

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7245833号

(P7245833)

(45)発行日 令和5年3月24日(2023.3.24)

(24)登録日 令和5年3月15日(2023.3.15)

(51)国際特許分類

F I

G 0 6 F 9/38 (2018.01)

G 0 6 F 9/38 3 7 0 C

G 0 6 F 15/80 (2006.01)

G 0 6 F 15/80

請求項の数 17 (全16頁)

| | | | |
|-------------------|-------------------------------|----------|---------------------|
| (21)出願番号 | 特願2020-529105(P2020-529105) | (73)特許権者 | 520039124 |
| (86)(22)出願日 | 平成30年8月2日(2018.8.2) | | ネクスト シリコン リミテッド |
| (65)公表番号 | 特表2020-530175(P2020-530175 A) | | Next Silicon Ltd |
| (43)公表日 | 令和2年10月15日(2020.10.15) | | イスラエル国 ギヴァタイム デレク イ |
| (86)国際出願番号 | PCT/US2018/045008 | | ツハク ラビン ストリート 33 |
| (87)国際公開番号 | WO2019/028253 | | 33 Derech Yitshak R |
| (87)国際公開日 | 平成31年2月7日(2019.2.7) | | abin Street, Givata |
| 審査請求日 | 令和3年7月29日(2021.7.29) | (74)代理人 | yim, Israel |
| (31)優先権主張番号 | 62/540,849 | | 110002952 |
| (32)優先日 | 平成29年8月3日(2017.8.3) | (72)発明者 | 弁理士法人鷲田国際特許事務所 |
| (33)優先権主張国・地域又は機関 | | | ラズ エラド |
| | 米国(US) | | イスラエル国 ラマト ガン ピンヤミン |
| (31)優先権主張番号 | 62/558,090 | | ストリート 11 |
| (32)優先日 | 平成29年9月13日(2017.9.13) | 審査官 | 三坂 敏夫 |
| 最終頁に続く | | 最終頁に続く | |

(54)【発明の名称】 構成可能なハードウェアの実行時の最適化

(57)【特許請求の範囲】

【請求項1】

少なくとも1つの計算装置によって、プログラムコードの複数の部分を含むソフトウェアプログラムを実行するためのシステムであって、

少なくとも1つの計算装置に接続された計算グリッドであって、データ入力メッシュネットワーク、複数の論理要素、および複数のマルチプレクサ、を有し、前記複数の論理要素および前記複数のマルチプレクサは、複数の計算グループを形成し、前記論理要素の少なくとも1つおよび複数のマルチプレクサの少なくとも1つは、複数の計算グループのそれぞれを形成する、前記計算グリッドと、

プログラムコードの前記複数の部分のうち、プログラムコードの2つ以上の部分を含むプログラムコードパターンであって、各パターンは、前記ソフトウェアプログラムを実行するときにプログラムコードの前記2つ以上の部分のそれぞれの出現における関連性を示している少なくとも1つの前記プログラムコードパターンにアクセスし、

前記少なくとも1つのプログラムコードパターンのそれぞれに対し、

プログラムコードの前記2つ以上の部分の少なくとも一部を操作し、

プログラムコードの前記2つ以上の部分のうち、操作された少なくとも一部を実行するように、前記複数の計算グループの少なくともいくつかを操作し、

少なくとも1つの計算装置の代わりに、前記計算グリッドによって少なくとも1つのプログラムコードパターンを計算する、

ことを複数回実行することにより、前記プログラムコードの複数の部分の少なくともいく

10

20

つかを実行するように前記計算グリッドを操作するために適合された少なくとも 1 つの処理回路と、

を備える、システム。

【請求項 2】

前記少なくとも 1 つの計算装置によって前記ソフトウェアプログラムを実行することは、前記計算グリッドがプログラムコードの前記複数の部分から選択されたプログラムコードの部分の第 2 のセットを実行する間に前記プログラムコードの複数の部分から選択されたプログラムコードの部分の第 1 のセットをそれによって実行することを含む、

請求項 1 に記載のシステム。

【請求項 3】

前記少なくとも 1 つの処理回路は、
少なくとも 1 つの他の処理回路から前記少なくとも 1 つのプログラムコードパターンを受信すること、および、
プログラムコードの前記複数の部分のうちの 1 つをそれぞれが実行するための複数の呼出しを受信し、

前記複数の呼び出しにおいて、前記少なくとも 1 つのプログラムコードパターンを識別する、

ことにより、前記少なくとも 1 つのプログラムコードパターンを計算すること、
の少なくとも 1 つのためにさらに適合されている、

請求項 1 に記載のシステム。

【請求項 4】

前記少なくとも 1 つのプログラムコードパターンは、

プログラムコードの前記 2 つ以上の部分のそれぞれのうちのコードの第 1 の部分およびコードの 2 つ以上の部分のそれぞれのうちのコードの第 2 の部分であって、前記複数の呼出しにおいて、プログラムコードの前記第 1 の部分およびプログラムコードの前記第 2 の部分が次々に出現する第 1 の頻度が所定のしきい値頻度を越えているような、前記第 1 の部分および前記第 2 の部分を有している、第 1 のプログラムコードパターンと、

プログラムコードの前記 2 つ以上の部分のそれぞれのうちのコードの第 3 の部分およびコードの前記 2 つ以上の部分のそれぞれのうちのコードの第 4 の部分であって、前記複数の呼出しにおいて、コードの前記第 3 の部分がコードの前記第 4 の部分呼び出す第 2 の頻度が所定のしきい値頻度を越えているような、前記第 3 の部分および前記第 4 の部分を有している、第 2 のプログラムコードパターンと、

プログラムコードの前記 2 つ以上の部分のそれぞれのうちのコードの第 5 の部分をそれぞれ有している、第 3 のコードパターンおよび第 4 のコードパターンと、

の少なくとも 1 つを備える、

請求項 3 に記載のシステム。

【請求項 5】

前記複数の計算グループの前記少なくともいくつかを操作することは、前記第 1 のプログラムコードパターンのために、プログラムコードの前記第 1 の部分を実行するように、前記複数の計算グループのうちの少なくとも 1 つの第 1 の計算グループを操作することを含んでおり、

