシュの方がサイズが小さく、予測を比較的短期間格納するように適合されるのに対し、第2のキャッシュはサイズが大きく、予測を第1のキャッシュよりも長い期間格納するように適合される。本発明の一実施形態によれば、さらに第2のキャッシュは、その予測アドレスおよび対応する予測データを第1のキャッシュのみから受け取る。両方のキャッシュを並行して検査するために、特に第2のキャッシュの方が第1のキャッシュより大きい場合、キャッシュフィルタは、キャッシュ内で「ルックアップ」または検査された各アドレスの2つの表現を生成する。第1の表現が第1のキャッシュに使用され、第2の表現が第2のキャッシュに使用される場合、両方のキャッシュを同時に検査することができる。その理由の1つは、小さいキャッシュよりも大きいキャッシュの方が検査を必要とするアドレスおよびエントリが多いことである可能性がある。したがって、両方とも同時に検査される場合、大きい方のキャッシュのアドレスを検査するためには、小さい方のキャッシュよりも効率の良い技法が必要である。以下で論じる照会インターフェースが、これらの機能を実行する。

#### 【0112】

プリフェッチャ2100は、予測を生成するためのスペキュレータ2108も含む。具体的に言えば、スペキュレータ2108は、前方順次予測、逆順次予測、バックライン順次予測、バックセクタ順次予測などの順次予測を生成するための、順次予測器(「SEQ. 予測器」) 2102を含む。またスペキュレータ2108は、非順次予測を形成するための非順次予測器(「NONSEQ. 予測器」) 2104も含む。プリフェッチャ2100は、これらの予測を使用して、メモリ(図示せず)からのプログラム命令およびプログラムデータの両方を「フェッチ」し、その後、プロセッサ(図示せず)が命令またはデータを要求する前に、フェッチされたプログラム命令およびプログラムデータをマルチレベルキャッシュ2120内に格納する。使用に先立ってそれらをフェッチすること(す

なわち「プリフェッチ」)によって、プロセッサのアイドル時間(たとえば、プロセッサがデータ不足である間の時間)が最小化される。

【0113】

非順次予測器2104は、先行アドレスから、それぞれを非順次予測としてみなすことが可能な1つまたは複数の潜在的な非順次アドレスへの関連を格納するためのリポジトリとして、ターゲットキャッシュ(図示せず)を含む。ターゲットキャッシュは、迅速な方法で非順次予測を生成するために、そのコンテンツと、入ってくる検出されたアドレスとを比較するように設計され、それによってターゲットキャッシュは、たとえばマルチレベルキャッシュ2120におけるヒットにตอบสนองして、格納済みの非順次予測を優先順位付けするように構成される。具体的に言えば、マルチレベルキャッシュ2120がその要求に応じてプロセッサに予測アドレスを提供する場合、そのアドレスが属する格納済みのトリガ-ターゲット関連の優先順位が上がる。「トリガ」アドレスとは、非順次予測器2104が非順次予測を生成する元となる検出されたアドレスのことであり、結果として、その2つの間に形成されるパターン化不可能な関連の「ターゲット」と呼ばれる予測が生じる。トリガアドレスは、ターゲットアドレスと呼ぶことも可能な、順次予測を生じさせるアドレスと呼ぶことも可能であることに留意されたい。

10

【0114】

プリフェッチャ2100は、フィルタ2114、オプションの予測インベントリ2116、オプションのインベントリ後フィルタ2117、およびオプションのアービタ2118も含む。ここでフィルタ2114は、生成された予測と、予測インベントリ2116内に常駐する以前に生成された予測とを比較するための、インベントリフィルタ(図示せず)を含むように構成することができる。予測インベントリ2116は、アービタ2118がメモリにアクセスするための予測を選択するまで、生成された予測を格納するための一時ストレージを提供する。アービタ2118は、命令およびデータをプリフェッチする場合に、生成された予測のうちどの予測がメモリにアクセスするために発行されるかを決定するように構成される。いくつかの実施形態では、フィルタ2114は、生成された予測と、プログラム命令およびプログラムデータをすでにマルチレベルキャッシュ2120内に「プリフェッチ」された状態にした以前に生成された予測とを比較するように構成可能な、キャッシュフィルタ2110を含むことが可能である。したがって、任意の生成された予測が、マルチレベルキャッシュ2120に格納された任意の以前に生成された予測と重複する場合、管理が必要な予測の数を最小にするようにその冗長予測は無効(または無効化)とすることが可能であり、それによってプリフェッチャリソースが解放される。

20

30

【0115】

動作時に、プロセッサがメモリへのアクセスを要求した(読み取り要求)場合、スペキュレータ2108はシステムバスを監視する。プロセッサがプログラム命令を実行すると、スペキュレータ2108は、プロセッサによってまだ使用されていないプログラム命令およびプログラムデータを含むアドレスに関して、読み取り要求を検出する。考察のために、「アドレス」は、一般にメモリと、マルチレベルキャッシュ2120などのキャッシュメモリとの間で転送される、メモリのキャッシュラインまたは単位に関連付けられる。キャッシュラインの「アドレス」はメモリロケーションを表すことが可能であり、キャッシュラインはメモリの複数のアドレスからのデータを含むことができる。「データ」という用語は、プリフェッチ可能な情報の単位を表すのに対して、「プログラム命令」および「プログラムデータ」という用語は、それぞれ、プロセッサによってその処理中に使用される命令およびデータを表す。したがって、データ(たとえば任意のビット数)は、プログラム命令またはプログラムデータあるいはその両方を構成する情報を示す「予測情報」を表すことができる。また「予測」という用語は、「予測アドレス」という用語と同じ意味で使用することもできる。予測アドレスがメモリへのアクセスに使用される場合、典型的には、その予測アドレスならびに他の(予測されるかまたはされない)アドレスを含む1つまたは複数のキャッシュラインがフェッチされる。

40

【0116】

50

プリフェッチャ2100は予測を発行する場合、各予測への参照を添付または関連付けることができる。予測が非順次予測である場合、それらに関連付けられる参照は予測識別子(「PID」)および対応するターゲットアドレスを含むことができる。PID(図示せず)は、対応するターゲットアドレスを予測させるトリガアドレス(またはその表現)を識別する。この参照は、メモリがプリフェッチされたデータを戻す場合、マルチレベルキャッシュ2120によって受け取られる。その後、マルチレベルキャッシュ2120は、プロセッサが要求するまでなどの間、戻されたデータを一時的に格納する。マルチレベルキャッシュ2120は、プリフェッチされたデータを格納する間、生成された予測に対してフィルタリングするため、その中に格納されたデータの一貫性を確認するため、そのデータを短期または長期データとして分類するためなどに、データを管理する。しかしながら、プロセッサがプリフェッチされたデータ(すなわち予測情報)を要求すると、そのデータはプロセッサに送信される。マルチレベルキャッシュ2120内に配置されたデータが非順次予測の結果である場合、必要であればターゲットキャッシュ内に格納された非順次予測の優先順位を再調整するために、非順次予測器2104に参照を送信することが可能である。

10

#### 【0117】

図22は、本発明の一実施形態に従った、例示的マルチレベルキャッシュ2220を示す図である。マルチレベルキャッシュ2220は、キャッシュフィルタ2210、第1レベルの戻りデータキャッシュ(「DRC1」)2222、および第2レベルの戻りデータキャッシュ(「DRC2」)2224を含む。キャッシュフィルタ2210は、それぞれ、第1レベルDRC 2222および第2レベルDRC 2124を、プリフェッチャ2100の構成要素ならびにメモリプロセッサ(図示せず)の構成要素などの他の構成要素とインターフェースさせるための、DRC1照会インターフェース2204およびDRC2照会インターフェース2214を含む。こうしたメモリプロセッサ構成要素の1つが、図21の書き戻しキャッシュ2290であり、これは、よく知られたキャッシュ方法に従って動作し、それによってキャッシュ内のデータへの修正は、必要となるまでキャッシュソース(たとえばシステムメモリ)にコピーされない。書き戻しキャッシュ2290の構造および機能は当分野で周知のものと同様であるため、詳細に論じる必要はない。さらに、DRC1照会インターフェース2204はDRC1マッチャ2206およびDRC1ハンドラ2208を含み、DRC2照会インターフェース2214はDRC2マッチャ2216およびDRC2ハンドラ2218を含む。

20

30

#### 【0118】

第1レベルDRC 2222は、アドレス(たとえば予測アドレス)を格納するためのDRC1アドレスストア2230を含み、DRC1アドレスストア2230は、データ(すなわち予測情報)およびPIDを格納するDRC1データストア2232に結合される。たとえば、予測アドレス(「PA」)の結果として生じるプリフェッチされたデータは、PID 2232bと関連付けられたdata(PA)2232aとして格納することができる。この表記法は、予測アドレスPAが予測情報を表すデータのプリフェッチに寄与したことを示す。data(PA)2232aがプロセッサによって要求されると、対応する予測アドレスPA、および予測識別子PID 2232bが、必要であればその予測アドレスの優先順位を修正するために非順次予測器2104に送られることになる。予測識別子PID 2232bは、一般に、PAを生じさせるトリガアドレスを示す情報を含む。プロセッサ要求アドレス(および関連データ)をマルチレベルキャッシュ2220に格納することも可能であるため、非順次予測器2104によって生成されるPAもターゲットアドレスと呼ぶことが可能であることに留意されたい。さらに、data(PA)2232aは必ずしもPID 2232bを伴う必要がないことにも留意されたい。

40

#### 【0119】

さらに、DRC1アドレスストア2230およびDRC1データストア2232はどちらも、その機能および/または構造を管理するDRC1マネージャ2234とも通信的に結合される。第2レベルDRC 2224は、データ2232aおよびPID 2232

50

bと同様の形でデータを格納するDRC2データストア2242に結合された、DRC2アドレスストア2240を含む。DRC2アドレスストア2240およびDRC2データストア2242はどちらも、その機能および/または構造を管理するDRC2マネージャ2246と通信的に結合される。

