

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6236093号
(P6236093)

(45) 発行日 平成29年11月22日(2017.11.22)

(24) 登録日 平成29年11月2日(2017.11.2)

| (51) Int. Cl. | | | F I | | |
|---------------|-------------|------------------|------|------|------|
| G06F | 9/38 | (2006.01) | G06F | 9/38 | 370A |
| G06F | 9/32 | (2006.01) | G06F | 9/32 | 320A |
| G06F | 9/30 | (2006.01) | G06F | 9/30 | 350F |
| G06F | 9/45 | (2006.01) | G06F | 9/44 | 322G |

請求項の数 20 (全 20 頁)

| | | | |
|---------------|-------------------------------|-----------|---|
| (21) 出願番号 | 特願2015-555420 (P2015-555420) | (73) 特許権者 | 591016172 |
| (86) (22) 出願日 | 平成26年1月28日 (2014.1.28) | | アドバンスト・マイクロ・デバイス・ インコーポレイテッド |
| (65) 公表番号 | 特表2016-504699 (P2016-504699A) | | ADVANCED MICRO DEVI CES INCORPORATED |
| (43) 公表日 | 平成28年2月12日 (2016.2.12) | | アメリカ合衆国、94088-3453 |
| (86) 国際出願番号 | PCT/US2014/013455 | | カリフォルニア州、サニベイブル、ピー・ オウ・ボックス・3453、ワン・エイ・ エム・ディ・プレイス、メイル・ストップ ・68 (番地なし) |
| (87) 国際公開番号 | W02014/120690 | (74) 代理人 | 100108833 |
| (87) 国際公開日 | 平成26年8月7日 (2014.8.7) | | 弁理士 早川 裕司 |
| 審査請求日 | 平成28年2月3日 (2016.2.3) | (74) 代理人 | 100111615 |
| (31) 優先権主張番号 | 13/753,098 | | 弁理士 佐野 良太 |
| (32) 優先日 | 平成25年1月29日 (2013.1.29) | | |
| (33) 優先権主張国 | 米国 (US) | | |
| 早期審査対象出願 | | | |

最終頁に続く

(54) 【発明の名称】 並列パイプラインにおいてブランチを分岐するためのハードウェアおよびソフトウェアソリューション

(57) 【特許請求の範囲】

【請求項1】

コンピュータシステムの少なくとも1つのプロセッサによる実行のために構成された少なくとも1つのプログラムを格納するコンピュータ可読記憶媒体であって、

前記少なくとも1つのプログラムは、前記プロセッサによって実行されると、

複数のプログラム命令を解析することと、

前記複数のプログラム命令内でループおよび対応する基本ブロックを識別することと、

前記複数のプログラム命令内の識別されたループ内の所与の分岐点を識別したことに応じて、前記識別されたループ内の複数の命令を複数の超大命令語(VLIW)に配置することであって、少なくとも1つのVLIWは、前記所与の分岐点と対応する収束点との間の異なる基本ブロックから混ぜ合わされた命令を含む、ことと、

を前記プロセッサに実行させる命令を含む、

コンピュータ可読記憶媒体。

【請求項2】

前記少なくとも1つのプログラムは、前記識別されたループ内の前記所与の分岐点を検出したことに応じて、第1の命令を前記複数のプログラム命令に追加すること、を前記プロセッサに実行させるための命令をさらに含む、

前記第1の命令は、所与のVLIW内の命令を、単一命令複数データ(SIMD)マイクロアーキテクチャを含むターゲットプロセッサ内の複数の並列実行レーンに実行時に割り当てることを前記ターゲットプロセッサに実行させるためのものである、請求項1に記

載のコンピュータ可読記憶媒体。

【請求項 3】

前記第 1 の命令は、実行時に、前記所与の分岐点において所与のレーンに対して検出された分岐方向に少なくとも部分的に基づいて、前記複数の命令のうち何れかの命令を前記所与のレーンに割り当てることを前記ターゲットプロセッサに実行させる、請求項 2 に記載のコンピュータ可読記憶媒体。

【請求項 4】

前記少なくとも 1 つのプログラムの前記命令は、次のプログラムカウンタ (P C) に対応する V L I W の格納されたサイズを前記ターゲットプロセッサに更新させるように構成された第 2 の命令を追加することを前記プロセッサに実行させる、請求項 2 に記載のコンピュータ可読記憶媒体。

10

【請求項 5】

前記少なくとも 1 つのプログラムの前記命令は、前記ターゲットプロセッサに対して、前記検出された所与の分岐点と前記対応する収束点との間で第 1 の命令シーケンスが第 2 の命令シーケンスよりも小さいことに応じて、nop を、V L I W 内の前記第 2 の命令シーケンスに対応する命令とグループ化させるように構成された第 3 の命令を追加することを前記プロセッサに実行させる、請求項 4 に記載のコンピュータ可読記憶媒体。

【請求項 6】

実行時に、前記所与の V L I W 内の命令を前記複数の並列実行レーンに割り当てるために、前記第 1 の命令は、前記複数の実行レーンのうち対応するレーンと関連付けられているベクトルレジスタ内の特定のビット範囲にオフセットを書き込むように前記ターゲットプロセッサに実行させ、前記オフセットは、実行のために関連付けられたレーンに対してフェッチされた所与の V L I W 内の所与の命令を識別する、請求項 4 に記載のコンピュータ可読記憶媒体。

20

【請求項 7】

前記少なくとも 1 つのプログラムの前記命令は、
所与の命令シーケンスが、前記識別されたループの終わりに達していることに応じて、
実行時に、前記所与の命令シーケンスが、前記識別されたループの開始に分岐して戻るようにスケジューリングされているという判断に応じて、前記ベクトルレジスタ内の対応するビット範囲内にスリープ状態を書き込むことと、

30

実行時に、前記所与の命令シーケンスが、前記識別されたループの外部に分岐するようにスケジューリングされているという判断に応じて、前記ベクトルレジスタ内の前記対応するビット範囲内に終了状態を書き込むことと、

を前記ターゲットプロセッサに実行させるように構成された第 4 の命令を追加することを前記プロセッサに実行させる、

請求項 6 に記載のコンピュータ可読記憶媒体。

【請求項 8】

前記少なくとも 1 つのプログラムの前記命令は、
前記複数のレーンのうち、前記スリープ状態または前記終了状態にある前記所与の命令シーケンスに対応するレーンを検出したことに応じて、

40

前記所与の命令シーケンスの実行を停止することと、

少なくとも次のプログラムカウンタ (P C) および前記所与の命令シーケンスに対応する識別子 (I D) を格納することと、

を前記ターゲットプロセッサに実行させるように構成された第 5 の命令を追加することを前記プロセッサに実行させる、

請求項 7 に記載のコンピュータ可読記憶媒体。

【請求項 9】

前記少なくとも 1 つのプログラムの前記命令は、
レーンが前記スリープ状態にあるか前記終了状態にあることに応じて、前記レーン内の命令の実行を、個別に格納された次の P C で再開することを前記ターゲットプロセッサに

50

実行させるように構成された第 6 の命令を追加することを前記プロセッサに実行させる、請求項 8 に記載のコンピュータ可読記憶媒体。

【請求項 10】

前記少なくとも 1 つのプログラムの前記命令は、

スリープ状態にある命令のみの実行を、個別に格納された次の PC で再開することを前記ターゲットプロセッサに実行させるように構成された第 7 の命令を追加することを前記プロセッサに実行させる、請求項 8 に記載のコンピュータ可読記憶媒体。

【請求項 11】

単一命令複数データ (SIMD) マイクロアーキテクチャ内の複数の並列実行レーンと、

超大命令語 (VLIW) のサイズを格納するように構成されたサイズレジスタと、プログラムコードを実行するように構成された制御ロジックと、を備え、

前記プログラムコードは、

複数のプログラム命令内でループおよび対応する基本ブロックを識別することと、前記複数のプログラム命令内の識別されたループ内の所与の分岐点を識別したことに応じて、前記識別されたループ内の複数の命令を複数の超大命令語 (VLIW) に配置することと、少なくとも 1 つの VLIW は、前記所与の分岐点と対応する収束点との間の異なる基本ブロックから混ぜ合わされた命令を含む、ことと、により生成される、複数の VLIW を含むものである、

プロセッサ。

【請求項 12】

前記複数の並列実行レーンのうち対応するレーンと関連付けられている特定のビット範囲内にオフセットを格納するように構成されたベクトルレジスタをさらに備え、

前記オフセットは、実行のために関連付けられたレーンに対してフェッチされた所与の VLIW 内の所与の命令を識別する、請求項 11 に記載のプロセッサ。

【請求項 13】

前記ベクトルレジスタの前記ビット範囲に格納される前記オフセットの有効な値としてありうる、互いに異なる値の数は、前記サイズレジスタに格納されたサイズと等しい、請求項 12 に記載のプロセッサ。

【請求項 14】

オフセットが、前記 VLIW 内の複数の資源に依存しない命令に対応しているとの検出に応じて、前記複数の並列実行レーンのうち前記オフセットに関連付けられたレーンは、関連付けられたレーン内の前記複数の命令を同時に実行するようにさらに構成されている、請求項 12 に記載のプロセッサ。