前記少なくとも 1 つの第 1 の計算グループは、プログラムコードの前記第 2 の部分を実行するように構成された前記少なくとも 1 つの第 2 の計算グループにトポロジ的に近接している前記少なくとも 1 つの計算グリッド内に割り当てられている、

請求項 4 に記載のシステム。

【請求項 6】

前記複数の計算グループの前記少なくともいくつかを操作することは、前記第 2 のプログラムコードパターンのために、プログラムコードの前記第 3 の部分を実行するように、前記複数の計算グループのうちの少なくとも 1 つの第 3 の計算グループを操作することを含んでおり、

10

20

30

40

50

前記少なくとも 1 つの第 3 の計算グループは、プログラムコードの前記第 4 の部分を実行するように構成された前記少なくとも 1 つの第 4 の計算グループにトポロジ的に近接している前記少なくとも 1 つの計算グリッド内に割り当てられている、

請求項 4 に記載のシステム。

【請求項 7】

前記複数の計算グループの前記少なくともいくつかを操作することは、プログラムコードの前記第 5 の部分を実行するように、前記複数の計算グループのうちの少なくとも 2 つの計算グループを操作することを含む、

請求項 4 のシステム。

【請求項 8】

前記複数の計算グループの前記少なくともいくつかを操作することは、前記少なくとも 1 つのプログラムコードパターンの少なくとも 1 つに対して、

プログラムコードの前記 2 つ以上の部分のそれぞれのうちの 1 つの少なくとも一部を実行するようにそれぞれ構成された前記複数の計算グループの少なくとも 2 つを識別すること、

少なくとも 2 つの他の計算グループが互いにトポロジ的に近接した状態で前記計算グリッド内に割り当てられるように、前記複数の計算グループの少なくとも 2 つの他の計算グループを選択すること、および、

前記少なくとも 2 つの計算グループが実行されたプログラムコードの前記 2 つ以上の部分のそれぞれのうちの 1 つの少なくとも一部を実行するように、前記少なくとも 2 つの他の計算グループのそれぞれを操作すること、

を含む、請求項 1 に記載のシステム。

【請求項 9】

前記複数の計算グループの前記少なくともいくつかを操作することは、前記少なくとも 1 つのプログラムコードパターンの少なくとも 1 つのために、

プログラムコードの前記 2 つ以上の部分のそれぞれのプログラムコードの第 1 の部分を実行するように操作された前記複数の計算グループのうちの少なくとも 1 つの第 1 の計算グループを識別すること、および、

プログラムコードの前記 2 つ以上の部分のそれぞれのプログラムコードの第 2 の部分を実行するように、前記複数の計算グループのうちの少なくとも 1 つの第 2 の計算グループであって、前記少なくとも 1 つの計算グループにトポロジ的に近接している前記計算グループ内に割り当てられている前記第 2 の計算グループを操作すること、

を含む、請求項 1 に記載のシステム。

【請求項 10】

前記複数の計算グループの前記少なくともいくつかを操作することは、

前記複数の計算グループのいずれにもアクセスされておらず、かつ前記計算装置にもアクセスされていない前記複数の計算グループのうちの少なくとも 1 つの計算グループを解放すること、

をさらに含む、請求項 1 に記載のシステム。

【請求項 11】

前記少なくとも 1 つの処理回路は、前記ソフトウェアプログラムが前記少なくとも 1 つの計算装置および前記計算グリッドによって実行されている間、前記計算グリッドを操作するために適合される、

請求項 1 に記載のシステム。

【請求項 12】

前記少なくとも 1 つの計算装置は、

マルチコア中央処理装置 (CPU)、フィールドプログラマブルゲートアレイ (FPGA)、グラフィック処理装置 (GPU)、粗粒度再構成可能アーキテクチャ (CGRA)、ニューラルネットワークアクセラレータ、インテリジェンス処理ユニット (IPU)、特定用途向け集積回路 (ASIC)、量子コンピュータ、

10

20

30

40

50

のいずれか 1 つである、
請求項 1 に記載のシステム。

【請求項 1 3】

少なくとも 1 つの計算装置によって、プログラムコードの複数の部分を含むソフトウェアプログラムを実行するための方法であって、

プログラムコードの前記複数の部分のうちの少なくともいくつかを実行するように、データ入力メッシュネットワーク、複数の論理要素、および複数のマルチプレクサを有する計算グリッドの少なくとも一部分を操作するステップを備え、

前記複数の論理要素および前記複数のマルチプレクサは、複数の計算グループを形成しており、

前記論理要素の少なくとも 1 つおよび複数のマルチプレクサの少なくとも 1 つは、複数の計算グループのそれぞれを形成しており、

前記計算グリッドの前記少なくとも一部分を操作するステップは、

プログラムコードの前記複数の部分のうち、2 つ以上を有するプログラムコードパターンであって、各パターンは、前記ソフトウェアプログラムを実行するときにプログラムコードの前記 2 つ以上の部分のそれぞれの出現における関連性を示している少なくとも 1 つの前記プログラムコードパターンにアクセスするステップと、

前記少なくとも 1 つのプログラムコードパターンのそれぞれに対し、
プログラムコードの前記 2 つ以上の部分の少なくとも一部を操作し、

プログラムコードの前記 2 つ以上の部分のうち、操作された少なくとも一部を実行するように、複数の計算グループの少なくともいくつかを操作するステップと、

前記少なくとも 1 つの計算装置の代わりに、前記計算グリッドによって少なくとも 1 つのプログラムコードパターンを計算するステップと、

を複数回実行することを含む、

方法。

【請求項 1 4】

プログラムコードの複数の部分を有するソフトウェアプログラムを少なくとも 1 つの計算装置によって実行するプロセスを処理ユニットに実行させる命令を格納している非一過性のコンピュータ可読媒体であって、前記プロセスは、

プログラムコードの前記複数の部分のうちの少なくともいくつかを実行するように、データ入力メッシュネットワーク、複数の論理要素、および複数のマルチプレクサを有する計算グリッドの少なくとも一部分を操作するステップを備え、