【0120】

本発明の特定の実施形態では、第2レベルDRC2224は、DRC2アドレスストア2240とは別の有効ビット2244を維持するための、「有効ビット」2244のリポジトリも含み、各有効ビットは、格納された予測が有効である（ならびにデータに関するプロセッサ要求を処理するために使用可能である）か、無効である（ならびに使用不可である）かを示す。無効予測を有するエントリは、空エントリとみなすことができる。有効ビット2244のビットをアドレスから分離して維持することにより、DRC2アドレスストア2240が有効ビットを対応するアドレスと共に格納する場合よりも1つまたは複数の有効ビットのリセットまたはセットは計算的な負担が軽くなり、迅速になる。ほとんどの場合、DRC1のアドレスに関する有効ビットは、通常、それらのアドレスと共に、またはその一部として格納されることに留意されたい。

10

【0121】

動作時に、DRC1照会インターフェース2204およびDRC2照会インターフェース2214は、第1レベルのDRC2222および第2レベルのDRC2224のコンテンツが、「入力アドレス」として印加される1つまたは複数のアドレスのうちのいずれかを含むかどうかを判別するために、そのコンテンツをそれぞれ検査するように構成される。入力アドレスは、生成された予測としてスペキュレータ2108から、書き込みアドレスとして書き戻しキャッシュから、またはマルチレベルキャッシュ2220外部にある他の要素からのものとすることができる。本明細書で説明するように、一般に、入力アドレスは、冗長をフィルタリング除去するためにマルチレベルキャッシュ2220のコンテンツと比較される、生成された予測である。しかしながら、時に入力アドレスは、データが書き込まれるかまたは今後書き込まれることになるメモリのロケーションを識別する書き込みアドレスである。この場合、マルチレベルキャッシュ2220は、メモリ、DRC1データストア2222、およびDRC2データストア2224の間の一貫性を維持するためのアクションが必要であるかどうかを特定するために、検査される。

20

【0122】

DRC1マッチャ2206およびDRC2マッチャ2216は、入力/出力ポート（「I/O」）2250上の1つまたは複数の入力アドレスが、DRC1アドレスストア2230およびDRC2アドレスストア2240内にそれぞれ常駐しているかどうかを特定するように構成される。DRC1マッチャ2206またはDRC2マッチャ2216のいずれかが、入力アドレスと、第1レベルのDRC2222および第2レベルのDRC2224内の入力アドレスとがマッチすることを検出した場合、DRC1ハンドラ2208またはDRC2ハンドラ2218などの関連付けられたハンドラは、冗長予測をフィルタリング除去するか、またはマルチレベルキャッシュ2220内のデータがメモリと一貫していることを確認するように動作する。DRC1マッチャ2206およびDRC2マッチャ2216は、入力アドレスの領域と、第1レベルのDRC2222および第2レベルのDRC2224のコンテンツとを並行して（すなわち、マルチレベルキャッシュ2220の構造に応じて、1または2サイクル内の動作（たとえばクロックサイクル）、あるいは他の最小数のサイクルなどで、同時に、またはほとんど同時に）比較するように構成可能であることに留意されたい。キャッシュと並行して比較可能な入力アドレスの領域の一例が、アドレスA0（トリガアドレス）および予測アドレスA1、A2、A3、A4、A5、A6、およびA7であり、後者の7つは、順次予測器2102によって生成されることが可能である。

30

40

【0123】

同時に検査される場合、こうした比較を実行するマッチャ2206、2216は、「先読みルックアップ」を実行していると言われる。いくつかの実施形態では、先読みルック

50

アップは、プロセッサがアイドル状態の場合、またはプリフェッチャ2100にデータを要求していない場合に、実行される。DRC1マッチャ2206およびDRC2マッチャ2216は、機能は同様であるが、それぞれの構造は、DRC1アドレスストア2230およびDRC2アドレスストア2240とそれぞれ動作するように適合されるため、必ずしも同様の構造ではないことにも留意されたい。DRC1マッチャ2206およびDRCマッチャ2216の例について、それぞれ、本発明の少なくとも1つの特定の実施形態に従った図23Aおよび図24に関連して以下で論じる。

#### 【0124】

次に、照会インターフェース2204、2214がフィルタリング操作を実行している場合の状況について考えてみる。いくつかの入力アドレスをマルチレベルキャッシュ2220のコンテンツと比較すること、および、マッチしない入力アドレスを検出することによって、ハンドラ2208、2218は、マッチしない入力アドレスが、フィルタリングが実行されなかった場合よりも早く、生成された予測として予測情報のフェッチを進行できるようにしながら、マッチした予測（すなわち冗長予測）をフィルタリング除去するための適切なアクションを実行することができる。したがって、マルチレベルキャッシュ2220およびそのキャッシュフィルタ2210は、どのキャッシュラインがフェッチを開始するかをより迅速に特定することによって、待ち時間を削減する。これにより、第1レベルのDRC2222および第2レベルのDRC2224のキャッシュが、一般に、予測が並行して比較されないかまたはフィルタリング除去されない場合、あるいはその両方の場合よりも早く、プリフェッチされた予測情報を含む可能性が高いことから、プロセッサが経験する待ち時間がさらに削減される可能性がある。

#### 【0125】

DRC1アドレスストア2230およびDRC2アドレスストア2240はそれぞれ、DRC1データストア2232およびDRC2データストア2242に格納されたプリフェッチされたデータに関連付けられたアドレスをそれぞれ格納する。アドレスストア2230および2240はそれぞれ、アドレスまたはアドレスの代替表現のいずれかを格納する。本発明の一実施形態によれば、例示的DRC1アドレスストア2230は、完全連想型であり、完全に固有のアドレスを格納するように構成される。たとえば、各アドレスについて、それらのアドレスを固有に識別するためにビット35:6がDRC1に格納される。DRC1アドレスストア2230に格納されたアドレスは、共通部分（たとえばタグ）およびデルタ部分（たとえばインデックス）を含むものとみなすことが可能であり、その両方が、少なくとも1つの実施形態に従ってDRC1の先読みルックアップ中にアドレスを表すために使用されることに留意されたい。さらに、DRC1アドレスストア2230およびDRC1データストア2232は、それぞれ、データのアドレスエン트리ごとに、32エントリのアドレスおよび64バイトのキャッシュラインを格納するように構成される。プリフェッチされたデータは、一般に、動的ランダムアクセスメモリ（「DRAM」）などのメモリからのものであるが、DRC1データストア2232内のデータが更新を必要とする場合は、書き戻しキャッシュからのものとするのが可能である。

#### 【0126】

これとは対照的に、例示的DRC2アドレスストア2240は、4wayセットの関連エントリからなり、アドレスを表すための基本部分（たとえばタグ）を格納するように構成することができる。さらに、DRC2アドレスストア2240およびDRC2データストア2242は、それぞれ、データのアドレスエン트리ごとに、1024エントリのアドレスおよび64バイトのキャッシュラインを格納するように構成される。DRC2データストア2242は、DRC1データストア2232からのプリフェッチデータを格納し、いくつかの実施では、任意数のメモリバンク（たとえば4つのバンク：0、1、2、および3）からなるものとするのができる。

#### 【0127】

予測情報のプリフェッチ元であるメモリは、通常、DRAMメモリ（たとえば、デュアルインラインメモリモジュール、すなわち「DIMM」内に配置構成される）であるが、

10

20

30

40

50

メモリは任意の他の知られたメモリ技術のものとする事ができる。通常、メモリは、特定の行アドレス内で使用可能なメモリのセクションである「ページ」に細分される。特定のページにアクセスする、すなわち「オープン」すると、他のページがクローズされ、このページのオープンおよびクローズのプロセスは完了するまでに時間を要する。したがって、プロセッサがDRAMメモリの様々なメモリロケーションからの命令およびデータのフェッチに関して、やや散漫な様式でプログラム命令を実行している場合、メモリへのアクセスは非順次的である。したがって、読み取り要求のストリームはページ領域をまたがって延在する可能性がある。次のページの次のアドレスが使用できない場合、通常、プロセッサはプログラム命令およびプログラムデータをメモリから直接フェッチしなければならない。これにより、こうした命令およびデータの取り出し待ち時間が増加する。そこで、マルチレベルキャッシュ2220内の複数ページにまたがる予測情報をプリフェッチおよび格納することによって、ページのオープンに関する待ち時間が本発明に従って削減される。さらに、プリフェッチされているデータがキャッシュからのものであるため、アクセスされたページはオープンされたままで、プロセッサによって認識されるかまたはプロセッサに関する待ち時間が削減される。

#### 【0128】

たとえば、非順次予測器2104が、アドレス「00100」のプロセッサ読み取りに続いてアドレス「00200」にアクセスすることになると、正しく予測するものと考えてみる。したがって、非順次予測器2104は、プロセッサが実際にアドレス「00200」にアクセスするより前に、アドレス「00200」（ならびに、バッチが4の場合、アドレス00201、00202、00203、および00204）で開始するラインの範囲（たとえば、1つのターゲットアドレスおよび4つの予測アドレスであり、生成する予測の数はバッチ「b」によって構成可能かつ定義される）をフェッチさせる。プロセッサが実際にアドレス「00200」に関する読み取りを実行する場合、マルチレベルキャッシュ2220の先読みルックアップは、アドレス「00200」に続く指定された範囲内で、どのキャッシュラインがすでにプリフェッチされているかを即時に特定する。読み取りアドレスストリーム内の非順次移行は、DRAMページオープン動作を伴うことが可能であるため、先読みルックアップは、プリフェッチャ2100が、読み取り要求のストリーム内で即時に先読みすること、および、どのアドレスまたはキャッシュラインをフェッチする必要があるかを特定することを可能にする。フェッチを即時に開始することにより、プリフェッチャ2100は、DRAMページオープン動作の待ち時間をしばしば隠し、その後、プロセッサ上での待ち時間の損失を招くことなく、キャッシュラインの順次ストリームを提供する（ターゲットアドレスに関する基準を形成するトリガアドレスとは非順次的であるが）ことができる。