【請求項 15】

所与の命令シーケンスが前記識別されたループの終わりに達しており、且つ、前記所与の命令シーケンス、及び、対応するレーンが、前記識別されたループの開始に分岐して戻るようにスケジューリングされていることを示すスリープ状態、または、前記識別されたループの外側に分岐するようにスケジューリングされていることを示す終了状態のいずれかであることに応じて、前記制御ロジックは、

前記所与の命令シーケンスの実行を停止することと、

少なくとも次のプログラムカウンタ (PC) および前記所与の命令シーケンスに対応する前記レーンの識別子 (ID) を格納することと、

を行うようにさらに構成されている、請求項 12 に記載のプロセッサ。

【請求項 16】

前記複数の並列実行レーンの各レーンが、前記スリープ状態または前記終了状態にあることに応じて、前記制御ロジックは、個別に格納された次の PC に分岐することにより、レーンごとの実行を再開するようにさらに構成されている、請求項 15 に記載のプロセッサ。

10

20

30

40

50

【請求項 17】

前記複数の並列実行レーンの各レーンが、前記スリープ状態または前記終了状態にあることに応じて、前記制御ロジックは、個別に格納された次のPCに分岐することにより、スリープ状態にあるレーンのみの実行を再開するようにさらに構成されている、請求項15に記載のプロセッサ。

【請求項 18】

複数のプログラム命令内でループおよび対応する基本ブロックを識別することと、識別されたループ内の所与の分岐点に応じて、前記識別されたループ内の複数の命令を複数の超大命令語(VLIW)に配置することとであって、少なくとも1つのVLIWは、前記所与の分岐点と対応する収束点との間の異なる基本ブロックから混ぜ合わされた命令を含む、ことと、
を含む、方法。

10

【請求項 19】

前記識別されたループ内の前記所与の分岐点に応じて、前記所与の分岐点において所与のレーンの実行時に検出した分岐方向に基づいて、実行時に、所与のVLIW内の命令を、単一命令複数データ(SIMD)マイクロアーキテクチャを含むターゲットプロセッサ内の複数の並列実行レーンに割り当てることをさらに含む、請求項18に記載の方法。

【請求項 20】

実行時に、前記所与のVLIW内の命令を前記複数の並列実行レーンに割り当てるために、前記複数の並列実行レーンのうち対応するレーンと関連付けられている指示を格納することをさらに含む、

20

前記指示は、実行するために関連付けられたレーンに対して前記所与のVLIW内の所与の命令を識別する、請求項19に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、コンピューティングシステムに関し、より詳細には、プロセッサ内のハードウェア並列実行レーンで命令を効率的に処理することに関する。

【背景技術】

【0002】

コンピュータシステムのスループットを向上させるために、タスクの並列化が用いられている。この目的のために、コンパイラは、並列化されたタスクをプログラムコードから抽出して、システムハードウェア上で並行して実行し得る。ハードウェア上での並列実行を向上するために、プロセッサは、複数の並列実行レーン(例えば、単一命令複数語(SIMD)マイクロアーキテクチャ内など)を含み得る。このタイプのマイクロアーキテクチャは、特定のソフトウェアアプリケーションに対して、単一レーンのマイクロアーキテクチャまたは汎用マイクロアーキテクチャよりも高い命令スループットを提供し得る。SIMDマイクロアーキテクチャから恩恵を受けるタスクのいくつかの例は、ビデオグラフィックスレンダリング、暗号化、およびガベージコレクションを含む。

30

【0003】

多くの場合、特定のソフトウェアアプリケーションは、各作業項目の実行や並列関数呼出しがその内部のデータ依存となるようなデータ並列処理を有する。例えば、第1の作業項目が第2の作業項目から独立したデータである場合には、第1および第2の作業項目の各々は、SIMDマイクロアーキテクチャ内の別個の並列実行レーン上に同時にスケジューリングされている。しかし、第1および第2の作業項目の各々で実行されるある量の命令は、データ依存の場合がある。分岐命令として実装された条件テストは、第1の作業項目に対してパスし得るが、各作業項目に対するデータに依存する第2の作業項目に対しては不合格になり得る。

40

【0004】

第2の作業項目が、実行を停止して、第1の作業項目が進行中の実行を継続するのを待

50

機するので、並列実行の効率が低下し得る。パスしたテストに起因して2～3の作業項目だけが実行を継続し、他方、不合格になったテストのために、ほとんどの作業項目がアイドルである場合には、効率の悪さが増大する。

【発明の概要】

【課題を解決するための手段】

【0005】

プロセッサ内のハードウェア並列実行レーンで命令を効率的に処理するためのシステムおよび方法を検討する。様々な実施形態では、バックエンドコンパイラは、ソフトウェアアプリケーションのプログラム命令を検査して、ターゲットプロセッサ上で命令を効率的に処理するように、命令を配置してコードを生成する。ターゲットプロセッサは、単一命令複数データ(SIMD)マイクロアーキテクチャ内に複数の並列実行レーンを含んでもよい。コンパイラは、ループと、対応する基本ブロックとを識別してもよい。ループ内の分岐点は、分岐命令を含んでもよい。例えば、if-else-if-else構成、if-else構成、case構成などが、識別されたループ内のプログラム命令で使用されてもよい。分岐点と対応する収束点との間での、翻訳されコンパイルされたプログラム命令の実行中に、複数のトレースパスがトラバースされてもよい。

10

【0006】

コンパイル中、識別されたループ内の所与の分岐点に応じて、コンパイラは、識別されたループ内の命令を、1つ以上の超大命令語(VLIW: very large instruction word)に配置してもよい。少なくとも1つのVLIWは、所与の分岐点と対応する収束点との間の異なる基本ブロックから混ぜ合わされた命令を含んでもよい。例えば、4つの命令を有する基本ブロックAと、6つの命令を有する基本ブロックBとが、所与の分岐点と対応する収束点との間に存在する場合には、コンパイラは、命令を6つのVLIW内に配置してもよい。最初の4つのVLIWは、基本ブロックAおよび基本ブロックBの各々から1つの命令を含んでもよい。第1のVLIWは、基本ブロックAおよび基本ブロックBの各々からの第1の命令を含んでもよい。第2のVLIWは、基本ブロックAおよび基本ブロックBの各々からの第2の命令を含んでもよく、以下同様である。最後の2つのVLIWは、nop(ノーオペレーション)とともにグループ化された基本ブロックBからの命令を含んでもよい。コンパイラは、各VLIWをポイントするプログラムカウンタ(PC)値を追跡してもよい。

20

30

【0007】

コンパイラは、翻訳されコンパイルされた命令で、挿入するためのコードを生成してもよい。挿入されたコードは、実行される際に、所与のVLIW内の命令を、ターゲットプロセッサ内の複数の並列実行レーンに実行時に割り当ててもよい。所与のレーンに対する割り当ては、実行時に、所与の分岐点で、所与のレーンに対して見つかった分岐方向に基づいてもよい。前述の例を続けると、VLIWが、基本ブロックAおよび基本ブロックBから生成された第2のVLIWであって、所与のレーンに対する分岐命令が選択された場合には、所与のレーンは、第2のVLIW内の基本ブロックA内の第2の命令が割り当てられてもよい。分岐命令が選択されない場合には、所与のレーンは、第2のVLIW内の基本ブロックB内の第2の命令が割り当てられてもよい。様々な実施形態では、VLIWは可変長である。挿入されたコードは、実行される際に、次のPCに対応するVLIWのサイズを更新してもよい。

40

【0008】

いくつかの実施形態では、プロセッサは、単一命令複数データ(SIMD)マイクロアーキテクチャ内に複数の並列実行レーンを含む。プロセッサは、可変長VLIWのサイズを格納するためのサイズレジスタを含んでもよい。プロセッサ内の制御ロジックは、格納したサイズに等しい所与のVLIW内のいくつかの命令を、それぞれのサイクル内でフェッチして復号してもよい。複数の実行レーンは、所与のVLIW内でいくつかの命令を同時に実行してもよい。さらに、プロセッサは、複数の実行レーンのうち対応するレーンと関連付けられたビット範囲を有するベクトルレジスタを含んでもよい。ビット範囲はオフ