前記複数の論理要素および前記複数のマルチプレクサは、複数の計算グループを形成しており、

前記論理要素の少なくとも 1 つおよび複数のマルチプレクサの少なくとも 1 つは、複数の計算グループのそれぞれを形成しており、

前記計算グリッドの前記少なくとも一部分を操作するステップは、

プログラムコードの前記複数の部分のうち、2 つ以上を有するプログラムコードパターンであって、各パターンは、前記ソフトウェアプログラムを実行するときにプログラムコードの前記 2 つ以上の部分のそれぞれの出現における関連性を示している少なくとも 1 つの前記プログラムコードパターンにアクセスするステップと、

前記少なくとも 1 つのプログラムコードパターンのそれぞれに対し、
プログラムコードの前記 2 つ以上の部分の少なくとも一部を操作し、

プログラムコードの前記 2 つ以上の部分のうち、操作された少なくとも一部を実行するように、複数の計算グループの少なくともいくつかを操作するステップと、

前記少なくとも 1 つの計算装置の代わりに、前記計算グリッドによって少なくとも 1 つのプログラムコードパターンを計算するステップと、

を複数回実行することを含む、

非一過性のコンピュータ可読媒体。

【請求項 1 5】

10

20

30

40

50

計算装置の動作を加速させる装置であって、
前記装置は、
少なくとも1つの関数を実行するための複数の呼出しを含むソフトウェアプログラムの少なくとも一部と、前記計算装置により実行される前記ソフトウェアプログラムの少なくとも別の一部を実行する計算グリッドを備え、
前記計算グリッドは、
データ入力メッシュネットワークと、
複数の論理素子と、
複数のマルチプレクサであって、前記複数の論理素子と前記複数のマルチプレクサが複数の計算グループを形成し、前記論理素子の少なくとも1つと前記複数のマルチプレクサの少なくとも1つが、前記計算グループのそれぞれを形成する、前記複数のマルチプレクサと、
を備え、
前記計算装置によって実行される少なくとも1つの関数を実行するための複数の呼出しの間で少なくとも1つのパターンが識別された場合、前記計算装置は、
前記少なくとも1つの関数进行操作し、
前記少なくとも1つのパターンに基づいて、前記複数の計算グループの少なくとも1つを操作して、操作された前記少なくとも1つの関数を実行する、
装置。

10

【請求項16】

20

前記少なくとも1つのパターンが少なくとも1つの関数を実行するための前記複数の呼出しにおける少なくとも1つの繰り返し関数の再出現を示している場合、前記計算装置は、
前記複数の計算グループの少なくとも2つが前記少なくとも1つの繰り返し関数を実行するように前記複数の計算グループの少なくとも2つのそれぞれを操作する、

請求項15に記載の装置。

【請求項17】

前記少なくとも1つのパターンが、第1の関数および第2の関数を用いた演算の実行を示している場合、前記計算装置は、前記第1の関数を実行する前記複数の計算グループの少なくとも1つの第1計算グループ进行操作し、
前記少なくとも1つの第1計算グループは、前記計算グリッドにおいて、前記第2の関数を実行する前記複数の計算グループの少なくとも1つの第2計算グループにトポロジ的に近接するように割り当てられる、

30

請求項15または16に記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、全般的には、ハードウェアの実行時の最適化に関し、より詳細には、構成可能なハードウェアを実行時に最適化する手法に関する。

【背景技術】

【0002】

40

技術が進歩するにつれて、より高い処理能力を有する、より強力な処理システムのニーズが急速に高まる。最近では、プロセッサは、高い計算スループットを提供し、かつ高い電力効率であることが期待されている。しかしながら既存の処理システムは、メモリへの明示的な格納を通じて伝えられる命令の逐次的なストリームを実行し、したがってモデルの非省電力性という課題がある。

【0003】

現代の処理アーキテクチャにおいては、たとえプログラムの大部分がコードの小さな静的部分を繰り返す場合でも、各動的命令をフェッチしてデコードしなければならない。さらには、状態を明示的に記憶することが、命令間でデータを伝えるための唯一の経路であるため、機能ユニットとレジスタファイルの間で中間結果が繰り返し伝送される。現代の

50

計算アーキテクチャの制限としては、高い電力消費量、熱放散、ネットワークおよびＩ／Ｏのボトルネック、およびメモリパーティションが挙げられる。

【 0 0 0 4 】

例えば、フィールドプログラマブルゲートアレイ（ＦＰＧＡ）は、ソフトウェアによって構成されるハードウェア回路上で動作する。ＦＰＧＡでは、低レイテンシでの極めて高いスループット率が可能である。ＦＰＧＡの構成可能性は、マルチコアアーキテクチャにおけるコプロセッサとして使用することができ、またはシステムのＣＰＵの処理負荷をオフロードするためにクリティカルデータパス上に配置することができる。ＦＰＧＡの主たる欠点の１つは、柔軟なプログラミング性がないことである。さらに、ＦＰＧＡの計算能力は比較的低い。

10

【 0 0 0 5 】

プログラム可能な処理アーキテクチャの例は、マルチコアプロセッサである。マルチコアプロセッサのアーキテクチャは、２つ以上の独立した実際の処理ユニット（「コア」）を有する１つの計算要素を含み、処理ユニットはプログラム命令を読み取って実行するユニットである。これらの命令は、普通のＣＰＵ命令（例：加算、データの移動、分岐）である。このアーキテクチャでは、１つのプロセッサが、複数の命令を個別のコア上で並列に実行することができる。マルチコアプロセッサの主たる欠点は、高い電力消費量および低いスループットである。

【 0 0 0 6 】

処理アーキテクチャの別の例は、グラフィック処理装置（ＧＰＵ）である。ＧＰＵは、複数のタスクを同時に扱うように設計されている、より小さく、より効率的な何千ものコアからなる並列アーキテクチャに基づいている。ＧＰＵは、ディープラーニング、グラフィックスのレンダリング、機械学習のアプリケーションの計算タスクをアクセラレートさせるために利用することができる。ＧＰＵの主たる欠点は、高い電力消費量および高レイテンシである。さらにＧＰＵは、メモリー貫性（memory coherency）がなく、したがって共有メモリには課題がある。