#### 【0129】

図22は、DRC1マネージャ2234およびDRC2マネージャ2246を別のエンティティとして示すが、必ずしもそうである必要はない。すなわち、DRC1マネージャ2234およびDRC2マネージャ2246は、単一の管理エンティティに組み合わせるか、またはマルチレベルキャッシュ2220の外部に配置する、あるいはその両方とすることができる。第1レベルのDRC2222および第2レベルのDRC2224は、プロセッサ内に常駐する従来のL1およびL2キャッシュとは構造的および/または機能的に異なるため、マルチレベルキャッシュ2220内に格納された予測情報を管理する固有のポリシーが採用される。こうしたポリシーの例には、各戻りデータキャッシュ内のメモリを割り振るためのポリシー、短期から長期のデータストアへ情報をコピーするためのポリシー、および、マルチレベルキャッシュ2220と書き戻しキャッシュなどの他のエンティティとの間の一貫性を維持するためのポリシーが含まれる。

#### 【0130】

第1に、情報が短期情報から長期情報へと古くなる場合の、第1レベルのDRC2222から第2レベルのDRC2224への予測情報のコピーを管理するために使用されるコピーポリシーについて考えてみる。データが、ある一定の時間しきい値まで第

10

20

30

40

50

1レベルのDRC 2222にある場合、DRC 1マネージャ2234はDRC 2マネージャ2246と協働して、DRC 1データストア2232からDRC 2データストア2242へとそのデータを転送する。しきい値は一定であるか、またはそうでなければ動作時に変化してもよいことに留意されたい。通常、DRC 1内にある無効エントリ（すなわち使用可能）がN個未満の場合は必ず、古くなったデータは転送されるように構成され、ここでNはプログラム可能である。動作時に、データが短期ストレージから長期ストレージへとコピーされた場合、第1レベルのDRC 2222内のエントリは消去（すなわち無効化）される。

#### 【0131】

第2に、第1レベルのDRC 2222および第2レベルのDRC 2224に予測情報を挿入するための割り振りポリシーについて考えてみる。予測情報を第1レベルのDRC 2222に挿入する場合、DRC 1マネージャ2234は、候補としてロックされたエントリを除いて、DRC 1データストア2232内の任意の無効なエントリを選択する。DRC 1マネージャ2234が予測情報を格納することが可能ないずれの無効エントリも検出しない場合、最も古いエントリを使用してエントリ用のスペースを割り振ることができる。DRC 2データストア2242におけるエントリの割り振りについても同様に、DRC 2マネージャ2246は、第1レベルのDRC 2222から第2レベルのDRC 2224へとコピーされたデータを受け取るためのいくつかのwayのうちの一つ（たとえば4wayのうちの一つ）を使用することができる。たとえば、予測アドレスのインデックスは、データを格納する4つのエントリを含むことができる。初期に、DRC 2データストア2242は、使用されていない（すなわち無効化された）way数のうちの一つを割り振る。しかしながら、すべてのwayが割り当てられている場合、第1のinが第1のoutである（すなわち、最も古いものが上書きされる）。しかしながら、最も古いエントリが同じ古さであり、かつ有効な場合、DRC 2マネージャ2246はロックされていないエントリを割り振る。最後に、wayセット内のすべてのエントリがロックされている場合、DRC 2マネージャ2246は、第1レベルのDRC 2222内のエントリを有効として維持しながら、第1レベルのDRC 2222から第2レベルのDRC 2224への書き込みを抑制する。ここでも、典型的には第2レベルのDRC 2224は、第1レベルのDRC 2222からのみ、ストレージ用のデータを受け取ることに留意されたい。

#### 【0132】

DRC 1マネージャ2234およびDRC 2マネージャ2246が遵守可能な他のポリシーは、一貫性を維持することに関する。DRC 1マネージャ2234は、データが書き込まれることになる書き込みアドレスとマッチするアドレスを有する、任意のエントリのデータを更新することによって、第1レベルのDRC 2222の一貫性を維持する。典型的には、書き戻しキャッシュ2290（図21）は、書き込みアドレスを書き込むためにメモリ（たとえばDRAM）に送信するまで、書き込みアドレス（および対応するデータ）を一時的に格納する。書き戻しキャッシュ2290内の書き込みアドレスとマッチする読み取り要求のアドレスがあるいくつかのケースでは、マルチレベルキャッシュ2220がデータを第1レベルのDRC 2222に転送するのに先立って、書き込みアドレスのデータとメモリのそれとをマージすることに留意されたい。DRC 2マネージャ2246は、書き戻しキャッシュ2290内にロードされる場合に、そのアドレスが書き込みアドレスとマッチする任意のエントリを無効化することによって、第2レベルのDRC 2224の一貫性を維持する。第2レベルのDRC 2224がDRC 1からのデータのみを受け取るため、および、第1レベルのDRC 2222がメモリおよび書き戻しキャッシュ2290との一貫性を維持するため、第2レベルのDRC 2224は一般に、陳腐化したデータを含むことがない。さらに、DRC 1からDRC 2へとコピーされる予定の任意のアドレスを、第1に書き戻しキャッシュ（「WBC」）2290に照らしてチェックすることができる。WBC 2290内にマッチが見つかった場合、コピー操作は中止される。見つからなかった場合、そのアドレスのDRC 1からDRC 2へのコピーは実行

10

20

30

40

50

される。この追加のチェックが、一貫性の維持にさらに役立つ。

【0133】

図23Aは、本発明の特定の実施形態に従った、第1のアドレスストア2305に関する例示的DRC1の照会インターフェース2323を示す図である。この例では、トリガアドレス(「A0」)2300(たとえば、プロセッサ要求アドレス)は、入力アドレスとして、共通アドレス部分2302aおよびデルタアドレス部分2302bからなる。アドレス2300は、いくつかのケースでの予測アドレス、または他のケースでの書き込みアドレス(一貫性を維持する場合)の、いずれかとする 것도可能であることに留意されたい。アドレス2300が予測アドレスのグループを生成するトリガアドレスの場合、こうしたグループ2307は、アドレス(「A1」)2301からアドレス(「Am」)2303まで識別されているようなアドレスを含むことが可能であり、ここで「m」は、本発明の少なくとも一実施形態に従って「先読みルックアップ」を実行する際に使用可能な任意数の予測を表す。いくつかのケースでは、「m」はバッチサイズ「b」と等価に設定される。

10

【0134】

DRC1アドレスストア2305のエントリ2306は、それぞれ、第1のエントリ部分2306a(たとえばタグ)および第2のエントリ部分2306b(たとえばインデックス)を含む。特定の実施形態では、第1のエントリ部分2306aおよび第2のエントリ部分2306bは、それぞれ、共通アドレス部分2302aおよびデルタアドレス部分2302bに類似している。第2のエントリ部分2306bは、アドレスに関して、トリガアドレス(「A0」)2300からその特定のエントリ2306への移動を示す。したがって、DRC1マッチャ2312が、トリガアドレス(「A0」)2300などの入力アドレスとエントリ2306とを比較する場合、共通部分2302aを使用してグループ2307のアドレスの共通部分を表すことができる。さらに、アドレス2300の共通部分2302aが、一般にアドレス(「A1」)2301から(「Am」)2303までの共通部分と同様であるため、エントリ2306の1つまたは複数の第1のエントリ部分2306aと比較するために使用する必要があるのは共通部分2302aのみである。また、アドレス(「A1」)2301から(「Am」)2303までのデルタ部分2302bを、エントリ2306の複数の第2のエントリ部分2306bとマッチングさせることも可能である。

20

30

【0135】

一実施形態では、DRC1マッチャ2312は、共通アドレス部分と第1のエントリ部分とをマッチングするための共通コンパレータ2308、およびデルタアドレス部分と第2のエントリ部分とをマッチングするためのデルタコンパレータ2310を含む。具体的に言えば、エントリ0からn番目のエントリについて、共通部分2302aと第1の部分2306aとが同時に比較され、同じエントリについて、デルタ部分2302bと第2の部分2306bとが同時に比較される。いくつかの実施形態では、共通コンパレータ2308は、高位ビット(たとえば、36ビットアドレスのビット35:12)を比較するための「ワイド」コンパレータであり、デルタコンパレータ2310は、低位ビット(たとえば、36ビットアドレスのビット11:6)を比較するための「ナロー」コンパレータである。図23Aは、1つのデルタ部分2302bにつき1つのデルタコンパレータを示すが、いくつかのケースでは、デルタコンパレータ2310の数は $m \cdot n$ に等しく(図示せず)、ここで各デルタコンパレータは、入力として1つのデルタ部分2302bおよび1つの第2のエントリ部分2306bを受け取ることに留意されたい。コンパレータサイズは、これらの比較を実行するために必要な物理リソースの量を制限するため、並行してルックアップされるアドレスは、同じメモリページ内に存在するように構成される(たとえば、メモリページサイズは通常4Kバイトである)。これによって、交差するページ境界からの先読みルックアップのアドレスは減少するが、これらの構成は、物理リソースに関して先読みルックアップを実行するためのコストを削減する。ここでも、共通部分2302aおよびデルタ部分2302bはそれぞれ同時に、またはほぼ同時に、エントリ23

40

50



06と比較されることに留意されたい。

【0136】

共通コンパレータ2308およびデルタコンパレータ2310の出力は、それぞれ、Hbase(0)、Hbase(1)、...Hbase(m)、およびH0、H1、H2、...HNであり、ここでそれぞれは0(たとえばマッチなしを示す)または1(たとえばマッチを示す)のいずれかである。この結果は、フィルタリングしているかまたは一貫性を維持しているかに応じて、アクションを実行するためにDRC1ハンドラ2314に送信される、0および1のヒットベクトルを形成する。ヒットリスト生成器2313は、範囲「r」(すなわちグループ2307)内のどのアドレスがDRC1アドレスストア2305内に常駐するかを示す、ヒットのリスト(「ヒットリスト」)を生成する。アドレスがマッチした(すなわち、その中に予測が格納されている)場合、そのアドレスはヒットリストに含められ、マッチしないアドレス(すなわち予測が格納されていない)はヒットリストから除外される。このヒットリストは、予測を生成するため、またはDRC1アドレスストア2305内の一貫性を管理するために使用される。

