

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7400169号  
(P7400169)

(45)発行日 令和5年12月19日(2023.12.19)

(24)登録日 令和5年12月11日(2023.12.11)

|                       |                |         |      |         |
|-----------------------|----------------|---------|------|---------|
| (51)国際特許分類            | F I            |         |      |         |
| G 0 6 F               | 9/50 (2006.01) | G 0 6 F | 9/50 | 1 5 0 A |
| G 0 6 F               | 9/38 (2018.01) | G 0 6 F | 9/38 | 3 1 0 J |
| G 0 6 F               | 9/48 (2006.01) | G 0 6 F | 9/38 | 3 7 0 C |
|                       |                | G 0 6 F | 9/38 | 3 7 0 X |
|                       |                | G 0 6 F | 9/48 | 3 0 0 G |
| 請求項の数 14 外国語出願 (全40頁) |                |         |      |         |

|                   |                             |          |                         |
|-------------------|-----------------------------|----------|-------------------------|
| (21)出願番号          | 特願2020-104328(P2020-104328) | (73)特許権者 | 591003943               |
| (22)出願日           | 令和2年6月17日(2020.6.17)        |          | インテル・コーポレーション           |
| (65)公開番号          | 特開2021-34020(P2021-34020A)  |          | アメリカ合衆国 9 5 0 5 4 カリフォル |
| (43)公開日           | 令和3年3月1日(2021.3.1)          |          | ニア州・サンタクララ・ミッション カ      |
| 審査請求日             | 令和4年6月13日(2022.6.13)        |          | レッジ ブレーバード・2 2 0 0      |
| (31)優先権主張番号       | 16/542,012                  | (74)代理人  | 110000877               |
| (32)優先日           | 令和1年8月15日(2019.8.15)        |          | 弁理士法人R Y U K A国際特許事務所   |
| (33)優先権主張国・地域又は機関 | 米国(US)                      | (72)発明者  | マイケル ベハー                |
|                   |                             |          | アメリカ合衆国 9 5 0 5 4 カリフォル |
|                   |                             |          | ニア州・サンタクララ・ミッション カ      |
|                   |                             |          | レッジ ブレーバード・2 2 0 0 インテ  |
|                   |                             |          | ル・コーポレーション内             |
|                   |                             | (72)発明者  | モシェ マオル                 |
|                   |                             |          | アメリカ合衆国 9 5 0 5 4 カリフォル |
|                   |                             |          | ニア州・サンタクララ・ミッション カ      |
|                   |                             |          | 最終頁に続く                  |

(54)【発明の名称】 ワークロードのスタティックマッピングの順不同にパイプライン化された実行を可能にする方法及び装置

(57)【特許請求の範囲】

【請求項1】

第1のローカルクレジットマネージャーを有する第1のコンピューティングユニットであって、前記第1のコンピューティングユニットは、第1のバッファにデータを書き込むタスクを処理する、前記第1のコンピューティングユニットと、

第2のローカルクレジットマネージャーを有する第2のコンピューティングユニットであって、前記第2のコンピューティングユニットは、第2のバッファからデータを読み出すタスクを処理する、前記第2のコンピューティングユニットと、

前記第1のコンピューティングユニットと前記第2のコンピューティングユニットに結合された少なくとも1つのファブリックと、

前記少なくとも1つのファブリックに結合された中央クレジットマネージャーとを備え、前記中央クレジットマネージャーは、

前記第1のローカルクレジットマネージャーへの第1のクレジットの送信を生じさせることであって、前記第1のクレジットは、前記第1のバッファに格納される第2のデータを生成すべく前記第1のコンピューティングユニットにより処理される第1のデータに対応する、前記送信を生じさせることと、

前記第1のコンピューティングユニットの前記第1のローカルクレジットマネージャーから前記第1のクレジットにアクセスすることと、

前記第2のコンピューティングユニットに対するクレジットのカウントを減らすこととを行う

装置。

【請求項 2】

前記中央クレジットマネージャーは、前記第 1 のコンピューティングユニットが前記第 1 のデータを処理することに応じて、前記第 1 のコンピューティングユニットの前記第 1 のローカルクレジットマネージャーから前記第 1 のクレジットにアクセスする

請求項 1 に記載の装置。

【請求項 3】

前記中央クレジットマネージャーは、前記第 2 のバッファにおける前記第 2 のデータの利用可能性に応じて、前記第 2 のコンピューティングユニットに対する前記クレジットのカウントを減らす