50

セットを格納してもよい。所与のオフセットは、実行するために関連付けられたレーンに対してフェッチされたVLIW内の所与の命令を識別してもよい。

【0009】

これらの実施形態および他の実施形態は、以下の記載および図面を参照することによって、さらに理解されるであろう。

【図面の簡単な説明】

【0010】

【図1】単一命令複数データ(SIMD)パイプライン実行フローの一実施形態の一般化されたブロック図である。

【図2】制御フローグラフの一実施形態の一般化されたブロック図である。

10

【図3】制御フローグラフに対する実行順序の一実施形態の一般化されたブロック図である。

【図4】プロセッサに対するSIMDマイクロアーキテクチャの論理レイアウトの一実施形態の一般化されたブロック図である。

【図5】オブジェクトコード配置の一実施形態の一般化されたブロック図である。

【図6】コンパイラ技術を用いて、プロセッサ内で複数の作業項目の並列実行を最適化するための方法の一実施形態の一般化されたフロー図である。

【図7】ハードウェア技術を用いて、プロセッサ内で複数の作業項目の並列実行を最適化するための方法の一実施形態の一般化されたフロー図である。

【図8】オブジェクトコード配置の別の実施形態の一般化されたブロック図である。

20

【発明を実施するための形態】

【0011】

実施形態は、様々な修正および代替形式を受け入れる余地があるが、特定の実施形態が例として図面に示され、本明細書で詳細に説明されている。しかし、図面およびそれに関する詳細な説明は、実施形態が、開示された特定の形式に限定されることを意図しておらず、それとは逆に、添付の請求項によって定義されるように、実施形態の趣旨および範囲に含まれる全ての修正、均等物および代替手段を包含することが理解されるべきである。

【0012】

以下の記載では、実施形態の完全な理解を提供するために、多数の具体的な詳細が記載されている。しかし、当業者は、これらの具体的な詳細なしに、実施形態が実施され得ることを理解するはずである。いくつかの場合には、実施形態を曖昧にするのを避けるために、周知の回路、構造および技術が詳細に示されていない。

30

【0013】

図1を参照すると、単一命令複数データ(SIMD)パイプライン実行フロー100の一実施形態を示す一般化されたブロック図が示されている。命令102~108はフェッチされ、関連するデータと共にSIMDパイプラインに送信されてもよい。並列で垂直な実行レーン内の複数の計算ユニットが示されている。いくつかの計算ユニットは、アクティブな計算ユニット110である。他の計算ユニットは、所与のパイプステージ中に無効にされていることに起因する、非アクティブな計算ユニット112である。制御ロジックおよび記憶素子(例えば、パイプラインレジスタなど)は、説明を容易にするために示されていない。

40

【0014】

ハードウェア計算ユニットは、関連付けられたデータを使用して、所与の作業項目の所与の命令の実行を行うハードウェアを含む。このハードウェアは、加算、乗算、ゼロ検出、ビット単位シフト、除算、ビデオグラフィックスおよびマルチメディア命令、または、プロセッサ設計の当業者に周知の他の操作を実行するように構成された演算論理装置を含んでもよい。SIMDパイプライン内に並列実行レーンを有するプロセッサの例は、グラフィック処理装置(GPU)、デジタル信号処理(DSP)などを含む。一実施形態では、SIMDパイプラインは、ビデオカード上に配置されてもよい。別の実施形態では、SIMDパイプラインは、マザーボード上に統合されてもよい。

50

【 0 0 1 5 】

S I M Dパイプラインは、ゲーム、エンタテインメント、科学および医療分野で使用される多種多様なデータ並列アプリケーションに対する計算性能を向上させ得る。かかるアプリケーションは、一般に、多数のオブジェクトについて同じプログラムを実行することを伴う。各オブジェクトは、他のオブジェクトと関係なく処理されるが、同じ順序の操作が使用されるので、S I M Dマイクロアーキテクチャは、相当な性能強化を提供する。G P Uは、非グラフィック計算用にも想定されている。

【 0 0 1 6 】

ソフトウェアアプリケーションは、関数呼出しまたは計算カーネルの集合、および、内部関数の集合を含んでもよい。ソフトウェアプログラムは、関数呼出しを定義してもよく、他方、内部関数は所与のライブラリ内で定義されてもよい。例えば、ソフトウェアアプリケーションは、例えば画像ファイルなどの2次元(2D)配列のデータのデータ処理を実行し得る。ソフトウェアアプリケーションは、ソフトウェアプログラムによって開発されたアルゴリズムを、2D画像の画素ごとまたは2次元行列の要素ごとに、実行し得る。所与の関数呼出しは、インデックス空間を介して呼び出されてもよい。インデックス空間は、次元空間とも呼ばれ得る。データ並列ソフトウェアアプリケーションに対して、N次元計算領域は、1、2もしくは3次元空間、または、インデックス空間を定義してもよい。一例は、2D画像内の画素である。

【 0 0 1 7 】

関数呼出しは、データの1つ以上のレコードと照合されて、1つ以上の計算の作業項目を生成してもよい。従って、2つ以上の作業項目は、単一の関数呼出しの同じ命令を利用し得るが、データの異なるレコードについて動作し得る。関数呼出しは、フォーク(fork)を生成する制御フロー転送命令を含んでもよく、他方、コンピュータプログラム内のフォークは、通常、共通の定義によってソフトウェアスレッドを生成する。インデックス空間内の所与の時点における関数呼出しの所与のインスタンスが、「作業項目」と呼ばれてもよい。作業項目は、作業ユニットとも呼ばれてよい。前述の例を続けると、作業項目は、2D画像の所与の画素(所与のインデックス)に対応するデータのレコードについて、関数呼出し内の1つ以上の命令で動作してもよい。通常、作業項目は、関連付けられた一意の識別子(I D)を有する。

【 0 0 1 8 】

インデックス空間は、十分なハードウェアサポートがある場合に並行して実行する作業項目の総数を定義してもよい。例えば、インデックス空間は、280の数の作業項目を定義してもよいが、G P Uは、いつでも64の作業項目の同時実行をサポートし得る。作業項目の総数は、グローバルな作業サイズを定義し得る。作業項目は、さらに作業グループにグループ化され得る。各作業グループは、一意の識別子(I D)を有してもよい。所与の作業グループ内の作業項目は、相互に通信して、実行を同期させ、メモリアクセスを調整することが可能であってよい。いくつかの作業項目は、S I M D方式でG P U上での同時実行のためにウェーブフロント(wave front)にクラスタ化されてもよい。280の総作業項目に対する前述の例に関して、ウェーブフロントは64の作業項目を含んでもよい。

【 0 0 1 9 】

命令102~108は、フェッチされて、関連付けられたデータと共にS I M Dパイプラインに入ってもよい。命令104は、例えば条件分岐などの制御フロー転送命令であってよい。命令106は、条件が真の場合に実行されるパス内の第1の命令であってよい。命令108は、条件が偽の場合に実行されるパス内の第1の命令であってよい。例えば、分岐命令104は、高水準言語プログラムにおけるI F文と関連付けられてもよい。命令106は、高水準言語プログラムにおけるT H E N文と関連付けられてもよい。命令108は、高水準言語プログラムにおけるE L S E文と関連付けられてもよい。

【 0 0 2 0 】

所与の行内の各計算ユニットは、同じ計算ユニットであってよい。これらの計算ユニッ

10

20

30

40

50

トの各々は、同じ命令であるが、異なる作業項目と関連付けられた異なるデータについて動作してもよい。図に示すように、いくつかの作業項目は、条件分岐命令104によって提供されたテストをパスし、他の作業項目はテストに不合格になる。SIMDパイプライン内の制御ロジックは、利用可能なパスの各々を実行して、現在のパスを選択しなかった作業項目に対応する実行ユニット（例えば、計算ユニットなど）を選択的に無効にしてもよい。例えば、If-Then-Else構文の実行中、SIMDアーキテクチャの各列内には、「Then」（パスA）および「Else」（パスB）のパスを実行するように構成された実行ユニットがある。

【0021】

第1および第2の作業項目が実行を停止して、第3の作業項目が進行中の実行を継続するのを待機すると、並列実行の効率が低下し得る。従って、分岐命令104の実行後、所与の行における全ての計算ユニットがアクティブな計算ユニット110というわけではない。図に示すように、1つ以上の計算ユニットは、実行に関して無効にされた非アクティブな計算ユニット112である。多数の計算ユニットが所与のパイプステージ中に非アクティブである場合には、SIMDコアの効率およびスループットが低下する。一実施形態では、「Else」パスは、関数呼出しに対するリターン（return）である。関数呼出しの実行が終了して、対応する作業項目がアイドルになる。しかし、SIMDコア内の隣接する作業項目は、実行を継続してもよい。

【0022】

ここで図2を参照すると、制御フローグラフ200の一実施形態を示す一般化されたブロック図が示されている。一般的に言えば、制御フローグラフは、コンパイラ最適化器および静的解析ツールによって使用され得る。制御フローグラフ200は、プログラムまたはプログラムの一部が、その実行中にトラバースされ得る全てのパスを表し得る。制御フローグラフでは、グラフ内の各ノードは、基本ブロックを表している。ほとんどの表現は、制御が制御フローグラフに入るための入口ブロックと、制御が制御フローグラフを出るための出口ブロックと、を含む。

【0023】

コンパイル中、ソフトウェアアプリケーションは、基本ブロック0（BB 0）から基本ブロック7（BB 7）まで番号付けされた8つの基本ブロック（BB）を有する制御フローグラフ200を提供し得る。8つの基本ブロックが示されているが、他の例では、別の数の基本ブロックが使用され得る。制御フローグラフ200において、基本ブロック1が入口ブロックであり、基本ブロック6が出口ブロックである。基本ブロック0～7の各々は、1つの入口点と1つの出口点とを有する命令のストレートラインシーケンスである。制御フローグラフ200は、ループを表してもよい。ループの内部では、制御フローグラフ200は、基本ブロック1～4を有するIF-THEN-ELSE構成と、基本ブロック4～6を有するIF構成と、を表してもよい。