10

【0137】

図23Bは、特定の実施形態に従った、図23AのDRC1照会インターフェース2323を使用して並行して検査可能な任意数の例示的入力アドレス2352を示す図である。ここでDRC1照会インターフェース2350は、DRC1アドレスストア2305とマッチングするために、任意の範囲のアドレス2352を受け入れることができる。図23Aのマッチャ2312は、いくつかの入力アドレスにわたって、並行先読みルックアップを実行するために必要な回数だけ複製される。一例として、パッチサイズ「b」が27に設定された前方順次予測の場合、DRC1照会インターフェース2350はマッチャに、基本(またはトリガ)アドレスとしてのA0と、グループ2307としての予測アドレスA1からA7とを、並行してマッチさせるように要求する。ブラインドバック予測の場合、A(-1)のみが、グループ2307として基本アドレスA0以外のマッチングを必要とするが、逆順次予測の場合、アドレスA(-1)からA(-7)がマッチングを必要とする。アドレス2352の範囲は、DRC1およびDRC2の両方の照会インターフェースにも同時に、並行して印加できることに留意されたい。

20

【0138】

図24は、本発明の特定の実施形態に従った、DRC2アドレスストア2404に関する例示的DRC2照会インターフェース2403を示す図である。DRC2照会インターフェース2403は、DRC2アドレスストア2404のコンテンツとアドレスを比較するために入力アドレス2402を受け取るように構成される。この例で、入力アドレス2402は、tag(A0)などのアドレスの基本部分(たとえばタグ)である。さらにこの例を見ると、DRC2アドレスストア2404は、バンク0、1、2、および3というメモリの4つのバンク2406からなり、それぞれのバンクがエントリ2410を含んでいる。このケースでは、エントリ2410を4つのway(W0、W1、W2、およびW3)のうちのいずれかが1つに配置できることに留意されたい。

30

【0139】

DRC2マッチャ2430は、tag(A0)をエントリ2410と比較するためのいくつかのコンパレータを含む。一般に、DRC2アドレスストア2404内の任意のマッチングアドレスは、同じtag(A0)を共有するが、他のビットグループとの関係は異なる(たとえばインデックスごと)場合がある。本発明の特定の実施形態では、タグがDRC2アドレスストア2404内のいずれかのエントリとマッチするかどうかの特定が、一般に以下のように実行される。第1に、各バンク2406について、そのバンク内のインデックスのうちの1つが、潜在的マッチングアドレスを探索するために選択される。これは、図25Aに示されるように、バンクが特定アドレス(たとえばA0)のあるインデックスビットによって識別できる場合、探索用を選択されるバンクは、特定のアドレス(たとえば図25のA0)がどのバンクに常駐するかによって異なるため、バンクごとに異なる場合がある。第2に、各バンク2406について選択されたインデックスの4つのw

40

50

ayすべてにアクセスする。次に、4つのway（たとえばW0からW3）に関して格納されたタグが、この例では基本アドレス2402であるtag（A0）と比較される。一般に、tag（A1）などの他のタグと比較することなく、tag（A0）と比較するだけで十分である。これは、これらのタグが一般に等しい（たとえば、tag（A0）=tag（A1）=tag（A2））と想定されるためである。予測に関する同時探索は、通常、4kバイトページなどの同じページ内にある予測に限定され、これによってタグが同じになることに留意されたい。第3に、DRC2マッチャ2430によってアドレスマッチが実行されると、ヒットベクトルおよび有効ビットの形の結果を使用して、図27および28に関連して説明するのと同様に、最終ヒットベクトルが取得される。

#### 【0140】

DRC2照会インターフェース2403のヒット生成器2442は、タグの比較結果（「TCR」）2422をDRC2マッチャ2430から受け取り、さらにそれらの結果を対応する有効ビット2450と比較して、順序付けされた予測のセット（「順序付けされた予測」）を生成する。ここで、バンク1、2、3、および4からのタグの比較結果は、それぞれTCR（a）、TCR（b）、TCR（c）、TCR（d）とラベル付けされ、それぞれが、タグが1つまたは複数のエントリ2410とマッチするかどうかを表す1つまたは複数のビットを含む。順序付けされた予測は、入力アドレス2402とマッチする（またはマッチしない）予測の順序付けされたセットとすることができる。または、順序付けされた予測それぞれを、入力アドレスがDRC2アドレスストア2404内に存在するアドレスを有するかどうかを表す、ビットのベクトルとすることができる。追加のDRC2マッチャ2430が含まれる場合、任意数の入力アドレス2402が同様にDRC2照会インターフェース2403とマッチングできることに留意されたい。図25Aから28は、本発明のいくつかの実施形態に従った例示的ヒット生成器を示す図である。

#### 【0141】

図25Aは、本発明の一実施形態に従った、DRC2アドレスストア2404に格納されたアドレス（またはその表現）の可能な配置構成を示す図である。以下の考察を簡単にするために、way W0、W1、W2、およびW3は示されていないことに留意されたい。入力アドレスA0、A1、A2、およびA3は、DRC2アドレスストア2404に格納される。一例として、順次予測器2102（図示せず）は、トリガアドレスA0（たとえば4つのwayのうちのいずれかにある）に基づいて、順次予測A1、A2、およびA3を生成することができる。第1の配置構成2502は、A0がバンク0に格納された結果である。同様に、第2の配置構成2504、第3の配置構成2506、および第4の配置構成2508は、それぞれアドレスA0をバンク1、2、および3に格納した結果であり、後続のアドレスは続くトリガアドレス内に順番に格納される。したがって、これらのアドレス（または、タグの形などのその一部）は、一般に、特定の順序のないDRC2アドレスストア2404からの出力である。

#### 【0142】

図25Bは、本発明の実施形態に従った、順序付けされていないアドレスおよび対応する有効ビットに基づいて結果を生成する、例示的ヒット生成器2430を示す図である。この例では、順次予測器2102がトリガアドレスA0に基づいて順次予測A1、A2、A3、A4、A5、A6、およびA7を生成し、そのすべてが、図に示された特定の配置構成で格納される（すなわち、トリガアドレスA0がバンク1に格納され、その他がそれに続く）。ヒット生成器2430は、順序付けされていないアドレスA2、A6、A1、A5、A0、A4、A3、A7、および順序付けされた有効ビットVB0からVB7を受け取り、それらを順序付けし、それらを比較した後、ビットベクトルまたはアドレスのリスト（マッチするかまたはマッチしないのいずれか）とすることが可能な、結果R0からR7を生成する。予測が無効であることを示す有効ビットは、格納された無効予測がマッチングされないようにすることに留意されたい。これは、有効ビットをアドレスストアのコンテンツとマッチングするための1つの理由である。本発明の特定の実施形態によれば、アドレスA2、A1、A0、およびA3、またはアドレスA6、A5、A4、およびA

10

20

30

40

50



することによって、それぞれいくつかのビットからなるアドレスを再順序付けするよりも、必要なハードウェアが少なくすむ。順序付け器 2702 および結果順序付け器 2810 の順序付けは例示的なものであり、ビットを順序付けおよび再順序付けするための他のマッピングも本発明の範囲内であることに留意されたい。

#### 【0146】

本発明の特定の実施形態では、非順次予測器 2104 およびマルチレベルキャッシュ 2120 を含む図 21 のプリフェッチャ 2100 は、ノースブリッジチップの同じ機能のうちの少なくともいくつかを有するメモリプロセッサ内などの、ノースブリッジ - サウスブリッジチップセットアーキテクチャ内に配置される。メモリプロセッサは、CPU、グラフィックスプロセッサユニット（「GPU」）などの 1 つまたは複数のプロセッサによって、少なくともメモリアクセスを制御するように設計される。ノースブリッジの実施では、プリフェッチャ 2100 は AGP / PCI Express インターフェースを介し GPU に結合することもできる。さらに、プロセッサとメモリとの間のシステムバスとして、フロントサイドバス（「FSB」）を使用することもできる。また、メモリはシステムメモリとすることもできる。別法として、メモリプロセッサが実行するのと同様にメモリへのアクセスを制御する働きをする、任意の他の構造、回路、デバイスなどで、マルチレベルキャッシュ 2120 を採用することができる。さらに、マルチレベルキャッシュ 2120 およびその要素、ならびにプリフェッチャ 2100 の他の構成要素は、ハードウェアまたはソフトウェアモジュールのいずれか、あるいはその両方からなるものとするのが可能であり、さらに、任意の様式で分散または結合することも可能である。

#### 【0147】

説明のために、前述の記述では、本発明を完全に理解するための特定の命名法を使用した。しかしながら、当業者であれば、本発明を実施するために特定の細部が必要でないことを理解されよう。したがって、本発明の特定の実施形態の前述の記述は、例示および説明のために提示されたものである。本発明を網羅するか、または開示された精密な形に限定することは意図されておらず、前述の教示に鑑みて、多くの修正形態および変形形態が可能であることは明らかである。実際、この記述は、本発明のいずれかの特徴または態様をいずれかの実施形態に限定するものとして読むべきではなく、むしろ一実施形態の特徴および態様は他の実施形態と容易に交換可能である。諸実施形態は、本発明の原理およびその実際の応用例を最もよく説明するために選択および説明されたものであり、それによって他の当業者が、企図された特定の使用に好適となるような様々な修正形態によって、本発明および様々な実施形態を最適に利用できるようにするものである。添付の特許請求の範囲およびそれらの等価物が本発明の範囲を画定することが意図されている。