請求項 1 又は 2 に記載の装置。

【請求項 4】

前記第 2 のコンピューティングユニットに対する前記クレジットのカウントは、クレジットの第 1 のカウントであり、

前記中央クレジットマネージャーは、

前記第 1 のコンピューティングユニットに対するクレジットの第 2 のカウントを初期化し、

前記第 2 のコンピューティングユニットに対するクレジットの前記第 1 のカウントを初期化する

請求項 1 から 3 のいずれか一項に記載の装置。

【請求項 5】

前記中央クレジットマネージャーは、前記第 1 のデータが前記第 1 のコンピューティングユニットに割り当てられるタスクに関連付けられることに基づいて、前記第 1 のローカルクレジットマネージャーへの前記第 1 のクレジットの送信を生じさせる

請求項 1 から 4 のいずれか一項に記載の装置。

【請求項 6】

コンピューティングシステムにおける少なくとも一つのプロセッサによって命令が実行されたとき、前記コンピューティングシステムにおける中央クレジットマネージャーが、第 1 のコンピューティングユニットの第 1 のローカルクレジットマネージャーへ、第 1 のクレジットを送信する段階であって、前記第 1 のコンピューティングユニットは、第 1 のバッファにデータを書き込むタスクを処理し、前記第 1 のクレジットは、前記第 1 のバッファに格納される第 2 のデータを生成すべく前記第 1 のコンピューティングユニットにより処理される第 1 のデータに対応する、段階と、

前記少なくとも一つのプロセッサによって命令が実行されたとき、前記中央クレジットマネージャーが、前記第 1 のコンピューティングユニットの前記第 1 のローカルクレジットマネージャーから前記第 1 のクレジットにアクセスする段階と、

前記少なくとも一つのプロセッサによって命令が実行されたとき、前記中央クレジットマネージャーが、第 2 のローカルクレジットマネージャーを有する第 2 のコンピューティングユニットに対するクレジットのカウントを減らす段階であって、前記第 2 のコンピューティングユニットは、第 2 のバッファからデータを読み出すタスクを処理する、段階と、

を備える

方法。

【請求項 7】

前記第 1 のコンピューティングユニットが前記第 1 のデータを処理することに応じて、前記中央クレジットマネージャーが、前記第 1 のコンピューティングユニットの前記第 1 のローカルクレジットマネージャーから前記第 1 のクレジットにアクセスする段階を更に

備える

請求項 6 に記載の方法。

【請求項 8】

前記第 2 のバッファにおける前記第 2 のデータの利用可能性に応じて、前記中央クレジ

10

20

30

40

50

ットマネージャーが、前記第 2 のコンピューティングユニットに対する前記クレジットの  
カウントを減らす段階を更に備える

請求項 6 又は 7 に記載の方法。

【請求項 9】

前記第 2 のコンピューティングユニットに対する前記クレジットのカウントは、クレジ  
ットの第 1 のカウントであり、

前記方法は、

前記中央クレジットマネージャーが、前記第 1 のコンピューティングユニットに対するク  
レジットの第 2 のカウント、及び前記第 2 のコンピューティングユニットに対するクレジ  
ットの前記第 1 のカウントを初期化する段階を更に備える

請求項 6 から 8 のいずれか一項に記載の方法。

【請求項 10】

前記第 1 のデータが前記第 1 のコンピューティングユニットに割り当てられるタスクに  
関連付けられることに基づいて、前記中央クレジットマネージャーが、前記第 1 のローカ  
ルクレジットマネージャーへ、前記第 1 のクレジットを送信する段階を更に備える

請求項 6 から 9 のいずれか一項に記載の方法。

【請求項 11】

メモリと、複数の命令と、請求項 6 から 10 のいずれか一項に記載の方法を実行するべ  
く、前記複数の命令を実行する前記少なくとも一つのプロセッサとを備える装置。

【請求項 12】

命令を備え、前記命令は、実行されると、請求項 6 から 10 のいずれか一項に記載の方  
法を前記少なくとも一つのプロセッサに実行させる

非一時的コンピュータ可読媒体。

【請求項 13】

請求項 6 から 10 のいずれか一項に記載の方法を実行するための手段を備える装置。

【請求項 14】

請求項 6 から 10 のいずれか一項に記載の方法を前記少なくとも一つのプロセッサに実  
行させるコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この開示は概して、処理に関し、より詳細にはワークロードのスタティックマッピング  
の順不同にパイプライン化された実行を可能にする方法及び装置に関する。