【0024】

ここで図3を参照すると、制御フローグラフに対する実行順序300の一実施形態を示す一般化されたブロック図が示されている。実行順序310は、既に示した制御フローグラフ200がSIMDパイプラインに割り当てられた場合の典型的な実行順序を表している。単一のループの繰返しに対する実行時間は、ループ内の各基本ブロック（例えば、BB 1～BB 6など）の実行時間の合計である。しかし、所与の作業項目および対応するハードウェア実行レーンに対して、BB 2およびBB 3のうち1つのみが実行される。同様に、BB 5が所与の作業項目に対してスキップされ得る。特定の基本ブロックが所与の作業項目に対して実行されない可能性があるが、関連付けられた実行時間は、ループの繰返しに対する実行時間の総合計に寄与する。

【0025】

実行順序320は、既に示した制御フローグラフ200が、修正されたSIMDパイプラインに割り当てられた場合の代替の実行順序を表している。単一のループの繰返しに対する実行時間は、ループ内の各基本ブロックの実行時間の合計ではなく、単一のループの

10

20

30

40

50

繰返し内で実際に実行された基本ブロックの実行時間の合計である。実行順序320は、制御フローグラフ200の実行を変換する。コンパイラは、ソースコードのオブジェクトコードへのコンパイル中に、この変換を実行し得る。

【0026】

いくつかの実施形態では、コンパイラは、各基本ブロックの終わりにコードを生成し、そのコードは、実行される際に、次に実行する基本ブロックを識別する。生成コードは、基本ブロックの終わりに挿入され得る。あるいは、制御を次の基本ブロックに転送する前に、制御フローを追加のコードに転送するために、分岐命令が基本ブロックの終わりに挿入され得る。追加のコードは、中央基本ブロックを示すBB Cによって表される。実行時、分岐する基本ブロックの各々(例えば、BB 1、BB 4、BB 6など)は、制御の転送先となる次の基本ブロックを識別する。識別は、分岐解決に基づくものであり、分岐解決は、データの特定レコードおよび分岐命令にさらに基づいている。実行時、BB Cは、ターゲットの基本ブロックのアドレスを各作業項目から受信して、それぞれのターゲットの基本ブロックを実行するためのスレッドレジスタをセットアップする。SIMDパイプライン内で実行されている異なる作業項目にわたって、所与の作業項目は、分岐、ジャンプおよびケース文などのように、制御フロー転送命令に対する単一のターゲットを有する。

10

【0027】

実行順序320では、BB 4の完了時に、第1の作業項目は、BB 5に分岐して、対応するアドレスをBB Cに渡し得る。BB 4の完了時に、第2の作業項目は、BB 6に分岐して、対応するアドレスをBB Cに渡し得る。コンパイラは、同時に実行される各ターゲット基本ブロックからの命令を含む、動的超大命令語(DVLIW)を生成し得る。実行される際に、コンパイラによって生成されたコードは、次のプログラムカウンタ(PC)値に応じて実行する次のDVLIWのサイズを更新し得る。加えて、実行される際に、生成されたコードは、所与の並列実行レーン内で実行されている所与の作業項目と、実行する次のDVLIW内の命令へのポインタとの間のマッピングを更新し得る。ポインタは、次のDVLIW内の命令のうち関連付けられた命令であって、フェッチされる命令を識別しているオフセットであってよい。

20

【0028】

図4を参照すると、プロセッサに対するSIMDマイクロアーキテクチャの論理レイアウト400の一実施形態を示す一般化されたブロック図が示されている。プロセッサは、データおよび命令を格納するためのダイナミックランダムアクセスメモリ(DRAM)450を有する。いくつかの実施形態では、所与のレベルのキャッシュメモリサブシステムが、DRAMに加えて使用される。図に示すように、プロセッサは、計算ユニットの行ごとに、制御ロジック420と一緒にグループ化された比較的小規模のキャッシュメモリサブシステム430を有してもよい。説明を簡略にするために、パイプラインレジスタなどの記憶素子は示されていないが、プロセッサ内のデータフローは、パイプライン化されてもよい。所与のパイプラインのステージでは、このステージ内の関連付けられた命令が、既に不合格になったテスト(例えば、選ばれなかった分岐など)に基づいて実行されない場合には、計算ユニットは使用されない可能性がある。

30

40

【0029】

SIMDパイプラインは、レーンA~Fを有する作業項目460を含む。レーンA~Fの各々は、計算ユニットを含む垂直で並列なハードウェアレーンの各々に対応し得る。さらに、パイプラインは、ベクトルレジスタ462を含んでもよい。ベクトルレジスタ462は、並列実行レーンの各々に対して、エントリ、フィールドまたはビット範囲を含んでもよい。各エントリは、それぞれの作業項目上で実行している所与のトレースを識別するための第1のビット数と、特別コードをサポートするための第2のビット数と、を含むビット総数を含んでもよい。特別コードは、待ち状態またはスリープ状態、ループ終了状態、ループを終了するため以外の実行を停止するためのバリア識別子、イベント識別子などを識別してもよい。特別コードが所与のエントリ内に格納されない場合には、格納された

50

値は、実行するために関連付けられたレーンについてD V L I W内の命令の各々を識別してもよい。

【 0 0 3 0 】

プログラムカウンタ (P C) レジスタ 4 6 6 は、 i - キャッシュなどのメモリからフェッチするために、次のD V L I Wをポイントしているポインタ値またはアドレスを格納してもよい。プロセッサは、D V L I Wのサイズ、すなわち長さを格納するサイズレジスタ 4 6 8 をさらに含んでもよい。いくつかの実施形態では、サイズは、可変長D V L I W内の命令数を表す整数であってよい。

【 0 0 3 1 】

D V L I W 4 6 4 内の命令 I n s t r A ~ I n s t r G の各々は、制御フローグラフ内の実行トレースを表している。コンパイラは、D V L I Wを、 i - キャッシュなどのメモリ内に配置してもよい。一例では、作業項目 4 6 0 内のレーン B は、S I M D パイプライン内の左から 2 番目の垂直な実行レーンに対応し得る。ベクトルレジスタ 4 6 2 内に格納されたオフセット B は、レーン B と関連付けられて、D V L I W 4 6 4 内の最初の命令をポイントし得るが、その命令は、I n s t r A である。従って、レーン B は、I n s t r A を受信して処理し得る。同様に、作業項目 4 6 0 内のレーン A は、S I M D パイプライン内の最も左の垂直な実行レーンに対応し得る。ベクトルレジスタ 4 6 2 内に格納されたオフセット A は、レーン A と関連付けられて、D V L I W 4 6 4 内の最後の命令 (I n s t r G) をポイントし得る。従って、レーン A は、I n s t r G を受信して処理し得る。

【 0 0 3 2 】

図示していないが、命令キャッシュ (i - キャッシュ) は、D V L I W をサポートするための複数の実施態様のうち 1 つを含んでもよい。 i - キャッシュは、D V L I W に対応する所与の単一の P C に対する 1 つ以上の命令をフェッチするための複数の小型のキャッシュを含んでもよい。同じ P C は、D V L I W のサイズに応じて、小型のキャッシュのうち 1 つ以上のキャッシュ内の有効な命令をインデックスしてもよい。 i - キャッシュは、P C レジスタ 4 6 6 内に格納されたポインタまたはアドレス値に加えて、サイズレジスタ 4 6 8 内に格納されたサイズを受信し得る。代替として、 i - キャッシュは、同じ有効なキャッシュラインまたはキャッシュセット内の 1 つ以上の命令にアクセスするための複数のデータポートを有してもよい。この場合もやはり、フェッチする有効な命令の数は、サイズレジスタ 4 6 8 からの受信したサイズに等しい可能性がある。

【 0 0 3 3 】

ここで図 5 を参照すると、オブジェクトコード配置 5 0 0 を示す一実施形態の一般化されたブロック図が示されている。コード配置 5 0 0 は、図 2 および図 3 にそれぞれ示す、制御フローグラフ 2 0 0 および付随の実行順序 3 2 0 に対してコンパイラにより生成され得るオブジェクトコードレイアウトの一実施形態を示している。基本ブロックコード 5 0 4 は、各基本ブロックに対するコードを表す。例えば、基本ブロック 0 ~ 3 , 7 に対するコードの 1 つのコピーが、レイアウト 5 0 0 内に配置されて示されている。基本ブロック 4 ~ 5 に対するコードの 2 つのコピーが、レイアウト 5 0 0 内に配置され示されている。基本ブロック 6 に対するコードの 4 つのコピーが、レイアウト 5 0 0 内に示されている。