#### 【図面の簡単な説明】

#### 【0148】

【図 1】本発明の特定の実施形態に従った、メモリプロセッサと共に実施される例示的スペキュレータを示すブロック図である。

【図 2】本発明の一実施形態に従った、例示的スペキュレータを示す図である。

【図 3 A】本発明の特定の実施形態に従った、例示的前方順次予測器を示す図である。

【図 3 B】本発明の特定の実施形態に従った、例示的ブラインドバック順次予測器を示す図である。

【図 3 C】本発明の特定の実施形態に従った、例示的バックセクタ順次予測器を示す図である。

【図 3 D】本発明の特定の実施形態に従った、例示的逆順次予測器の挙動を示す図である。

【図 4】本発明の一実施形態に従った、例示的非順次予測器を示す図である。

【図 5】本発明の一実施形態に従った、インタリーブされた順次アドレスのストリームに対して非順次予測を抑制する例示的技法を示す図である。

【図 6】本発明の一実施形態に従った、複数のスレッドにわたってインタリーブされた順次アドレスに対して非順次予測を抑制する例示的技法を示す図である。

【図 7】本発明の特定の実施形態に従った、基本アドレスおよび非順次アドレスの着信時間に基づいて非順次予測を抑制するための他の技法を示す図である。

【図 8】本発明の特定の実施形態に従った、予測の生成を促進するための例示的技法を示す図である。

【図 9】本発明の一実施形態に従った、予測フィルタを含む他の例示的スペキュレータを示す図である。

【図 10】本発明の特定の実施形態に従った、例示的非順次予測器を実施するプリフェッチャを示すブロック図である。

【図 11】本発明の一実施形態に従った、例示的非順次予測器を示す図である。

【図 12】本発明の実施形態に従った、例示的予測生成器を示す図である。

10

【図 13】本発明の特定の実施形態に従った、例示的優先順位調整器を示す図である。

【図 14】本発明の特定の実施形態に従った、非順次予測を形成する場合に非順次予測器生成器を動作させるための例示的パイプラインを示す図である。

【図 15】本発明の特定の実施形態に従った、非順次予測を優先順位付けするように優先順位調整器を動作させるための例示的パイプラインを示す図である。

【図 16】本発明の特定の実施形態に従った、メモリプロセッサ内の例示的予測インベントリを示すブロック図である。

【図 17】本発明の一実施形態に従った、例示的予測インベントリを示す図である。