20

【 0 0 0 7 】

したがって、上述した欠陥を克服する処理アーキテクチャを提供することは有利であろう。

【 発明の概要 】

30

【 課題を解決するための手段 】

【 0 0 0 8 】

以下では、いくつかの例としての実施形態の要約を開示する。この要約は、そのような実施形態を読み手が基本的に理解できるように便宜的に提供されるものであり、本開示の範囲を完全に定義するものではない。この要約は、考えられるあらゆる実施形態の広範な概要ではなく、すべての実施形態の主要または重要な要素を識別することも、いずれかの態様またはすべての態様の範囲を示すことも意図していない。その唯一の目的は、１つまたは複数の実施形態のいくつかのコンセプトを、後から提示するさらに詳細な説明の導入部として簡略的な形で提示することである。本明細書では、本開示の１つの実施形態または複数の実施形態を指す目的で、便宜上、語「いくつかの実施形態」を使用することがある。

40

【 0 0 0 9 】

開示されている実施形態のさまざまな態様は、構成可能な処理アーキテクチャを実行時に最適化する方法を含む。本方法は、少なくとも１つの関数を実行するための複数の呼出しを受け取るステップと、受け取った複数の呼出しの間で少なくとも１つのパターンを識別するステップと、少なくとも１つの関数を計算するために、この少なくとも１つのパターンに基づいて、構成可能な処理アーキテクチャの少なくとも一部を操作するステップと、を含む。

【 0 0 1 0 】

開示されている実施形態のさまざまな態様は、構成可能な処理アーキテクチャをさらに

50

含む。本システムは、計算グリッドと、処理ユニットと、処理ユニットに接続されているメモリであって、メモリが命令を含み、命令が処理ユニットによって実行されたとき、命令が、少なくとも1つの関数を実行するための複数の呼出しを受け取り、受け取った複数の呼出しの間で少なくとも1つのパターンを識別し、少なくとも1つの関数を計算するために、この少なくとも1つのパターンに基づいて、構成可能な処理アーキテクチャ内の計算グリッドを操作する、ように、処理ユニットを構成する、メモリと、を備えている。

【0011】

本明細書に開示されている主題は、本明細書の最後における請求項に具体的に示され、明確に特許請求されている。本発明の上記およびそれ以外の目的、特徴、および利点は、添付の図面を参照しながら行われる以下の詳細な説明から明らかであろう。

10

【図面の簡単な説明】

【0012】

【図1A】一実施形態に係る構成可能な処理アーキテクチャの概略図である。

【図1B】一実施形態に係る構成可能な処理アーキテクチャの計算グリッドの概略図である。

【図2A】一実施形態に係る、それぞれのパターンの識別および関数の構成を示している概略的なシミュレーションである。

【図2B】一実施形態に係る、それぞれのパターンの識別および関数の構成を示している概略的なシミュレーションである。

【図2C】一実施形態に係る、それぞれのパターンの識別および関数の構成を示している概略的なシミュレーションである。

20

【図2D】一実施形態に係る、それぞれのパターンの識別および関数の構成を示している概略的なシミュレーションである。

【図3】一実施形態に係る構成可能なハードウェアの実行時の最適化の流れ図である。

【発明を実施するための形態】

【0013】

重要な点として、本明細書に開示されている実施形態は、本明細書における革新的な教示内容の多数の有利な使用の例にすぎない。本出願の明細書中に行われている記述は、一般的に、特許請求されるさまざまな実施形態のいずれも必ずしも制限しない。さらに、いくつかの記述は、いくつかの独創的な特徴にあてはまるが、それ以外の特徴にはあてはまらないことがある。一般には、特に明記しない限り、一般性を失うことなく、1つの要素は複数の要素であってもよく、逆も同様である。図面においては、類似する数字は、いくつかの図を通じて類似する部分を指している。

30

【0014】

図1Aは、一実施形態に係る構成可能な処理アーキテクチャ100の例としての概略図を示している。処理アーキテクチャ100は、計算グリッド110と、メモリ（プログラムメモリ102など）に結合されている処理ユニット（回路）101とを含む。

【0015】

処理ユニット101は、計算グリッド110でのプログラムコードの一部の実行を最適化するための処理を実行するように構成されている。プログラムコードの一部は、関数、基本ブロック、またはその両方（まとめて関数と称する）を含むことができる。基本ブロックとは、ブロックの中からのジャンプ、またはブロックの中へのジャンプが存在しないような、連続する命令を有するプログラムコードの断片である。

40

【0016】

一実施形態においては、構成可能な処理アーキテクチャ100は、計算装置の動作をアクセラレートさせるように構成されている。このような装置としては、例えば、マルチコア中央処理装置（CPU）、フィールドプログラマブルゲートアレイ（FPGA）、グラフィック処理装置（GPU）、特定用途向け集積回路（ASIC）、量子コンピュータ、光コンピューティング、ニューラルネットワークアクセラレータ、またはこれらの組合せ、が挙げられる。

50

【 0 0 1 7 】

開示されている実施形態によれば、アクセラレーションは、例えば、計算装置（図示していない）上ではなく、計算グリッド 1 1 0 上で関数を実行することによって達成される。計算グリッド 1 1 0 は、本明細書の中で後から図 1 B を参照しながらさらに説明するように、論理素子のアレイおよびマルチプレクサ（M U X ）を含む。

【 0 0 1 8 】

計算グリッド 1 1 0 による関数の実行は、射影グラフ（projection graph）を使用して関数を計算グループに射影することによって行われる。計算グループとは、グリッド 1 1 0 内で M U X を介して接続されている論理素子のセットである。グリッド 1 1 0 内で論理素子を正しく割り当てて選択することによって、関数の最適化された実行を達成することができる。

10

【 0 0 1 9 】

具体的には計算グリッド 1 1 0 は、グリッド 1 1 0 の動作を最適化するように、一部分が処理ユニット 1 0 1 によって構成される。この目的のため、処理ユニット 1 0 1 は、複数の最適化処理を実行する。