【 0 0 3 4 】

コード 5 0 2 は、ループなどの領域のエントリに対して、コンパイラによって生成および挿入され得る。後に実行される場合、コード 5 0 2 は、次のD V L I Wのサイズを更新し、フェッチされたD V L I W内の命令と、ターゲットプロセッサ内の並列実行レーンとの間のマッピングを更新し得る。例えば、ターゲットプロセッサ内のサイズレジスタおよびベクトルレジスタは、それらに格納された内容を、実行されたコード 5 0 2 によって更新させてもよい。図に示すように、コード 5 0 2 は、基本ブロック 0 , 1 , 7 の開始時に挿入され得る。

【 0 0 3 5 】

コード 5 0 6 は、基本ブロック 1 , 4 , 6 などの分岐点に移行するために、コンパイラ

によって生成および挿入され得る。後に実行される場合、コード506は、DVLIWのサイズ変更、および、フェッチされたDVLIW内の命令とターゲットプロセッサ内の並列実行レーンとの間の対応するマッピング変更を判断し得る。従って、サイズおよびマッピングが、制御フローグラフ内の分岐点および収束点において更新される。コンパイラは、DVLIWのサイズが変わる点、および、マッピングが変わる点を識別する。第1のインデックスがトレース識別子(ID)を示し、第2のインデックスが基本ブロック(BB)IDを示す表記法BBC(0,1)を使用すると、コード506を挿入するための識別された点は、BBC(0,1)、BBC(0,4)、BBC(0,6)、BBC(1,4)、BBC(1,6)、BBC(2,6)およびBBC(3,6)に存在し得る。この例におけるトレースIDは、対応するオフセットと同じであってよい。

10

【0036】

オブジェクトコード配置500の開始時に、コード502は、初期化ステップを実行して、DVLIWサイズを1に設定してもよい。ベクトルレジスタの各エン트리内のオフセットは、例えば0のオフセットなどのように、BB0内の同じ命令をポイントするように設定され得る。従って、PCは、0または別の適切な開始アドレスに設定され得る。図に示すように、4つの可能な並列トレースが存在しているが、作業項目の数は独立であってよい。例えば、SIMDパイプラインは、割り当てられた作業項目を処理するための、8、16、64または別の数の並列実行レーンを有してもよい。SIMDパイプライン内の各作業項目は、ベクトルレジスタ内に格納された0のオフセットを有しており、同じ命令を実行し得る。各作業項目に対するこの同じ命令は、BB0からの命令である。BB0内の命令は、各作業項目によって1つずつ実行され、各命令フェッチの後にPCが増加する。

20

【0037】

BB0の実行が完了した後、ループ入口ブロックであるBB1が次に処理される。コード502は、BB1の開始時に、DVLIWサイズを1として保持し、各作業項目に対するオフセットを0として保持する。ベクトルレジスタのエントリの各々のオフセットは、例えば0のオフセットなどのように、BB1内の同じ命令をポイントするように設定され得る。PCは、BB0の完了時に増加された値のままであってよい。SIMDパイプライン内の各作業項目は、ベクトルレジスタ内に格納された0のオフセットを有しており、同じ命令を実行することになる。作業項目ごとのこの同じ命令は、BB1からの命令である。BB1内の命令は各作業項目によって1つずつ実行され、各命令のフェッチ後にPCが増加する。

30

【0038】

実行される際に、コード506は、BB1の終わりにあるBBC(0,1)において、格納されたDVLIWサイズを1から2に変更する。ここで、BB3内の命令は、フェッチされたDVLIWに追加される。さらに、実行される際に、コード506は、BBC(0,1)にて、BB3に分岐する作業項目に対するベクトルレジスタ内のエント리를、値1を格納するように設定する。BB2に分岐する作業項目に対するベクトルレジスタ内のエント리는、0を継続して格納することにより、変更されないままである。0および1の値がこのように使用されるが、対応する指示およびマッピングを設定するために他の数値が使用されてもよい。この時点で、DVLIWは2つの命令を有し、これらは、2つの別個の基本ブロックBB2およびBB3から混ぜ合わされている。PCが継続して増加されるので、フェッチされたDVLIWは、BB2の処理が完了するまで、これらの2つの基本ブロックからの混ぜ合わされた命令を含み続ける。コンパイラは、DVLIW内の命令の並列実行をサポートするために、メモリ内でこのように混ぜ合わされるように命令を配置してもよい。

40

【0039】

BB2の完了時に、DVLIWサイズは2のままである。ベクトルレジスタ内に格納されたオフセットもそれらの値のままである。しかし、この時、オフセット0は、BB2ではなく、BB4内の命令に対応する。BBC(0,4)の完了時に、コード506

50

は、トレース 0 内の B B 4 の終わりにおいて、3 を格納するようにサイズレジスタを更新し、B B 6 に分岐する作業項目に対するエントリを、2 を格納するように更新する。単一の P C および格納されたサイズを i - キャッシュに送信した後に、長さ 3 の D V L I W が i - キャッシュからフェッチされる。D V L I W は、B B 3 または B B 4 と、B B 5 と、B B 6 とから混ぜ合わされた命令を含む。ベクトルレジスタ内に関連する格納された 0 のオフセットを有する作業項目は、B B (0 , 5) からフェッチされた命令を得る。ベクトルレジスタ内に関連する格納された 1 のオフセットを有する作業項目は、どの程度まで P C が増加されているかに応じて、B B (1 , 3) または B B (1 , 4) の何れかから結果を得る。ベクトルレジスタ内に関連する格納された 2 のオフセットを有する作業項目は、B B (2 , 6) からフェッチされた命令を得る。単一の P C および格納されたサイズが、フェッチする D V L I W のタイプを i - キャッシュに対して示すように、コンパイラは、命令をこの方式でメモリ内に既に配置している。

10

【 0 0 4 0 】

B B C (1 , 4) の完了時に、コード 5 0 6 は、トレース 1 内の B B 4 の終わりにおいて、4 を格納するようにサイズレジスタを更新し、B B (3 , 6) に分岐する作業項目に対するエントリを、3 を格納するように更新する。単一の P C および格納されたサイズを i - キャッシュに送信した後に、長さ 4 の D V L I W が i - キャッシュからフェッチされる。D V L I W は、B B 6 の第 1 のコピーと、B B 5 の単一のコピーと、B B 6 の第 2 のコピーと、B B 6 の第 3 のコピーと、から混ぜ合わされた命令を含む。ベクトルレジスタ内に関連する格納された 0 のオフセットを有する作業項目は、B B (0 , 6) からフェッチされた命令を得る。ベクトルレジスタ内に関連する格納された 1 のオフセットを有する作業項目は、B B (1 , 5) から結果を得る。ベクトルレジスタ内に関連する格納された 2 のオフセットを有する作業項目は、B B (2 , 6) からフェッチされた命令を得る。ベクトルレジスタ内に関連する格納された 3 のオフセットを有する作業項目は、B B (3 , 6) からフェッチされた命令を得る。単一の P C および格納されたサイズが、フェッチする D V L I W のタイプを i - キャッシュに対して示すように、コンパイラは、命令をこの方式でメモリ内に既に配置している。

20

【 0 0 4 1 】

B B (0 , 6) 、 B B (1 , 6) 、 B B (2 , 6) および B B (3 , 6) の各々に対して、制御フローは、B B 6 の終わりにおいて、ループの別の繰り返しのために B B 1 に戻ってもよいし、ループを終了してもよい。関連付けられた分岐命令および対応するレコード内のデータは、実行時に制御フローの方向を判断するであろう。いくつかの作業項目は、別の繰り返しを継続してもよく、他の作業項目は、ループを終了してもよい。特別コード状態は、ベクトルレジスタ内の対応するエントリ内に格納され、どのパスが選ばれるかを示してもよい。所与の作業項目が別のループの繰返しを継続し、複数の作業項目のうち少なくとも 1 つの他の作業項目が、関連付けられた基本ブロックに対するコードを処理しているという判断に応じて、スリープ状態符号化は、所与の作業項目に対するベクトルレジスタ内の関連付けられたエントリに格納され得る。

30

【 0 0 4 2 】

所与の作業項目がループを終了するという判断に応じて、終了状態符号化は、所与の作業項目に対するベクトルレジスタ内の関連付けられたエントリに格納され得る。スリープ状態符号化および終了状態符号化の各々は、ループ繰返し中に使用されるオフセットから一意であって、互いに一意である。いくつかの実施形態では、スリープ状態または終了状態にある所与の作業項目に対して、コード 5 0 6 は、所与の作業項目の実行を停止して、少なくとも次のプログラムカウンタ (P C) および作業項目識別子 (I D) を、例えば高速読み出しのためのスタックメモリなどのメモリに格納する。

40

【 0 0 4 3 】

コード 5 0 6 は、B B 6 の終わりにおいて、各作業項目の状態をチェックし得る。各作業項目がスリープ状態であるか、または各作業項目が終了状態であるという判断に応じて、プロセッサは、コード 5 0 6 を実行している間に、それぞれの格納された次の P C に