【図 18】本発明の特定の実施形態に従った、インベントリフィルタの例を示す図である。

20

【図 19 A】本発明の特定の実施形態に従った、冗長をフィルタリング除去する例示的技法を示す図である。

【図 19 B】本発明の特定の実施形態に従った、冗長をフィルタリング除去する例示的技法を示す図である。

【図 20】本発明の一実施形態に従った、プリフェッチャ内に配置される他の例示的予測インベントリを示す図である。

【図 21】本発明の特定の実施形態に従った、例示的キャッシュメモリを含むプリフェッチャを示すブロック図である。

【図 22】本発明の一実施形態に従った、例示的マルチレベルキャッシュを示す図である。

30

【図 23 A】本発明の特定の実施形態に従った、第 1 のアドレスストアに関する例示的の照会インターフェースを示す図である。

【図 23 B】図 23 A の第 1 の照会インターフェースを使用して並行して検査可能な任意数の入力アドレスを示す図である。

【図 24】本発明の特定の実施形態に従った、第 2 のアドレスストアに関する例示的の照会インターフェースを示す図である。

【図 25 A】本発明の一実施形態に従った、第 2 のアドレスストアに格納される場合の例示的アドレス（またはその表現）の可能な配置構成を示す図である。

【図 25 B】本発明の実施形態に従った、順序付けされていないアドレスおよび順序付けされた有効ビットに基づいて結果を生成する、例示的ヒット生成器を示す図である。

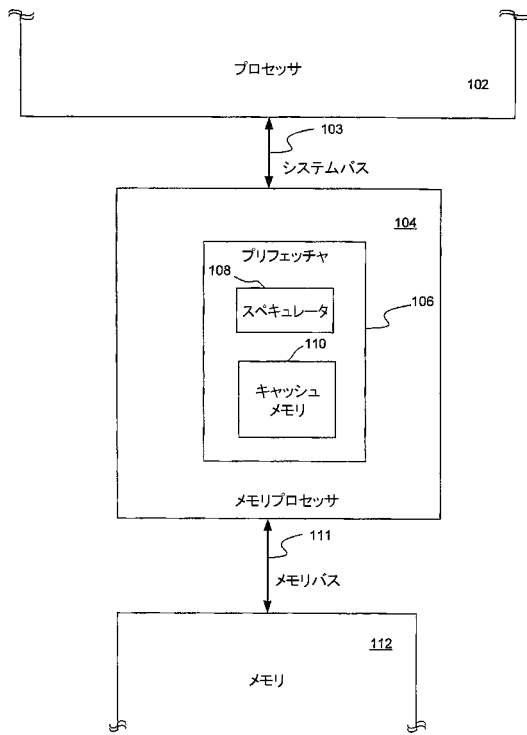
40

【図 26】本発明の実施形態に従った、図 25 のヒット生成器の 1 つの結果、R を生成するための構成要素を示す概略図である。

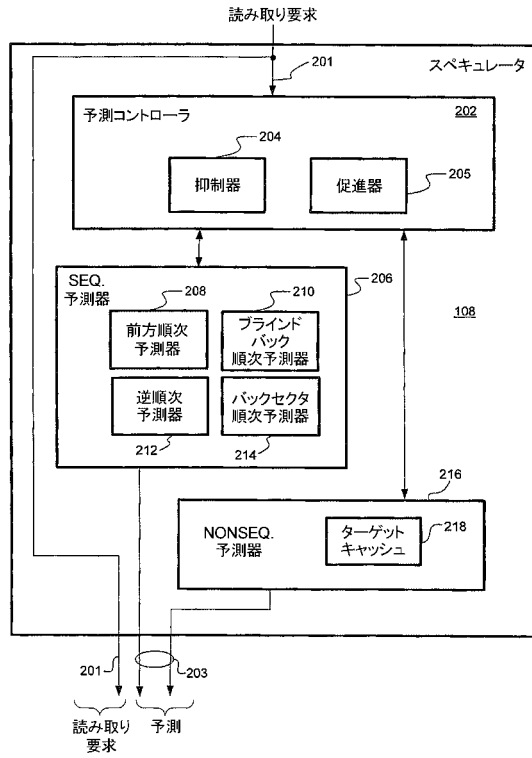
【図 27】本発明の特定の実施形態に従った、ヒット生成器の一例を示す図である。

【図 28】本発明の他の実施形態に従った、ヒット生成器の他の例を示す図である。

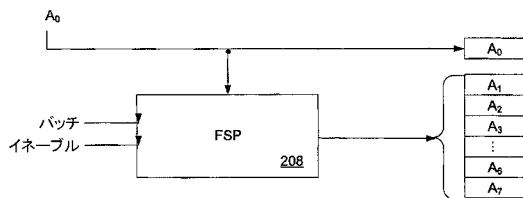
【図 1】



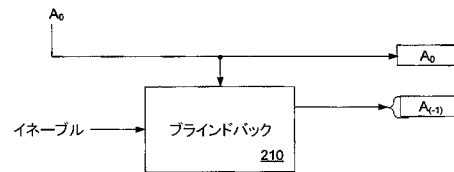
【図 2】



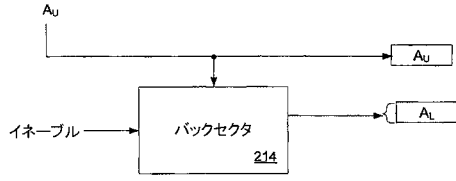
【図 3 A】



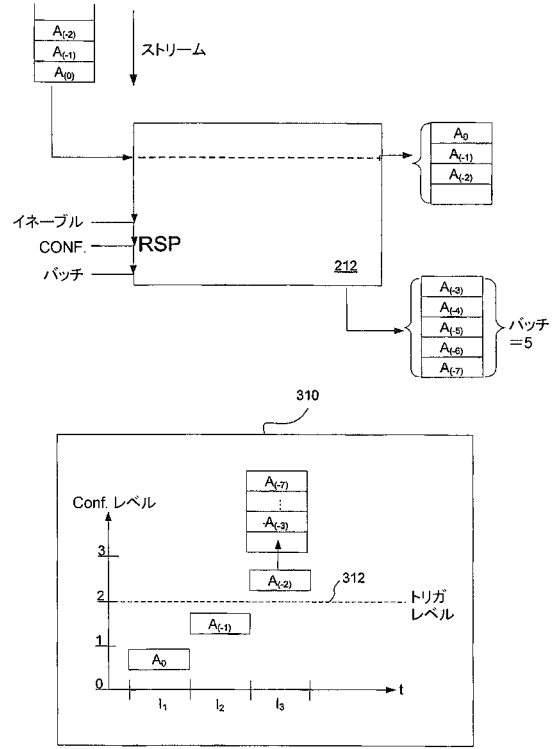
【図 3 B】