【 0 0 2 0 】

一実施形態においては、処理ユニット 1 0 1 は、計算装置内の計算要素によって動作する関数のための複数の呼出しを受け取るように構成されている。この呼出しは、複数の異なる関数、同じ関数、または両方の呼出し、メモリの異なる部分、メモリの同じ部分、または両方の呼出し、などとすることができる。関数は、算術演算、論理演算、またはその両方（ただしこれらに限定されない）を含む計算演算を実行することができる。

20

【 0 0 2 1 】

処理ユニット 1 0 1 は、受け取った複数の呼出しの間で少なくとも 1 つのパターンを識別するように構成されている。別の実施形態によれば、識別を外部処理によって行うことができ、その後その識別を処理ユニット 1 0 1 に転送することができる。パターンは、一実施形態においては、2 つ以上の関数が、ある所定のしきい値を超える関連性で再出現することである。しきい値は、特定の関数の再出現（re-occurrence）の特定の統計分析などに基づいて、経時的に動的に変更することができる。すなわち、例えば 2 つの関数が比較的頻繁に次々と実行されるとき、これら 2 つの関数の間の関連性のパターンが求められる。

30

【 0 0 2 2 】

処理ユニット 1 0 1 は、関数の実行を最適化するため、求められたパターンに基づいて計算グリッド 1 1 0 を構成する、または別の方法で操作する。関数の操作は、複数の処理（複製、拡張、縮小、インライン化、修正、またはこれらの組合せなどを含み、ただしこれらに限定されない）を通じて、達成することができる。

【 0 0 2 3 】

例としての一実施形態においては、再配置処理を実行するとき、処理ユニット 1 0 1 は、頻繁に互いを呼び出す 2 つの計算グループがグリッド 1 1 0 内でトポロジ的に互いに近くに配置されるように、計算グリッド 1 1 0 を構成する。

【 0 0 2 4 】

別の実施形態においては、処理ユニット 1 0 1 は、それまでの使用頻度が高い計算グループのセットが多くのインスタンスに複製されるように、計算グリッド 1 1 0 を構成する。さらなる実施形態においては、計算グループの呼出し元は、複製された新しい計算グループの間で負荷分散する。

40

【 0 0 2 5 】

さらに別の実施形態においては、インライン化処理を実行するときには、処理ユニット 1 0 1 は、頻繁に互いを呼び出す 2 つの論理上の計算グループが、これら 2 つの論理上の計算グループを結合する 1 つの論理上の計算グループとして構成および再構成されるように、計算グリッド 1 1 0 を構成する。

【 0 0 2 6 】

50

さらに別の実施形態においては、拡張処理を実行するときには、処理ユニット 101 は、分岐およびループを含む論理上の計算グループが、インライン化処理および再配置処理を使用してループが展開された状態で再構成されるように、計算グリッド 110 を構成する。

【0027】

さらに別の実施形態においては、縮小処理を実行するときには、処理ユニット 101 は、最適化された実行をもはや提供しない複製または拡張された計算グループが、再構成により解放されるように、計算グリッド 110 を構成する。縮小された計算グループの呼出し元は、別の計算グループに参照される。参照がなくなった時点で、計算グループを再構成により解放することができる。すなわち、その計算グループに関連付けられるリソースを解放することができる。

【0028】

なお、処理ユニット 101 は、電力消費量、メモリアクセス、レイテンシ、ゲート使用回数 (gate usage count)、およびスループット、のうちの少なくとも 1 つに関連する必要性に従って、より良好な性能が達成されるように計算グリッド 110 を構成することに留意されたい。

【0029】

処理ユニット 101 は、1 つまたは複数のハードウェア論理要素および論理回路として実施することができる。使用することのできるハードウェア論理要素のタイプの実例としては、例えば、以下に限定されないが、汎用マイクロプロセッサ、マイクロコントローラ、マルチコア CPU、FPGA、GPU、ASIC、量子プロセッサ、光コンピューティングプロセッサ、ニューラルネットワークアクセラレータ、粗粒度再構成可能アーキテクチャ (CGRA: coarse-grained configurable architecture)、インテリジェンス処理ユニット (intelligence processing unit) など、または、情報の計算もしくは別の操作を実行することのできる任意の他のハードウェア論理要素、が挙げられる。いくつかの実施形態においては、アクセラレートさせる対象の計算装置は、処理ユニット 101 の役割を果たす。

【0030】

いくつかの実施形態においては、メモリバウンドである関数を含む、メモリの厳しいコード断片の分析が実行される。このような関数は、計算装置のメモリの近くに再配置される計算グループにマッピングされる。このような分析は、I/O 関連の動作に適用することもできる。例えば、ネットワークポート、記憶装置、PCI-e、ビデオセンサーなどに関する動作である。

【0031】

次に図 1B を参照し、例としての一実装形態においては、計算グリッド 110 は、複数の論理素子 120 (まとめて 1 つ以上の LE 120 と称する、または個別に LE 120 と称する) を含む。一実施形態においては、LE 120 は、論理演算子 (AND、OR、NOT、XOR など) またはこれらの組合せとすることができる。さらに別の実施形態においては、LE 120 を、ルックアップテーブル (LUT) 演算を実行するように構成することができる。さらに別の実施形態においては、LE 120 を、高レベルの算術関数 (固定小数点数または浮動小数点数の加算、減算、乗算、除算、指数など) を実行するように構成することができる。さらに別の実施形態においては、LE 120 は、シフト演算 (左シフト、右シフトなど) を実行することができる。

【0032】

なお、計算グリッド 110 内の各 LE 120 は、上述した演算関数のいずれかが 1 つ以上に対応するように構成することができることに留意されたい。いくつかの構成においては、第 1 のグループの 1 つ以上の LE 120 が第 1 の演算を実行することができ (例: 論理演算子)、第 2 のグループの 1 つ以上の LE 120 が別の演算 (例: シフト演算) を実行することができ、以下同様である。特定の構成においては、すべての LE 120 が、同じ演算関数を実行することができる。

10

20

30

40

50

【0033】

一実施形態においては、計算グリッド110は、マルチプレクサ、デマルチプレクサ、スイッチなど複数のデータルーティング中継部（data routing junction）130（まとめて1つ以上のMUX 130と称する、または個別にMUX 130と称する）をさらに含む。MUX 130は、データを1つ以上のLE 120にルーティングする、および1つ以上のLE 120からのデータをルーティングする、ように構成されている。