【背景技術】

【0002】

コンピュータハードウェア製造者は、コンピュータプラットフォームの様々なコンポー  
ネントに用いられるハードウェアコンポーネントを開発する。例えば、コンピュータハー  
ドウェア製造者は、マザーボード、マザーボード用のチップセット、中央処理ユニット（  
CPU）、ハードディスクドライブ（HDD）、ソリッドステートドライブ（SSD）及  
び、他のコンピュータコンポーネントを開発する。更に、コンピュータハードウェア製造  
者は、アクセラレータとして知られる、ワークロードの処理を加速する処理要素を開発す  
る。例えば、アクセラレータは、CPU、グラフィック処理ユニット（GPU）、ビジョ  
ン処理ユニット（VPU）及び/又はフィールドプログラマブルゲートアレイ（FPGA  
）などであり得る。

【図面の簡単な説明】

【0003】

【図 1】異種システムのアクセラレータで実行されるワークロードを表したグラフの図で  
ある。

【0004】

【図 2】パイプライン及びバッファを実装した異種システムのアクセラレータで実行され

10

20

30

40

50

るワークロードを表したグラフの図である。

【 0 0 0 5 】

【図 3】本開示の教示に従って構築された例示的なコンピューティングシステムを示すブロック図である。

【 0 0 0 6 】

【図 4】例示的な 1 又は複数のスケジューラを含む例示的なコンピューティングシステムを示すブロック図である。

【 0 0 0 7 】

【図 5】図 3 及び 4 の 1 又は複数のスケジューラを実装し得る例示的なスケジューラのブロック図である。

【 0 0 0 8 】

【図 6】図 5 のバッファクレジット記録装置のさらなる詳細を示す例示的なスケジューラのブロック図である。

【 0 0 0 9 】

【図 7】パイプライン及びバッファを実装した異種システムのアクセラレータで実行するワークロードを表す例示的なグラフの図である。

【 0 0 1 0 】

【図 8】図 5 のスケジューラ及び / 又は図 6 のスケジューラを実装するために実行できる機械可読命令によって実装され得る処理を表すフローチャートである。

【 0 0 1 1 】

【図 9】図 5 のスケジューラ及び / 又は図 6 のスケジューラの 1 又は複数のインスタンス化を実装するための図 8 の命令を実行するよう構築された例示的なプロセッサプラットフォームのブロック図である。

【 0 0 1 2 】

図は縮尺通りではない。概して、同じもの又は一部のようなものを指すべく、図面及び付随する記述説明全体で同じ参照が用いられるであろう。接続についての言及（例えば、取り付け、結合、接続及び結合）は広く解釈されるべきであり、そうでないと示していない限り、要素の集合の間の中間部材及び要素の間の相対的な移動を含んでよい。従って、接続についての言及は、2 つの要素が直接接続されたり互いに固定された関係であることを必ずしも推論されるものではない。

【 0 0 1 3 】

「第 1」、「第 2」、「第 3」等の記述子は、別個に称される複数の要素又はコンポーネントを識別する場合に本明細書で用いられる。用いられるそれらの文脈に基づいて特定又は理解されるのでない限り、そのような記述子は、優先性、物理的順序若しくはリストの配置、又は、時間的な順序のいかなる意味を負わせることを意図しておらず、開示された例の理解の簡略化のために、複数の要素又はコンポーネントを別個に参照するためのラベルとして、単に用いられている。いくつかの例において、ある要素を指すのに「第 1」という記述子が詳細な説明で用いられる一方で、同じ要素が請求項で「第 2」又は「第 3」のような異なる記述子で称されてよい。このような例において、そのような記述子は単に、複数の要素又はコンポーネントの参照を簡略化するために用いられていると理解されるべきである。

【発明を実施するための形態】

【 0 0 1 4 】

多くのコンピュータハードウェア製造者は、アクセラレータとして知られる、ワークロードの処理を加速する処理要素を開発する。例えば、アクセラレータは、中央処理ユニット（CPU）、グラフィック処理ユニット（GPU）、ビジョン処理ユニット（VPU）及び / 又はフィールドプログラマブルゲートアレイ（FPGA）であり得る。さらに、アクセラレータは、ワークロードの任意のタイプを処理可能であるが、ワークロードの特定のタイプを最適化するように設計される。例えば、CPU 及び FPGA はより汎用の処理を扱うよう設計され得るが、GPU は、ビデオ、ゲーム及び / 又は他の物理及び数学に基