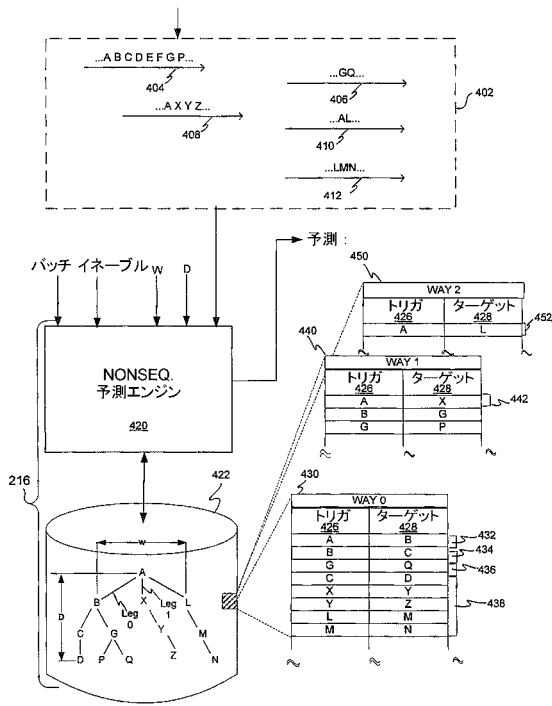
【図3C】



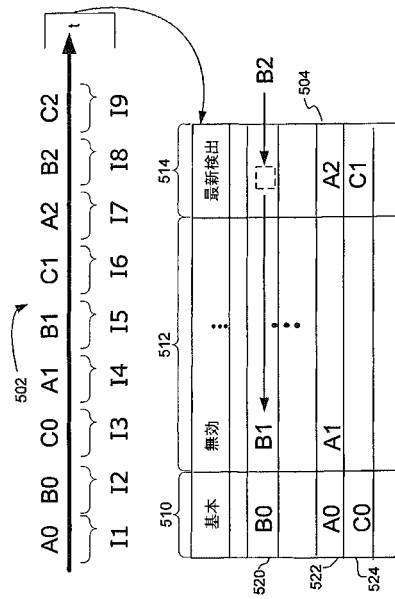
【図3D】



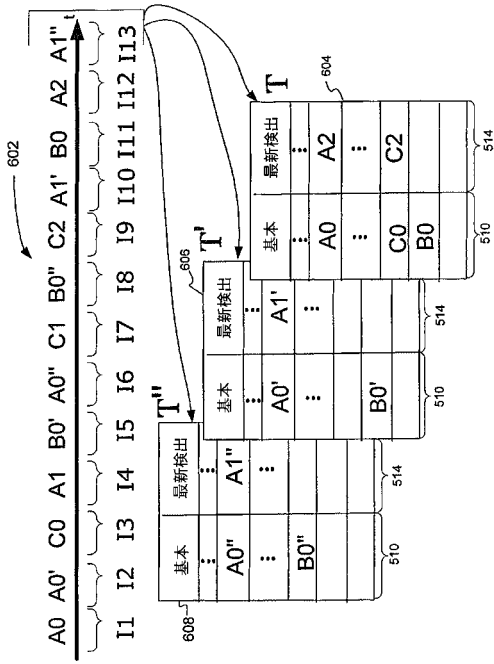
【図4】



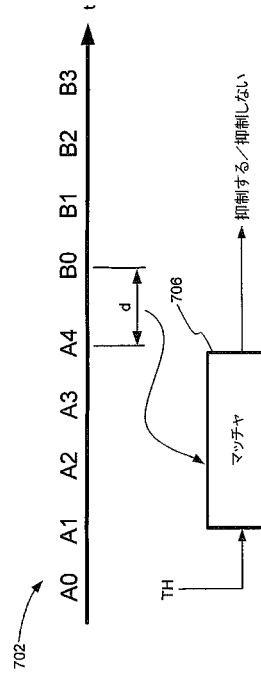
【図5】



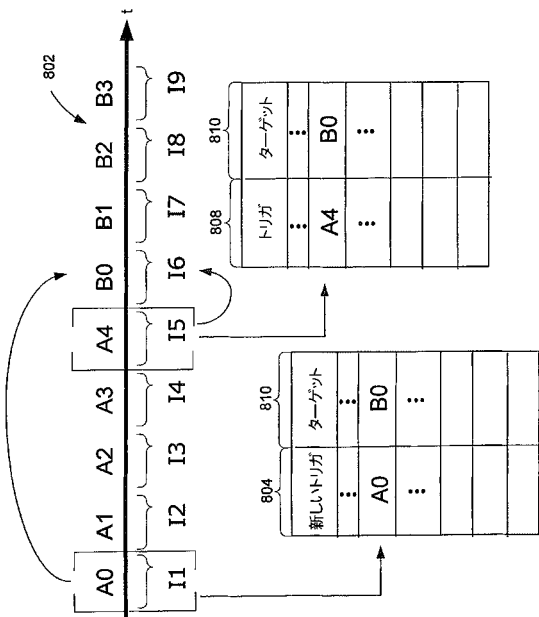
【図6】



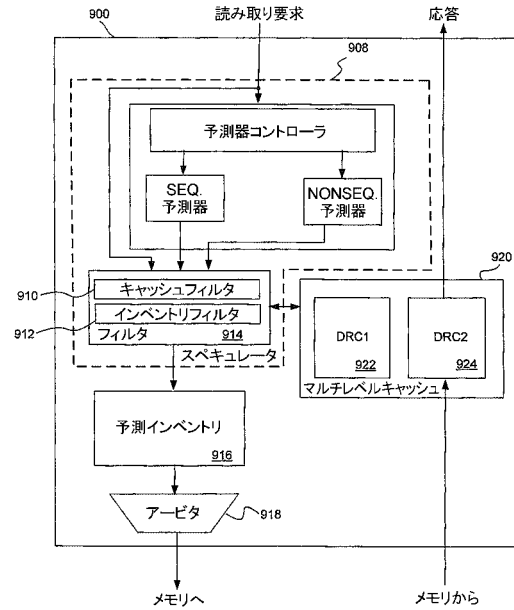
【図7】



【図8】

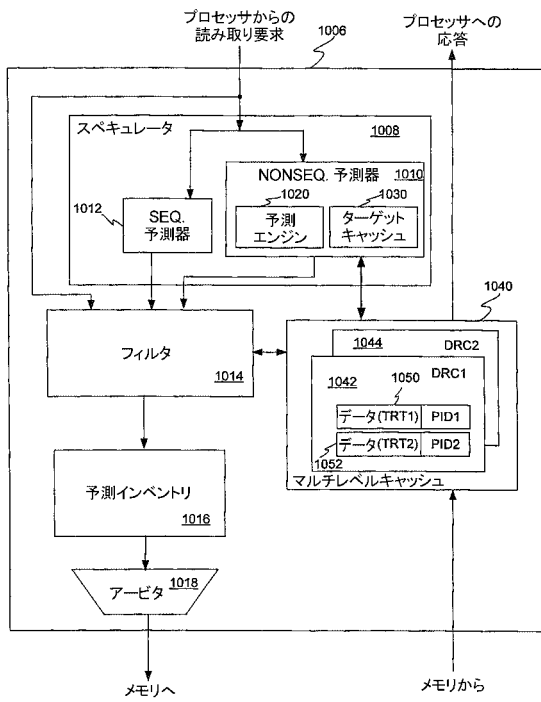


【図9】

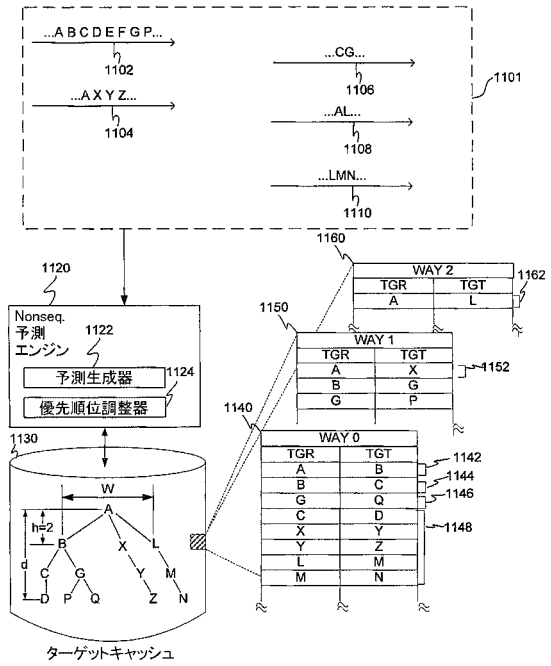




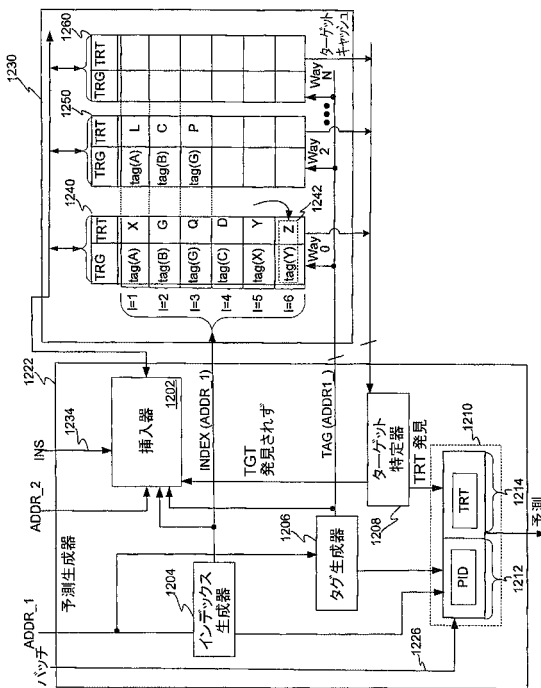
【図10】



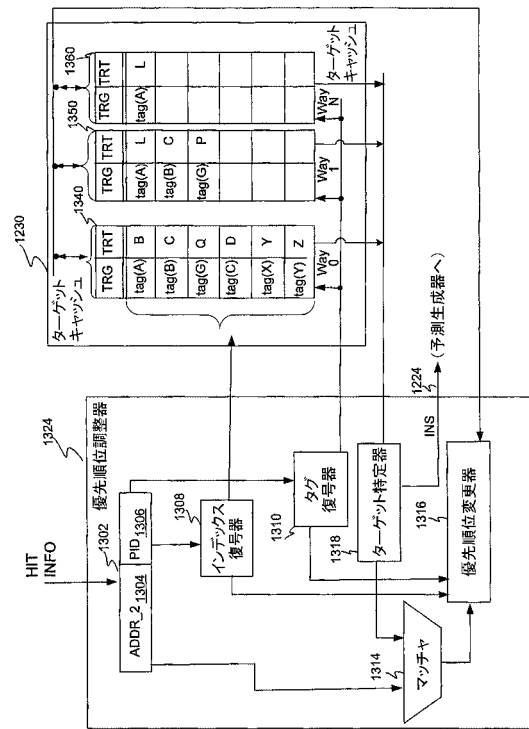
【図11】



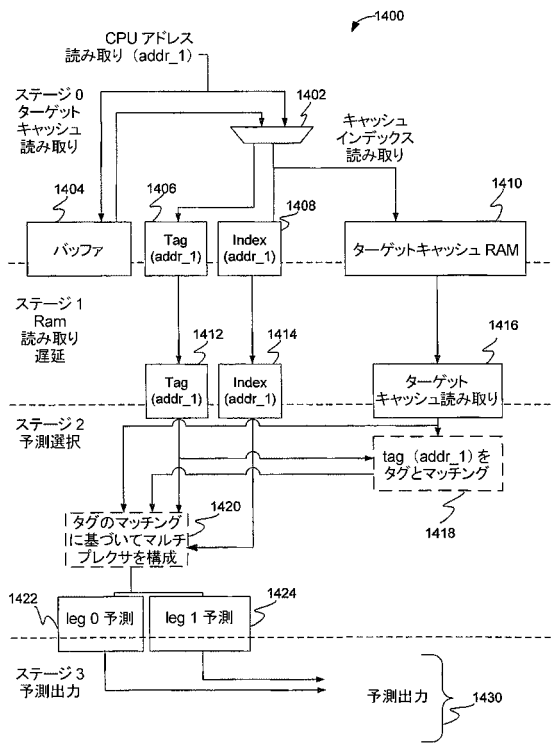
【図12】



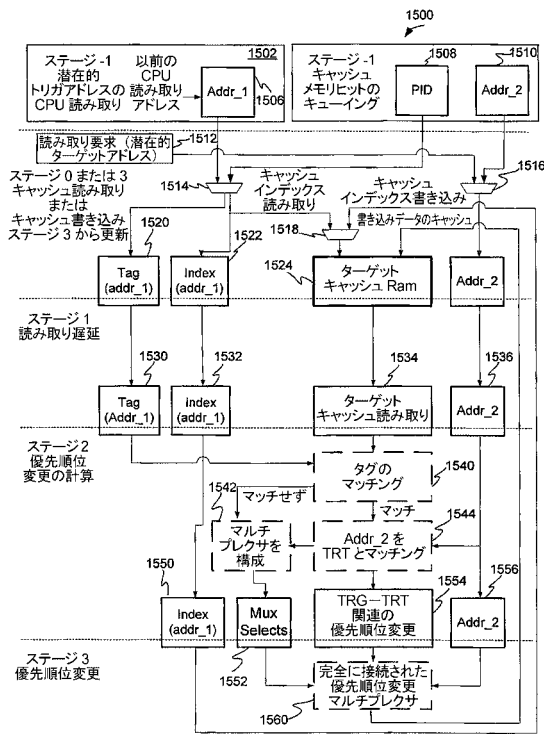
【図13】



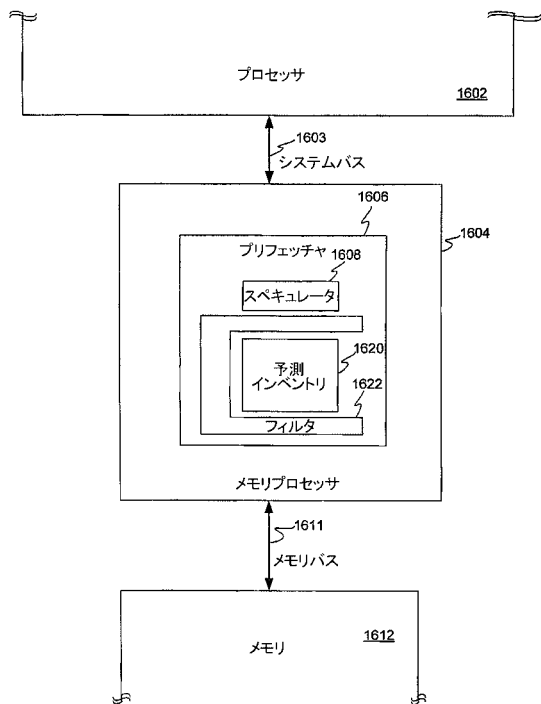
【図14】



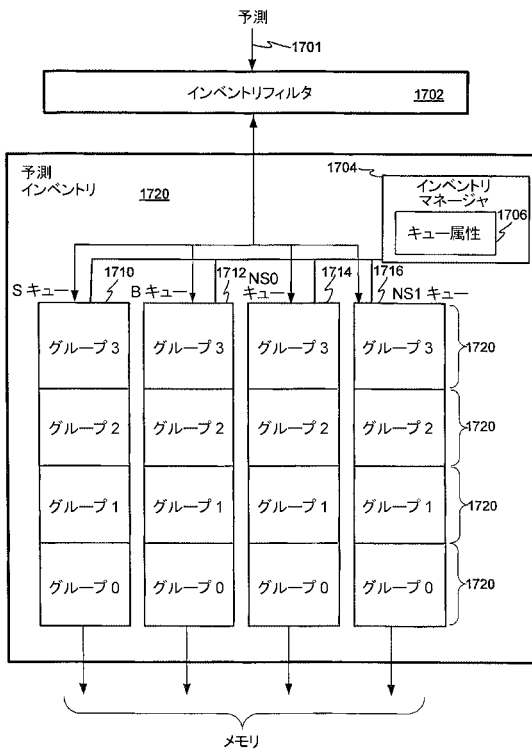
【図15】



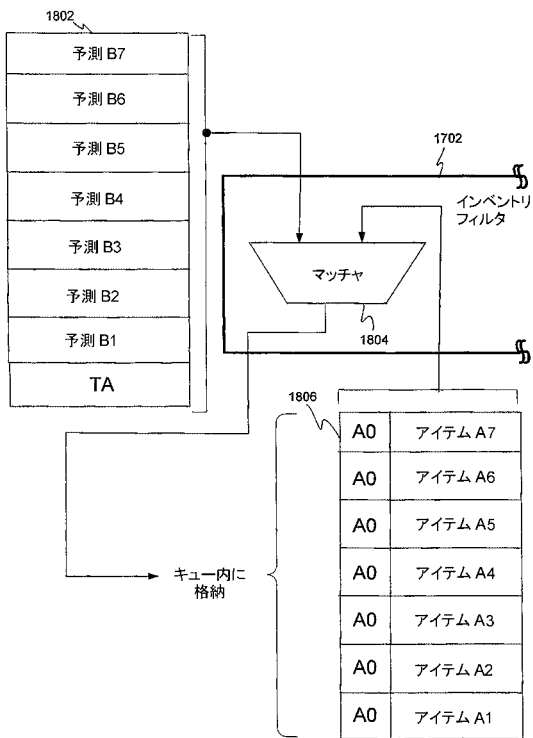
【図16】



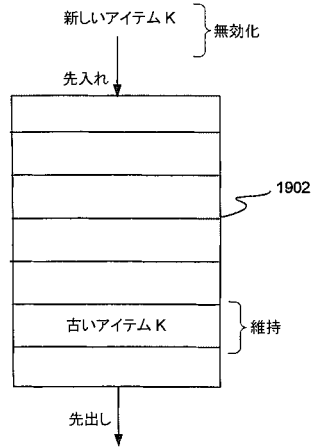
【図17】



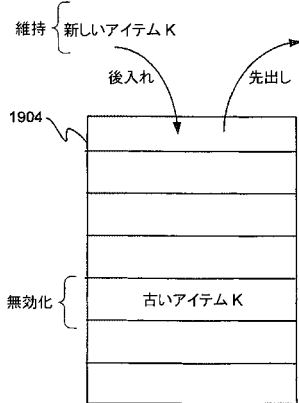
【図18】



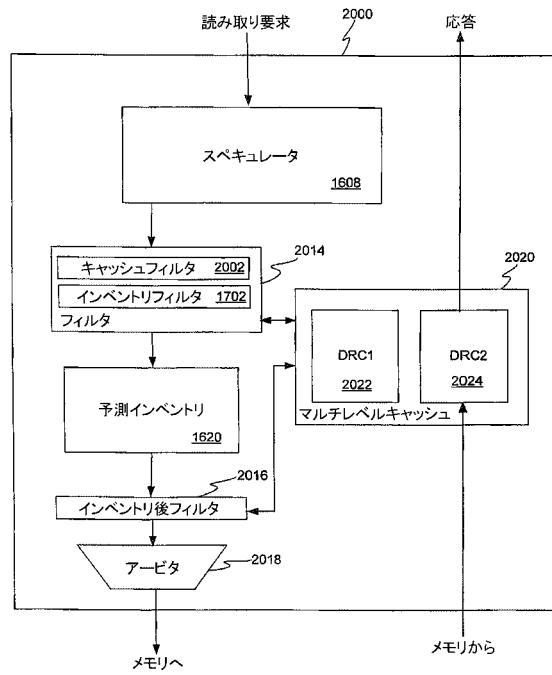
【図19A】



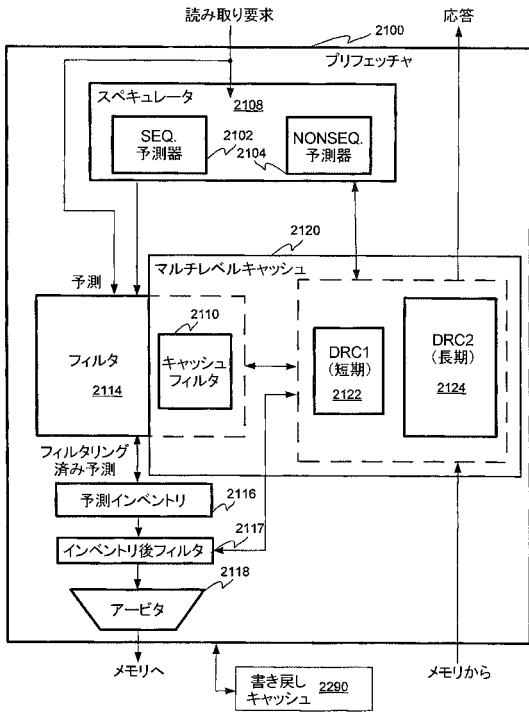
【図19B】



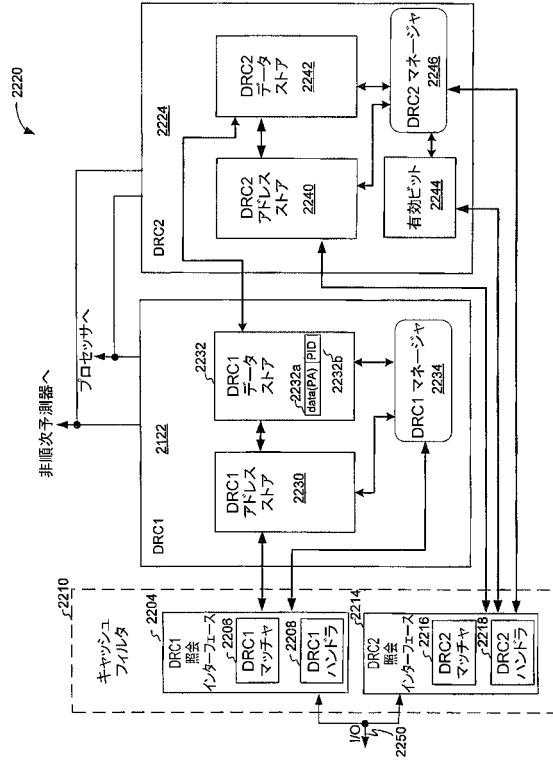
【図20】



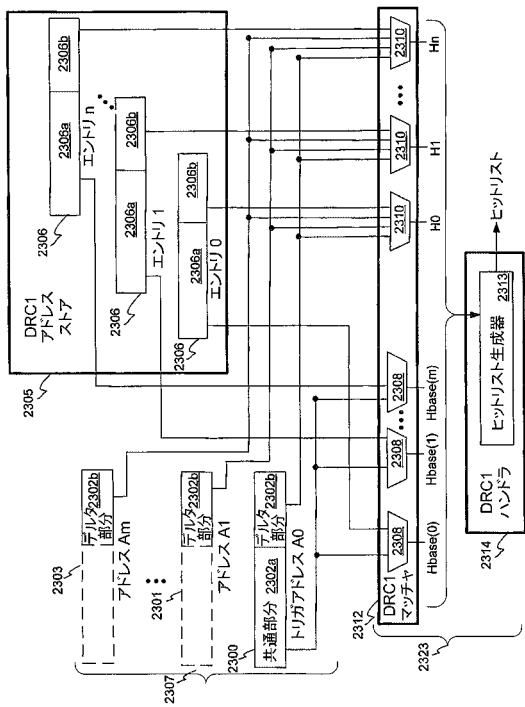
【図 2 1】



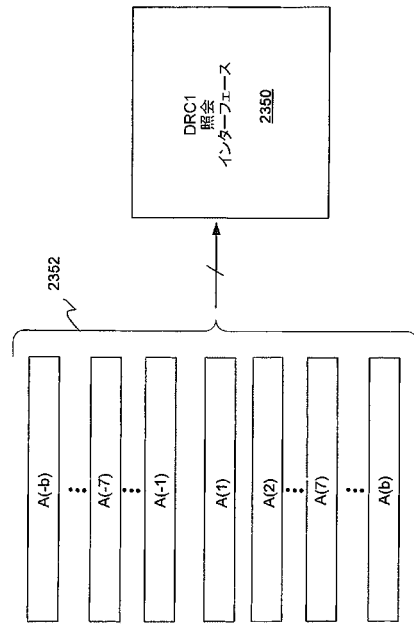