【0034】

計算グリッド110は、複数の計算グループから構成されている。各計算グループは、1つのMUX 130を介して接続されている $N \times M$ 個のLE 120を含む。「N」および「M」の値は、異なるかまたは等しい整数である。一実施形態においては、Mは、N - 1に等しい（ $M = N - 1$ ）。MUX 130は、計算グリッド110内の2つの隣接する行からの1つ以上のLE 120を接続している。

10

【0035】

一実施形態においては、各MUX 130は、リテラル値を維持するための複数のレジスタ（図示していない）を含む。MUX 130は、実行される演算の条件、接続テーブル、およびコード有するようにさらに構成される。1つ以上のMUX 130は、バス（図示していない）を介して処理ユニット101によって構成される。

【0036】

別の実施形態においては、計算グリッド110は、データ入力メッシュネットワーク（data input mesh network）の役割を果たす、フィーダ（feeder）（Fr）140のメッシュネットワークを含む。フィーダ140は、1つ以上のLE 120によって処理するために入力（例えば関数のパラメータ、ただしこれに限定されない）を送り込むように構成されている。一実施形態においては、フィーダ140は、内部チップ相互接続（ICI：internal chip interconnection）として実施される。

20

【0037】

なお、計算グリッド110の操作は、アクセラレートさせる対象の計算装置の動作と同時にリアルタイムで実行されることを理解されたい。一実施形態によれば、操作は、1つまたは複数の関数を別の実行装置に移動させるステップをさらに含むことができ、すなわち、ある関数がアクセラレータ内に必要ではなく、したがってその関数がCPUにおいて実行されるように再配置されることを決定することができる。

30

【0038】

なお、1つ以上のLE 120、1つ以上のMUX 130、および1つ以上のフィーダ140は、ハードウェア、ソフトウェア、ファームウェア、またはこれらの任意の組合せにおいて実施してよいことをさらに理解されたい。例示的な一実施形態においては、計算グリッド110（およびそのさまざまな素子）は、半導体デバイスとして実施される。別の実施形態においては、計算グリッド110は、一部が半導体デバイスとして実施され、一部がソフトウェア、ファームウェア、またはこれらの組合せとして実施される。一例として、より高い確率で、または多くの回数出現するものと判定される論理経路を、半導体デバイスにおいて実施することができ、その一方で、より低い確率で出現する、またはめったに出現しないと判定される論理経路を、ソフトウェアにおいて実施することができる。なお、このような実施形態は、同種の実施に頼る別の実施形態よりも、例えば実行時間や電力消費量における全体的な予測見返りコストが低くなりうることに留意されたい。論理経路は、処理ユニット101によって、または計算グリッド110内の専用ハードウェア素子によって検出することができる。

40

【0039】

図2Aおよび図2Bは、それぞれ、一実施形態に係る、特定の関数の再配置の例としてのシミュレーション200Aおよび200Bである。再配置処理は、頻繁に互いを呼び出す2つの論理上の計算グループ（基本ブロックまたは関数など）が、計算グリッド内で互いに物理的に近くに位置するように、計算グリッドを構成または操作するステップを含む。

【0040】

50

図 2 A および図 2 B に示した例としてのシミュレーションには、複数の LE 220 - 1 ~ 220 - N を有する計算グリッド 210 を示してあり、N は 1 より大きいかまたは 1 に等しい整数である（グリッド 210 には、単に図面を簡略化する目的で MUX は含まれていない）。計算グループが、射影グラフに基づいてそれぞれの LE 220 にマッピングまたは射影されている。LE 220 - 1, 220 - 2, 220 - 3 を含む計算グループが、LE 220 - 4, 220 - 5, 220 - 6 を含む計算グループに関連付けられることを識別した時点で、図 2 B に示したように、LE 220 - 1, 220 - 2, 220 - 3 を含む計算グループが、LE 220 - 4, 220 - 5, 220 - 6 を含む計算グループの近傍に再配置される。

【0041】

図 2 C および図 2 D は、それぞれ、特定の関数の複製の、例としてのシミュレーション 200 C および 200 D である。複製処理は、計算の速度を高め、かつボトルネックを広げる目的で、頻繁に呼び出される計算グループにマッピングされている 1 つ以上の LE 220 がグリッド 210 上で何度か複製されるように、計算グリッド 210 を構成するステップを含む。

【0042】

例としての一実施形態においては、複製処理は、新たに複製されたインスタンスを対象とする負荷分散処理をさらに含むことができ、これは例えば、呼出し元からの元のブロックの参照を、呼出し元がすべての複製を均一に呼び出すように再構成することによる。別の例としては、呼出し元からの参照が、複製を負荷分散する、新たに構成可能な MUX に再構成される。

【0043】

頻繁に呼び出される関数の計算グループを複製することによって計算グリッド 210 を最適化することができると判定されたとき、そのような各グループが複製され、元のグループから遠く離して再配置される。図 2 C および図 2 D に示した例では、LE 220 - 4, 220 - 5, 220 - 6 を含む計算グループが複製される。複製されたグループは、このような関数の呼出しを強化する目的で、計算グリッド 210 上で元のグループの LE から遠く離して配置される。このことは図 2 D に示してあり、図 2 D は、図 2 C に示した配置との比較における、複製されたグループを含む。複製されたグループを、他のグループより低い頻度で呼び出される関数のグループに置き換えることができる。これに代えて、複製された論理計算要素が性能に必要なものと識別されたとき、操作は縮小処理、すなわち複製解放処理を含むことができ、したがってグリッドから削除される。

【0044】

図 3 は、一実施形態に係る、構成可能な処理アーキテクチャにおいて動作する計算グリッドを実行時に最適化する方法を示す、例としての流れ図 300 である。S310 においては、少なくとも 1 つの関数を実行するための複数の要求を受け取る。これらの要求は、演算を実行するためにインタフェースを介して受け取る呼出しとすることができる。一例として、要求された演算は、例えばプログラムメモリ（例：図 1 B のプログラムメモリ 102）からのメモリのロードとすることができる。