50

分岐することにより、各作業項目に対する実行を再開し得る。各作業項目が停止されており、且つ、少なくとも1つの作業項目が別の作業項目と異なる状態にあるという判断に応じて、プロセッサは、コード506を実行している間に、それぞれの格納された次のPCに分岐することにより、スリープ状態にある作業項目のみについて実行を再開し得る。少なくとも1つの作業項目が、依然として、ループ内の基本ブロックにある命令を処理している場合には、その少なくとも1つの作業項目について実行を継続する一方で、特別な状態にある他の作業項目を待機する。スリープ状態を終えている作業項目は、BB 1に分岐して戻る。また、コード502は、BB 1の開始時に、ベクトルレジスタおよびサイズレジスタを初期化する。終了状態を終えている作業項目は、BB 7に分岐する。また、コード502は、BB 7の開始時に、それに応じて、ベクトルレジスタおよびサイズレジスタを再初期化する。

10

【0044】

前述の例では、ループは単一の出口を有する。複数の出口を有する他の場合には、例えば、少なくとも次のPCおよび作業項目IDなどのような対応する状態情報が、例えばスタックなどのメモリに格納され得る。後に、状態情報は、再開のために、例えばスタックからポップするなどのように取得され得る。スリープ状態または終了状態にある両方の作業項目は、例えばスタックなどのメモリに格納された状態情報を有し得る。異なる作業項目が異なるループ繰返しでループを終了し得るので、状態情報を有する複数のエントリが、例えばスタックなどのメモリに置かれ得る。再開時において、実行される際に、コンパイラ生成コードが状態情報をポップして、同じ次のPCから再開する作業項目に対する情報を組み合わせ得る。

20

【0045】

ここで図6を参照すると、コンパイラ技術を用いて、プロセッサ内で複数の作業項目の並列実行を最適化するための方法600の一実施形態が示されている。議論を進めるために、本実施形態、および後述する方法の後続の実施形態におけるステップは、連続した順序で示されている。しかし、他の実施形態では、いくつかのステップは、示されたものとは異なる順番で起こってもよく、いくつかのステップは同時に実行されてもよく、いくつかのステップは他のステップと組み合わせられてもよく、また、いくつかのステップは存在しなくてもよい。

【0046】

ブロック602では、ソフトウェアプログラムまたはサブルーチンが検出され解析され得る。プログラムコードは、設計者により、例えばCまたは別の言語などの高水準言語で書かれてもよい。このソフトウェアプログラムは、ゲーム、ビジネス、医療および他の分野などにおいて、並列データアプリケーションのコンパイルおよび実行のために書かれてもよい。プログラムコードは、ソフトウェアアプリケーション、サブルーチン、ダイナミックリンクライブラリ、または他の任意の部分の指してもよい。パス名は、ユーザーによりコマンドプロンプトに対して入力され得る。あるいは、パス名は、ソースコードのコンパイルを開始するために、所与のディレクトリ位置、またはその他から読み込まれてもよい。プログラムコード内の命令は、検査され、翻訳され、最適化されて、コンパイル中にさらに処理されてもよい。

30

40

【0047】

いくつかの実施形態では、ソースコードが静的にコンパイルされる。かかる実施形態では、フロントエンドのコンパイル中に、ソースコードが中間表現(IR)に翻訳されてもよい。バックエンドのコンパイルステップは、IRを機械コードに翻訳してもよい。静的なバックエンドコンパイルは、さらなる変換および最適化を実行してもよい。他の実施形態では、ソースコードは、ジャストインタイム(JIT)方式でコンパイルされる。JIT方式は、システム構成を取得した後に、適切なバイナリコードを生成し得る。何れの方法でも、コンパイラは、関数呼出し、ループ、ループ内のトレースおよびプログラムコード内の基本ブロックを識別し得る。1つ以上の制御フローグラフが、プログラム解析中に構築され得る。

50

【 0 0 4 8 】

様々な実施形態では、プログラムコードは、例えば汎用プロセッサなどのプロセッサ上でコンパイルされる。プログラムコードは、ターゲットプロセッサに対してコンパイルされてもよく、ターゲットプロセッサは、例えばSIMDマイクロアーキテクチャなどの並列マイクロアーキテクチャを含む。データの1つ以上の関連付けられたレコードは、1つ以上の作業項目を生成するために、関数呼出しに割り当てられてもよい。

【 0 0 4 9 】

プログラムコード内の任意の分岐点を検出する前に、コンパイラは、解析され翻訳された命令を、それらがプログラムコード内に出現する通りにメモリ内に配置し得る。基本的に、コンパイラは、1のサイズをもつVLIWを生成している場合がある。コンパイラは、識別されたループ内で分岐点を検出すると(条件ブロック604:はい)、ブロック606において、超大命令語(VLIW)を生成し得る。コンパイラは、分岐点と対応する収束点との間の複数の基本ブロックから混ぜ合わされた命令をメモリに配置することにより、VLIWを生成し得る。データの1つ以上の関連付けられたレコードは、生成されたVLIW内の混ぜ合わされた命令を一緒にするために配置され割り当てられて、1つ以上の関連付けられた作業項目を生成し得る。

【 0 0 5 0 】

ブロック608では、生成されたコードが挿入されてもよく、生成されたコードは、実行される際に、VLIW内の命令をポイントしているオフセットを、ターゲットプロセッサ内の複数の並列実行レーンの所与のレーンにマッピングする。あるいは、マッピングは、オフセットと作業項目IDとの間であってよい。ブロック610では、実行される際に、フェッチするための次のVLIWのサイズを更新する、生成されたコードが挿入され得る。ブロック612では、実行される際に、スリープ状態または終了状態になるようにしている実行のレーンに対する状態情報を格納する、生成されたコードが、ループの終わりに挿入され得る。生成されたコードは、上述した例で説明したように、プログラムコード内の特定のポイントに挿入され得る。分岐点および収束点と関連付けられた基本ブロックは、マッピングおよびDVLIWサイズに対する更新を維持するための追加の挿入されたコードを有し得る。

【 0 0 5 1 】

ここで図7を参照すると、ハードウェア技術を使用して、プロセッサ内で複数の作業項目の並列実行を最適化するための方法700の一実施形態が示されている。議論を進めるために、本実施形態、および後述する方法の後続の実施形態におけるステップは、連続した順序で示されている。しかし、他の実施形態では、いくつかのステップは、示されたものとは異なる順番で起こってもよく、いくつかのステップは同時に実行されてもよく、いくつかのステップは他のステップと組み合わせられてもよく、また、いくつかのステップは、存在しなくてもよい。

【 0 0 5 2 】

ブロック702では、データの関連付けられたレコードが、複数の作業項目を生成するために、コンパイル済みコードに割り当てられる。ブロック704では、作業項目が、単一命令複数データ(SIMD)マイクロアーキテクチャを有するターゲットプロセッサに対してスケジューリングされる。ブロック706では、更新されたVLIWサイズおよび単一のプログラムカウンタ(PC)を使用して、更新されたVLIWサイズに等しい長さを有するVLIWが、例えばi-キャッシュなどのメモリからフェッチされる。VLIW内の命令は、ループ内の分岐点と収束点との間の別個の基本ブロックからのものであってよい。

【 0 0 5 3 】

ブロック708では、フェッチされたVLIW内の命令と、作業項目を実行しているプロセッサ内の並列実行レーンとの間のマッピング情報のために、ベクトルレジスタが読み取られる。マッピング情報は、所与の作業項目および対応する実行レーンに対して、VLIW内のどの命令を処理するかを判断し得る。ブロック710では、VLIW内の命令が

10

20

30

40

50

、並列実行レーンを使用して同時に実行される。所与の作業項目についてループの終わりに達し（条件ブロック712：はい）、アクティブであると検出された作業項目がない（条件ブロック714：いいえ）場合には、ブロック716において、それぞれの状態情報が、並列実行レーンに割り当てられた各作業項目に対して読み取られる。状態情報は、少なくとも次のPCおよび作業項目IDを含んでもよい。状態情報は、並列実行レーン内で実行を継続するために使用されてもよい。所与の作業項目についてループの終わりに達し（条件ブロック712：はい）、いずれかの作業項目がアクティブであると検出された（条件ブロック714：はい）場合には、ブロック718において、所与の作業項目に対する状態情報が後で使用されるために格納される。所与の作業項目は、実行を停止させてもよく、スリープ状態または終了状態に置かれ得る。

10

【0054】

ここで図8を参照すると、オブジェクトコード配置800を示す別の実施形態の一般化されたブロック図が示されている。コード配置800は、図2および図3にそれぞれ示す制御フローグラフ200および付随する実行順序320に対して、コンパイラによって生成され得るオブジェクトコードレイアウトの一実施形態を示している。コード502～506は、既に説明したものと同一機能を実行し得る。

【0055】

基本ブロックのサイズはコンパイル時に分かっている。コンパイラは、初期化コードを最小限にしてDVL I Wサイズを削減するようにプログラムコードの命令を配置してスケジューリングし得る。制御フローグラフ200およびオブジェクトコード配置500を使用する前述の例では、BB 6の4つのコピーが使用される。コンパイラは、ギャップを取り入れることにより、トレース数および付随のオフセットを削減し得る。ギャップは、nop操作を利用し得る。