10

20

30

40

50

づく計算の処理を向上するよう設計され得るとともに、VPUは、マシンビジョンタスクの処理を向上するよう設計され得る。

【0015】

更に、いくつかのアクセラレータは、人工知能(AI)アプリケーションの処理を特に向上するよう設計される。VPUはAIアクセラレータの特定のタイプであるが、多くの異なるAIアクセラレータが用いられ得る。実際、多くのAIアクセラレータは、特定用途向け集積回路(ASIC)によって実装され得る。このようなASICベースのAIアクセラレータは、機械学習(ML)、深層学習(DL)、及び/又は、サポートベクタマシン(SVM)、ニューラルネットワーク(NN)、リカレントニューラルネットワーク(RNN)、畳み込みニューラルネットワーク(CNN)、ロングショートタームメモリ(LSTM)、ゲートリカレントユニット(GRU)等を含む他の人工機械駆動ロジックのようなAIの特定のタイプに関するタスクの処理を向上するよう設計され得る。

10

【0016】

コンピュータハードウェア製造者は、1つより多いタイプの処理要素を含む異種システムもまた開発している。例えば、コンピュータハードウェア製造者は、CPUのような汎用の処理要素と、FPGAのような汎用アクセラレータ、及び/又は、GPU、VPU及び/又は他のAIアクセラレータのようなより調整されたアクセラレータのいずれかとの両方を組み合わせてよい。このような異種システムは、システムオンチップ(SoC)として実装され得る。

【0017】

20

開発者が異種システム上で機能、アルゴリズム、プログラム、アプリケーション及び/又は他のコードを動作させることを望む場合、開発者及び/又はソフトウェアは、コンパイル時に、機能、アルゴリズム、プログラム、アプリケーション及び/又は他のコードのためのスケジュールを生成する。一旦スケジュールが生成されると、スケジュールは、実行可能ファイルを生成するために(アヘッドオブタイム又はジャストインタイムのいずれかのパラダイムで)、機能、アルゴリズム、プログラム、アプリケーション及び/又は他のコードの仕様と組み合わせられる。さらに、機能、アルゴリズム、プログラム、アプリケーション及び/又は他のコードはノードを含むグラフとして表されてよく、ここで、グラフはワークロードを示し、各ノードはそのワークロードの特定のタスクを示す。さらに、グラフ内の異なるノード間の接続は、特定のノードが実行されるために必要なデータ入力及び/又は出力を示し、グラフの頂点はグラフのノード間のデータ依存性を示す。

30

【0018】

実行可能ファイルは多数の異なる実行可能なセクションを含み、ここで、各実行可能なセクションは特定の処理要素(例えば、CPU、GPU、VPU及び/又はFPGA)によって実行可能である。実行可能ファイルの各実行可能なセクションは実行可能なサブセクションをさらに含んでよく、ここで、各実行可能なサブセクションは特定の処理要素の計算ビルディングブロック(CBB)によって実行可能である。更に又は代替的に、本明細書で開示されるいくつかの例において、開発者及び/又はソフトウェア開発用ソフトウェアは、実行ファイルの実行の成功を判断する基準を定め得る(例えば成功基準)。例えば、このような成功基準は、異種システム及び/又は特定の処理要素の利用の閾値と満たし及び/又はそうでなければ満足するよう実行ファイルを実行することに対応してよい。他の例において、成功基準は、閾値量の時間に実行ファイルを実行することに対応してよい。しかしながら、異種システム及び/又は特定の処理要素でどのように実行ファイルを実行するかを判断する場合に任意の適切な成功機能が用いられてよい。このように、成功基準は、開発者、ソフトウェア及び/又は人工知能システムが成功基準を満たすよう最適化されたスケジュールを含む実行ファイルを生成するのに有益であり得る。

40

【0019】

図1は、異種システムのアクセラレータで実行されるワークロードのグラフ100を表す図である。グラフ100は、第1ワークロードノード102(WN[0])、第2ワークロードノード104(WN[1])、第3ワークロードノード106(WN[2])、

50

第4ワークロードノード108(WN[3])及び第5ワークロードノード110(WN[4])を含む。図1において、アクセラレータは、スタティックソフトウェアスケジューラでグラフ100によって表されるワークロードを行っている。スタティックソフトウェアスケジューリングは、アクセラレータの計算ビルディングブロック(CBB)上でグラフ100の異なるワークロードノードを実行するための予め定義された態様を決定することを含む。例えば、スタティックソフトウェアスケジューラは、第1ワークロードノード102(WN[0])を第1CBB112に、第2ワークロードノード104(WN[1])を第2CBB114に、第3ワークロードノード106(WN[2])を第3CBB116に、第4ワークロードノード108(WN[3])を第4CBB118に、第5ワークロードノード110(WN[4])を第2CBB114に割り当てる。