【0045】

S320 においては、受け取ったシステムコールを分析して、2 つ以上のシステムコールの間で少なくとも 1 つのパターンが識別されるかをチェックする。この識別は、呼出しを次々と受け取ったことの識別、または呼出しを次々と受け取らなかったことの識別とすることができる。一例として、パターンは、2 つの関数が次々に実行されることを示す。このようなパターンが識別される場合、実行は S330 に続き、そうでない場合、実行は S310 に戻る。

【0046】

S330 においては、識別されたパターンを分析して、計算グリッドの修正または再構成が必要であるかを判定する。例えば、分析は、例えばある関数を実行するためにすでに割り当てられた計算グループがグリッド内に存在することをパターンが示しているか判定

10

20

30

40

50

するステップを含む。別の例としては、分析は、関数の再出現をパターンが示しているか判定するステップを含む。別の例としては、分析は、1つの関数が別の関数を呼び出すことをパターンが示しているか判定する。

【0047】

S340においては、分析に基づいて、計算グリッドの操作が必要であることをチェックする。必要である場合、実行はS350に続く。そうでない場合、実行はS360に続く。

【0048】

S350においては、受け取った呼出しによって呼び出される関数を最適化するために、計算グリッドを修正する、または別の方法で再構成する。計算グリッドの修正または再構成の一部として実行することのできる処理のさまざまな例は、上に説明してある。

【0049】

S360においては、さらなる要求を受け取っているかをチェックし、受け取っている場合、実行はS310に続き、そうでない場合、実行は終了する。

【0050】

本明細書に開示されている実施形態は、ハードウェア、ファームウェア、ソフトウェア、またはこれらの任意の組合せとして実施することができる。任意の適切なアーキテクチャを備えたマシンに、アプリケーションプログラムをアップロードして実行することができる。マシンは、1つまたは複数の中央処理装置（「CPU」）、メモリ、および入力/出力インタフェースなどのハードウェアを有するコンピュータプラットフォーム上に実施されることが好ましい。

【0051】

コンピュータプラットフォームは、オペレーティングシステムおよびマイクロ命令コードをさらに含むことができる。本明細書に記載されているさまざまな処理および関数は、CPU（そのようなコンピュータまたはプロセッサが明示的に示されているか否かを問わない）によって実行することのできるマイクロ命令コードの一部またはアプリケーションプログラムの一部のいずれか、またはこれらの任意の組合せとすることができる。

【0052】

これに加えて、コンピュータプラットフォームには、追加のネットワークファブリック記憶装置（network fabric storage unit）や印刷装置など、さまざまな他の周辺装置を接続することができる。さらに、非一過性のコンピュータ可読媒体は、一過性の伝搬信号を除く任意のコンピュータ可読媒体である。

【0053】

本明細書に記載されているすべての例、および条件を表す語（conditional language）は、この技術分野を発展させるために本発明者によって提供される本開示およびコンセプトの原理を読み手が理解できるように支援するための教示目的を意図しており、そのような具体的に記載されている例および条件に限定されないものと解釈されたい。

【0054】

なお、本明細書において「第1の」、「第2の」などの指定を使用して要素が言及されている場合、一般的には、それによって、それらの要素の数量または順序が制限されることはない。そうではなく、本明細書では、これらの指定は、一般に2つ以上の要素、または要素の2つ以上のインスタンスを区別する簡便な方法として使用されている。したがって、第1の要素および第2の要素という表現は、2つの要素のみを採用してよいことを意味するのではなく、何らかの形で第1の要素が第2の要素に先行しなければならないことを意味する。また、特に明記されていない限り、要素のセットは、1つまたは複数の要素を含む。

【0055】

本明細書において使用されているとき、後ろに項目のリストが続く句「の少なくとも1つ」は、リストされている項目のいずれかを個別に利用することができる、またはリストされている項目の2つ以上の任意の組合せを利用できることを意味する。例えば、システムが、「A、B、Cの少なくとも1つ」を含むものと説明されている場合、そのシステム

10

20

30

40

50

は、Aのみを含む、Bのみを含む、Cのみを含む、AとBを組合せで含む、BとCを組合せで含む、AとCを組合せで含む、またはAとBとCを組合せで含む、ことができる。

【 0 0 5 6 】

本出願は、米国仮特許出願第 6 2 / 5 4 0 , 8 4 9 号（出願日：2 0 1 7 年 8 月 3 日）および米国仮特許出願第 6 2 / 5 5 8 , 0 9 0 号（出願日：2 0 1 7 年 9 月 1 3 日）の利益を主張し、これらの文書の内容は参照により本明細書に組み込まれている。