20

【0056】

コード配置500に関して前述したステップがここで使用され得る。BB 2の完了時に、DVL I Wサイズは2のままである。しかし、フェッチされたDVL I W内の2つの命令のうち1つの命令は、配置800内のギャップを提供するnop操作である。オフセット0はnop操作に対応しており、オフセット1はBB 3内の命令に対応している。BB 3の完了時に、DVL I Wは、2から1に減らされる。この時、オフセット0はBB 4内の命令に対応しており、オフセット1はnop操作に対応している。

30

【0057】

BB C (0 , 4) の完了時に、コード506は、トレース0内のBB 4の終わりににおいて、サイズレジスタを1から2に更新する。また、コード506は、BB 6に分岐している作業項目に対するエントリを、1を格納するように更新する。単一のPCおよび格納されたサイズをi - キャッシュに送信した後に、長さ2のDVL I Wがi - キャッシュからフェッチされる。DVL I Wは、BB 5およびBB 6からの混ぜ合わされた命令を含む。ベクトルレジスタ内に関連する格納された0のオフセットを有する作業項目は、BB (0 , 5) からフェッチされた命令を得る。ベクトルレジスタ内に関連する格納された1のオフセットを有する作業項目は、BB (1 , 6) からフェッチされた命令を得る。

40

【0058】

BB (1 , 6) の完了時に、コード506は、DVL I Wサイズを2から1に更新する。BB (1 , 6) に対応する作業項目は、実行を停止して、状態情報を格納し、次のPCで実行を再開するまで待機する。次のPCは、BB 1またはBB 7をポイントしてもよい。前述したように、他のステップが実行されてもよい。コード配置800は、検出された所与の分岐点と対応する収束点との間の第1のトレースパスが、所与の分岐点と対応する収束点との間の第2のトレースパスよりも小さいという判断に応じて、コンパイラが、第1のトレースパスの完了と対応する収束点との間で生成されたVL I W内の第2のトレースパスに対応する命令とともにnopをグループ化し得ることを示す。

【0059】

前述した実施形態の各々に対して、コードを、ループの外部でループ内のコードと並列

50

化することによって、さらなる並列化が生じてもよい。例えば、BB 7に対するプログラムコードは、ループを完了する作業項目に対するプログラムを終了するために、BB 1と並列化されてもよい。また、オフセットが、単一の命令ではなく、VLIW内の複数の資源に依存しない命令に対応していることの検出に応じて、関連付けられた作業項目および実行レーンが、実行レーン内の複数の資源に依存しない命令を同時に実行してもよい。

【0060】

さらに、コンパイラは、レジスタ割当てを使用して、DVLIWサイズを削減してもよい。プログラムコードは、次の文： $X = (A + B) + (B + M)$ を含み得る。ここでは、2つの演算が同じ演算コードを使用する。第1のトレースは、 $T1 = A + B$ などのADD演算を含む。第2のトレースは、 $T2 = C + D$ などのADD演算を含む。0などのオフセットを有する基本ブロックX (BB X)からのT1を使用する作業項目がある。1などのオフセットを有するBB YからのT2を使用する他の作業項目がある。第1のオペランド対「C」および「A」、第2のオペランド対「B」および「D」、ならびに、結果対「T1」および「T2」の各々が、BB XおよびBB Y内の同じレジスタに割り当てられる場合には、式 $r3 = r1 + r2$ が、1のサイズを有するDVLIWとして使用され得る。復号時間を節約するか、またはスロットを解放するために、対応するオフセットが0に設定され得る。

【0061】

前述の実施形態は、ソフトウェアを含み得ることに留意されたい。かかる実施形態では、方法および/または機構を実装するプログラム命令は、コンピュータ可読媒体上で伝達され、または格納され得る。プログラム命令を格納するように構成されている多数のタイプの媒体が利用可能であり、ハードディスク、フロッピー(登録商標)ディスク、CD-ROM、DVD、フラッシュメモリ、プログラマブルROM(PROM)、ランダムアクセスメモリ(RAM)、および、様々な他の形式の揮発性または不揮発性記憶装置を含む。一般的に言えば、コンピュータアクセス可能記憶媒体は、命令および/またはデータをコンピュータに提供するために、使用中にコンピュータによってアクセス可能な任意の記憶媒体を含み得る。例えば、コンピュータアクセス可能記憶媒体は、磁気または光媒体などの記憶媒体、例えば、ディスク(固定または取り外し可能)、テープ、CD-ROMもしくはDVD-ROM、CD-R、CD-RW、DVD-R、DVD-RW、または、Blu-Ray(登録商標)を含み得る。記憶媒体は、RAM(例えば、シンクロナスダイナミックRAM(SDRAM)、ダブルデータレート(DDR、DDR2、DDR3など)SDRAM、低電力DDR(LPDDR2など)SDRAM、ラムバスDRAM(RDRAM)、スタティックRAM(SRAM)など)、ROM、フラッシュメモリ、ユニバーサルシリアルバス(USB)インタフェースなどの周辺インタフェースを経由してアクセス可能な不揮発性メモリ(例えば、フラッシュメモリ)などの、揮発性または不揮発性メモリ媒体をさらに含み得る。記憶媒体は、微小電気機械システム(MEMS)、ならびに、ネットワークおよび/または無線リンクなどの通信媒体を経由してアクセス可能な記憶媒体を含み得る。

【0062】

さらに、プログラム命令は、例えばCなどの高水準プログラミング言語、例えばVerilog、VHDLなどの設計言語(HDL)、または、例えばGDS IIストリーム形式(GDSII)などのデータベース形式、におけるハードウェア機能の動作レベル記述またはレジスタ転送レベル(RTL)記述を含み得る。いくつかの場合には、記述は、合成ライブラリからのゲートのリストを含むネットリストを生成するために記述を合成し得る合成ツールによって読み取られてもよい。ネットリストは、システムを含むハードウェアの機能を表すゲートのセットを含む。ネットリストは、次いで、マスクに適用される幾何学形状を記述するデータセットを生成するために配置され、ルーティングされ得る。マスクは、次いで、システムに対応する半導体回路または複数の回路を製造するために様々な半導体作製ステップで使用され得る。あるいは、コンピュータアクセス可能記憶媒体

10

20

30

40

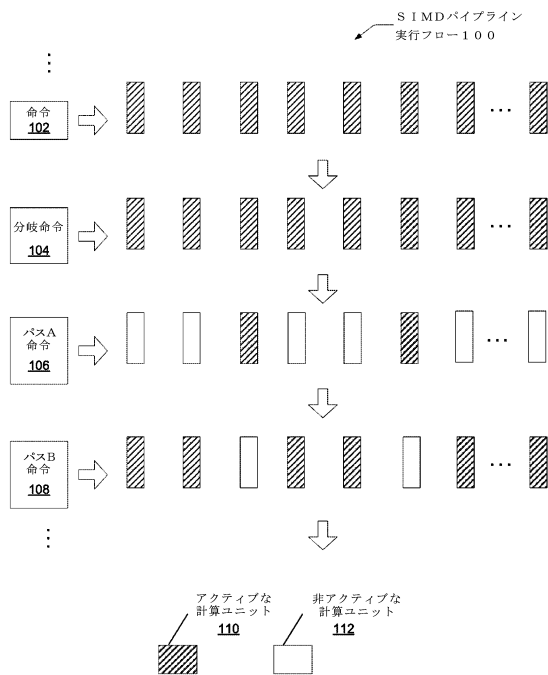
50

上の命令は、必要に応じて、ネットリスト（合成ライブラリの有無にかかわらず）またはデータセットであってよい。また、命令は、Cadence（登録商標）、EVE（登録商標）およびMentor Graphics（登録商標）などのベンダーからのハードウェアベースタイプのエミュレータによるエミュレーションのために利用され得る。

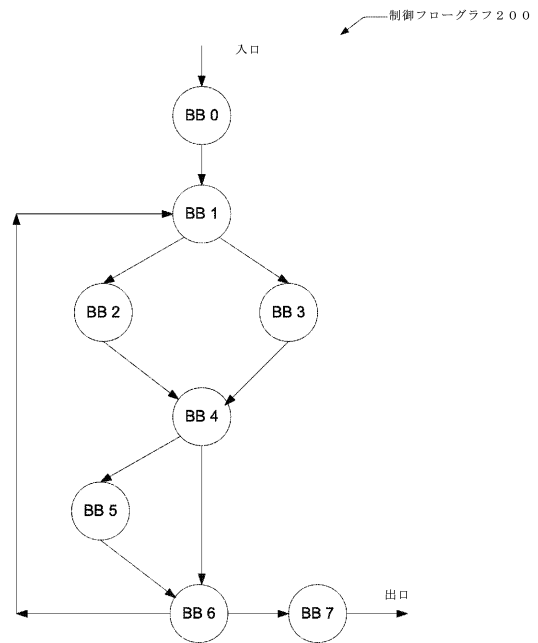
【0063】

上記実施形態はかなり詳細に説明されているが、上述した開示が完全に理解されると、多数の変形および修正が当業者において明らかになるであろう。以下の請求項は、かかる変形および修正の全てを包含すると解釈されることを意図する。