【図 2 2】



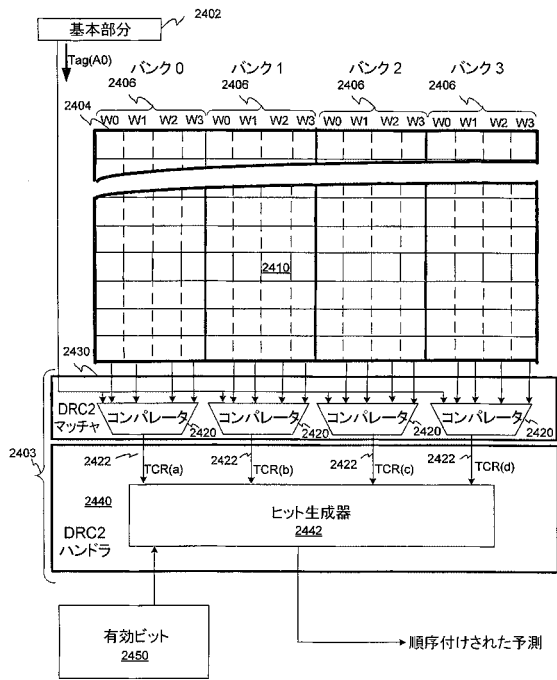
【図 2 3 A】



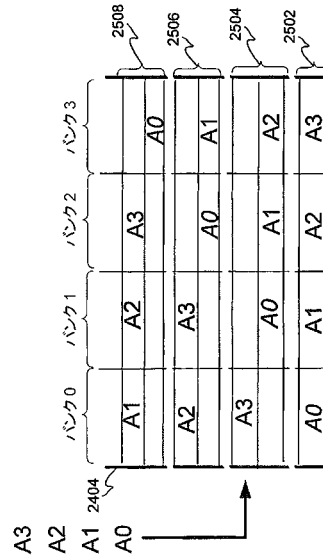
【図 2 3 B】



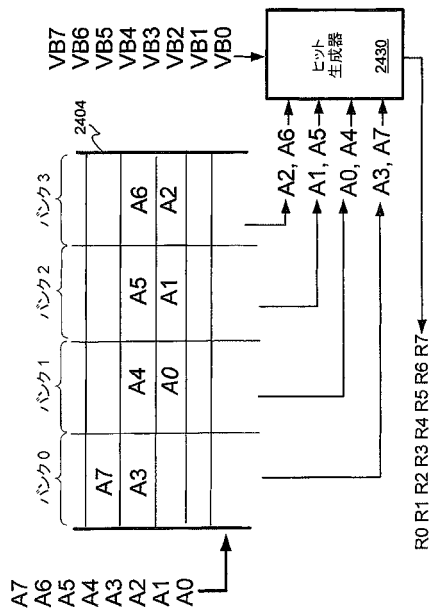
【図24】



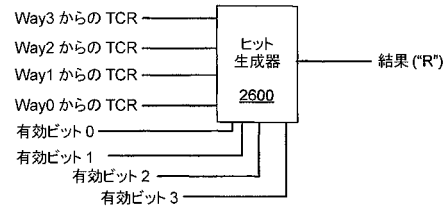
【図25A】



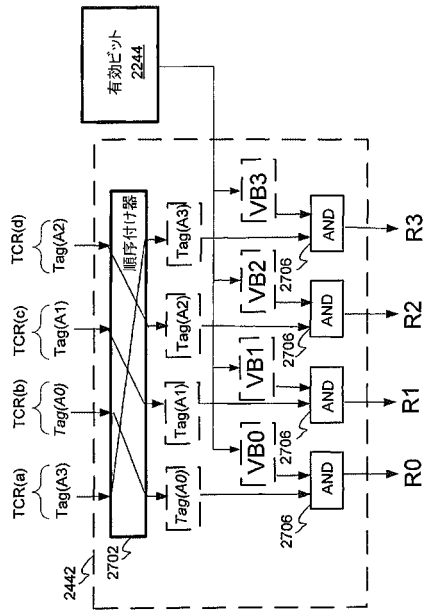
【図25B】



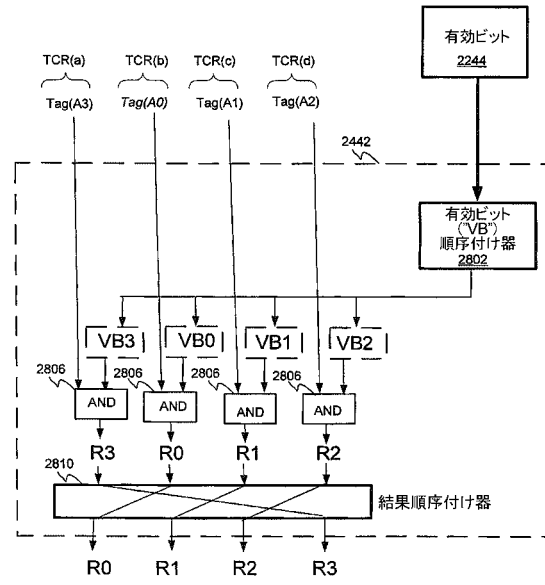
【図26】



【 図 27 】



【 図 28 】



## フロントページの続き

- (31)優先権主張番号 10/920,995  
(32)優先日 平成16年8月17日(2004.8.17)  
(33)優先権主張国 米国(US)
- (31)優先権主張番号 10/921,026  
(32)優先日 平成16年8月17日(2004.8.17)  
(33)優先権主張国 米国(US)
- (72)発明者 ダニラク, ラドスラフ  
アメリカ合衆国, カリフォルニア州, サンタ クララ, アパートメント ナンバー25,  
ブルネリッジ アヴェニュー 1850
- (72)発明者 シメラル, ブラッド, ダブリュー.  
アメリカ合衆国, カリフォルニア州, サン フランシスコ, コンド ナンバー3, ドロア  
ーズ ストリート 1049
- (72)発明者 ランゲンドーフ, ブライアン, キース  
アメリカ合衆国, カリフォルニア州, ベニシア, センプルズ クロッシング 272
- (72)発明者 ペスカドア, ステファノ, エー.  
アメリカ合衆国, カリフォルニア州, サニーヴェール, レイクミュアー ドライヴ 268
- (72)発明者 ヴィシェトスキー, ドミトリー  
アメリカ合衆国, カリフォルニア州, クパチーノ, シダー スプリング コート 1162  
7

審査官 中野 裕二

- (56)参考文献 特表2001-516919(JP,A)  
特開平08-212054(JP,A)  
特開平6-103169(JP,A)  
特開2001-166934(JP,A)  
特開平7-64862(JP,A)

## (58)調査した分野(Int.Cl., DB名)

G06F 12/08  
G06F 9/38