10

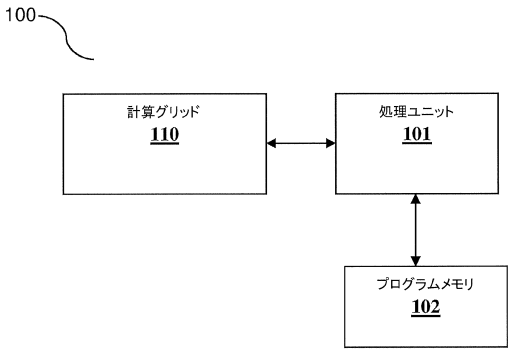
20

30

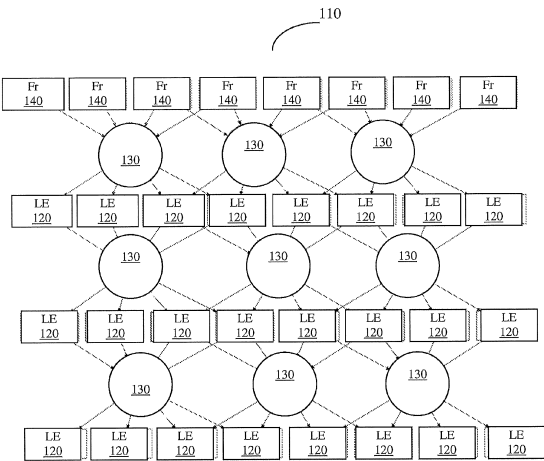
40

50

【図面】
【図 1 A】



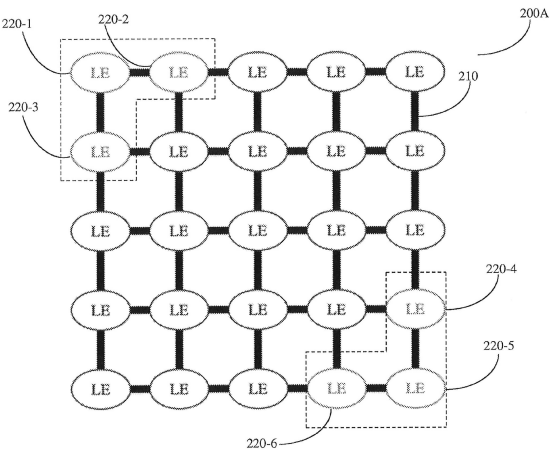
【図 1 B】



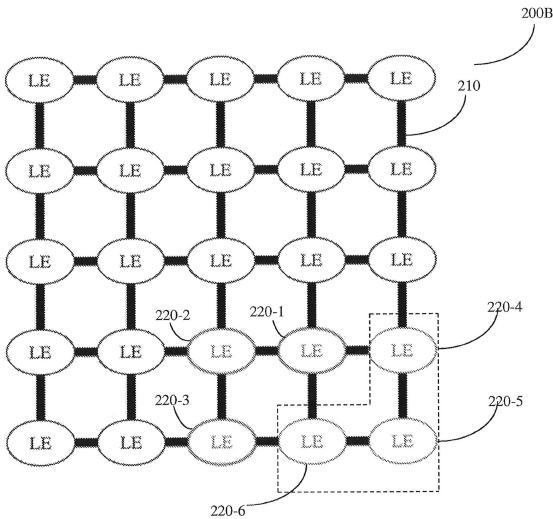
10

20

【図 2 A】



【図 2 B】

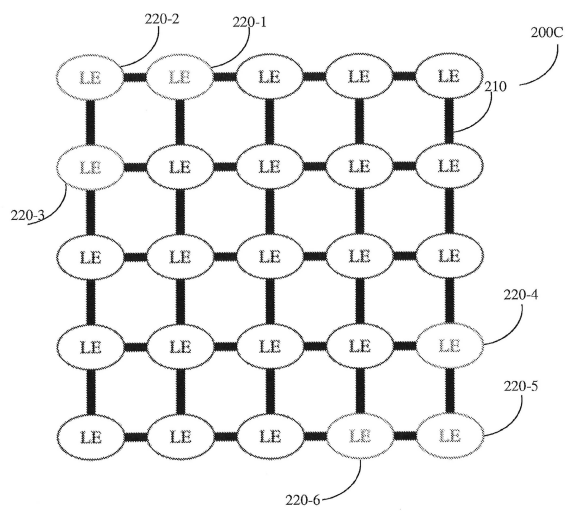


30

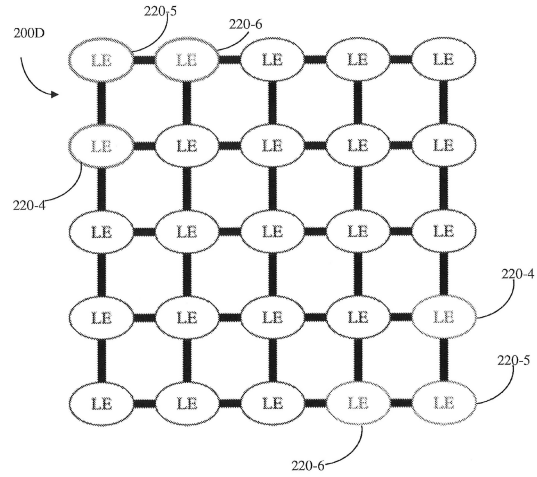
40

50

【図 2 C】



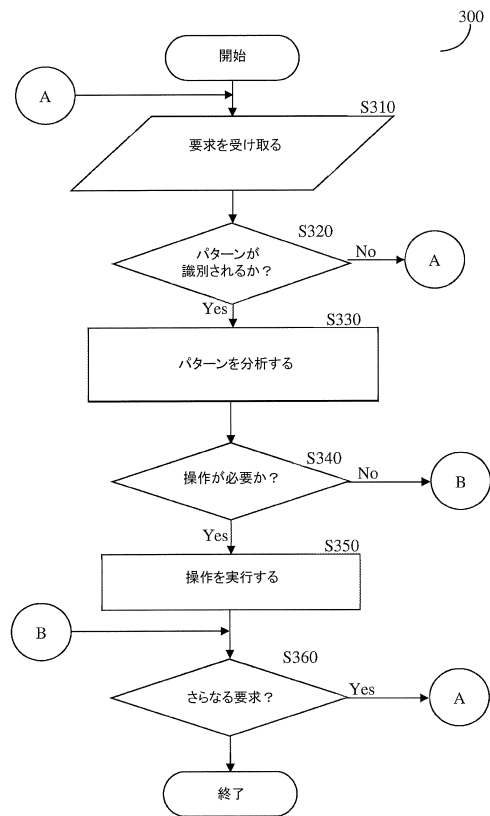
【図 2 D】



10

20

【図 3】



30

40

50

フロントページの続き

(33)優先権主張国・地域又は機関

米国(US)

(56)参考文献

特開 2 0 0 9 - 1 6 3 3 2 8 (J P , A)

特表 2 0 0 2 - 5 2 6 8 2 6 (J P , A)

国際公開第 2 0 1 7 / 0 2 9 7 4 3 (W O , A 1)

特開 2 0 0 9 - 2 3 8 2 2 1 (J P , A)

特表 2 0 0 5 - 5 0 5 0 3 0 (J P , A)

国際公開第 2 0 0 8 / 0 2 6 7 3 1 (W O , A 1)

佐野 雅彦 他, F P G A化のための k - a r y n - c u b e 型相互結合網用のルータの設計, 情報処理学会研究報告, 社団法人情報処理学会, 1996年05月16日, 第96巻 第39号, 第31頁 - 第36頁

(58)調査した分野 (Int.Cl., D B 名)

G 0 6 F 9 / 3 8

G 0 6 F 1 5 / 8 0