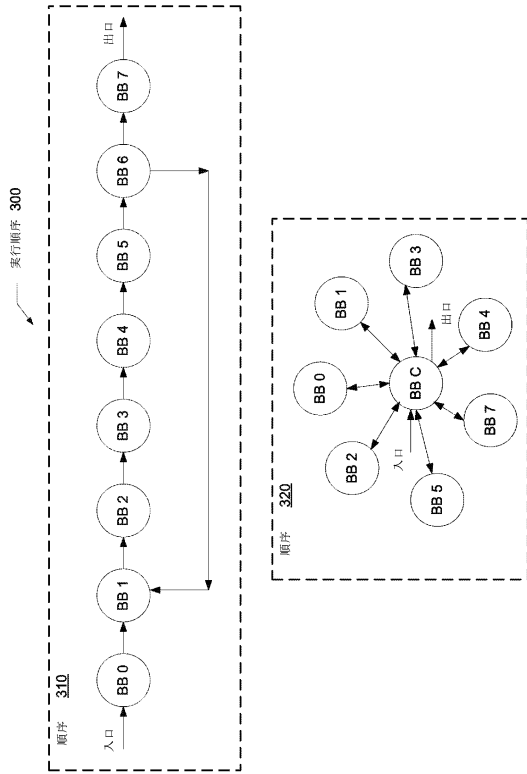
【図1】



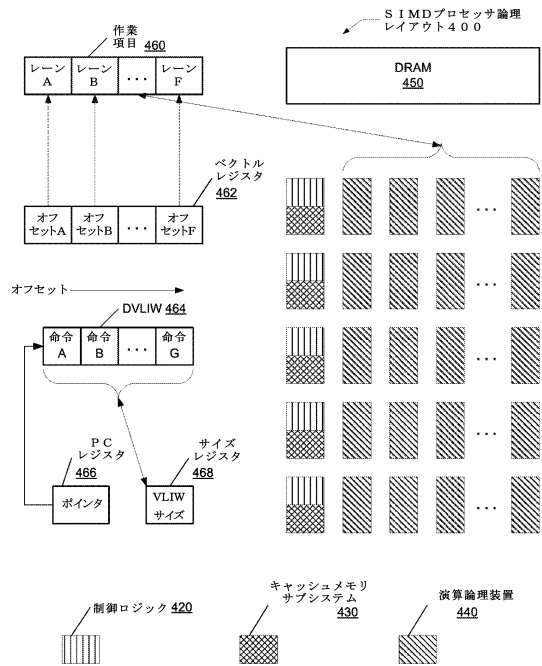
【図2】



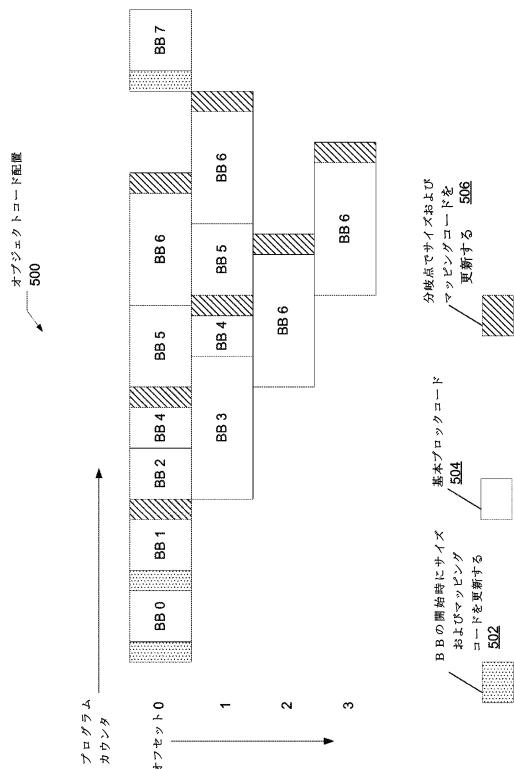
【図3】



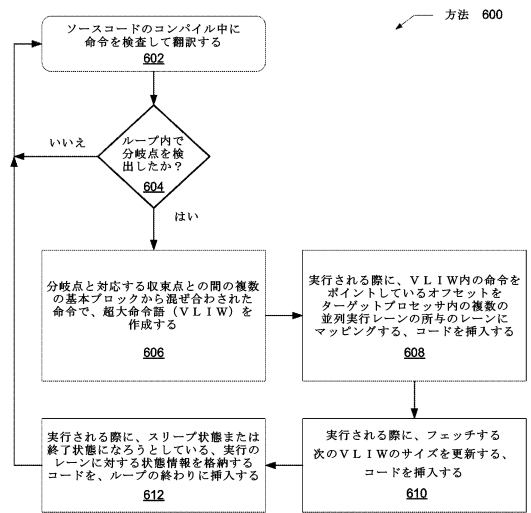
【図4】



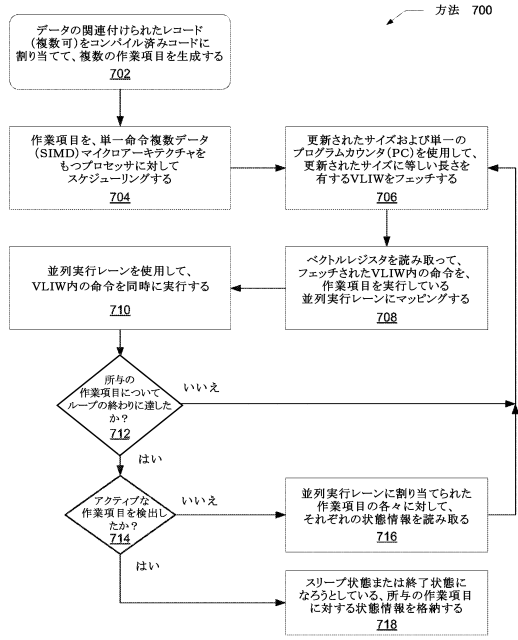
【図5】



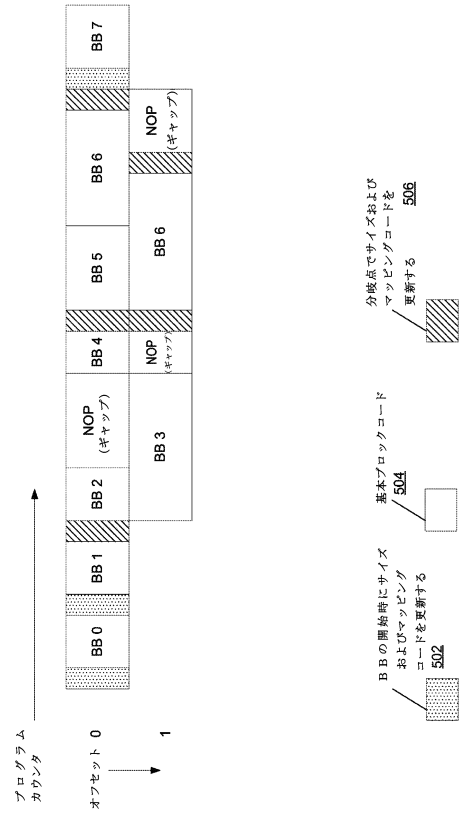
【図6】



【図 7】



【図 8】



フロントページの続き

(74)代理人 100162156

弁理士 村雨 圭介

(72)発明者 レザ ヤズダニ

アメリカ合衆国 9 4 0 2 4 カリフォルニア州、ロスアルトス、マッツ コート 9 1 1

審査官 清木 泰

(56)参考文献 米国特許第07287152(US, B2)

米国特許出願公開第2006/0242645(US, A1)

特開2008-090744(JP, A)

特表2003-520360(JP, A)

特表2000-509528(JP, A)

特開2000-259579(JP, A)

Wilson W. L. Fung, Ivan Sham, George Yuan, Tor M. Aamodt, Dynamic Warp Formation and Scheduling for Efficient GPU Control Flow, Proceedings of 40th IEEE/ACM International Symposium on Microarchitecture, IEEE, 2007年12月 1日, Pages:407-418

Yaohua Wang, et al., Instruction Shuffle: Achieving MIMD-like Performance on SIMD Architectures, IEEE Computer Architecture Letters, IEEE, 2011年12月 6日, Vol:11, No:2, Pages:37-40

Nicolas Brunie, Sylvain Collange, Gregory Diamos, Simultaneous Branch and Warp Interweaving for Sustained GPU Performance, Proceedings of 39th Annual International Symposium on Computer Architecture, IEEE, 2012年 6月 9日, Pages:49-60

Jiayuan Meng, David Tarjan, Kevin Skadron, Dynamic Warp Subdivision for Integrated Branch and Memory Divergence Tolerance, Proceedings of the 37th annual International Symposium on Computer Architecture (ISCA'10), ACM, 2010年 6月19日, Pages:235-246

(58)調査した分野(Int.Cl., DB名)

G06F 9/30 - 9/42

G06F 15/80

G06F 9/45

G06F 17/16

G06F 7/57 - 7/575

G06F 15/16 - 15/177