10

#### 【0020】

図1において、スタティックソフトウェアスケジューラは、第4CBB118で実行する第4ワークロードノード108(WN[3])と並列して第1ワークロードノード102(WN[0])が第1CBB112で実行するという枠組みを作っている。図1において、第1CBB112が第1ワークロードノード102(WN[0])を実行するより速く、第4CBB118は第4ワークロードノード108(WN[3])を実行する。スタティックソフトウェアスケジューラが、第2CBB114が第5ワークロードノード110(WN[4])を実行する前に第2CBB114が第2ワークロードノード104(WN[1])を実行するという枠組みを作っているように、第1CBB112が第1ワークロードノード102(WN[0])の実行を完了するまで、第2CBB114はアイドル状態である。さらに、次のワークロードノードの実行前にワークロードノードが全て実行されるまで待つことは、著しいメモリオーバーヘッドを必要とする。というのは、CBBが2番目のワークロードノード(例えば第2ワークロードノード104(WN[1]))を実行し得る前に、CBBによる1番目のワークロードノード(例えば第1ワークロードノード102(WN[0]))の実行で生成されたデータをアクセラレータに格納することが必要とされるからである。

20

#### 【0021】

図2は、パイプライン及びバッファを実装している異種システムのアクセラレータで実行するワークロードを表すグラフ200の図である。グラフ200は、第1ワークロードノード102(WN[0])、第2ワークロードノード104(WN[1])、第3ワークロードノード106(WN[2])、第4ワークロードノード108(WN[3])及び第5ワークロードノード110(WN[4])を含む。図2において、アクセラレータは、スタティックソフトウェアスケジューラでグラフ200によって表されるワークロードを行っている。図2のスタティックソフトウェアスケジューラは、パイプラインを実装するとともに第1バッファ202、第2バッファ204及び第3バッファ206を含むアクセラレータのCBBでのグラフ200の異なるワークロードノードに対する実行スケジュールの枠組みを作っている。更に、スタティックソフトウェアスケジューラは、第1ワークロードノード102(WN[0])を第1CBB112に、第2ワークロードノード104(WN[1])を第2CBB114に、第3ワークロードノード106(WN[2])を第3CBB116に、第4ワークロードノード108(WN[3])を第4CBB118に、第5ワークロードノード110(WN[4])を第2CBB114に割り当てる。第1バッファ202は第1CBB112及び第2CBB114と結合し、第2バッファ204は第2CBB114及び第3CBB116と結合し、第3バッファ206は第2CBB114及び第4CBB118と結合する。

30

40

#### 【0022】

バッファ202、204及び206によって、スタティックソフトウェアスケジューラが、ある時間間隔内でワークロードノードの全体を実行するよりむしろ各CBBがワークロードノードの一部(例えばタイル)をその時間間隔内で処理する枠組みを作るのが可能となる。同様に、スタティックソフトウェアスケジューラは、ワークロードのそのような一部が利用可能となった場合に、他のCBB(例えばコンシューマー)によって生成され

50

たデータを処理している C B B がワークロードノードの一部（例えばタイル）を実行し得る枠組みを作り得る。しかしながら、ワークロードノードを実行している C B B は利用可能なデータを処理して新たなデータをメモリに書き込むので、C B B で所与のワークロードノードを実行するためには、ランタイムにおいて閾値量のデータが利用可能でなければならない。ランタイムにおいて結果を書き込むメモリ内の閾値量のスペースがなければならない。バッファは基本的なスタティックソフトウェアスケジューリングによってメモリのオーバーヘッドを減少させるが、それはランタイムにおいてデータ利用可能性及びノ又は依存性に高く依存するので、バッファでスタティックソフトウェアスケジューリングの枠組みを作ることはますます難しい。さらに、アクセラレータ全体の負荷はアクセラレータ上の各 C B B の処理速度に影響し得るので、所与のアクセラレータの C B B を効果的に利用するスタティックソフトウェアスケジューリングを開発するのは難しい。

10

#### 【0023】

本明細書で開示された例は、ワークロードのスタティックマッピングの順不同にパイプライン化された実行を可能にする方法及び装置を含む。スタティックソフトウェアスケジューリングとは対照的に、本明細書で開示された例は、予め定められたスタティックソフトウェアスケジューリングには依存しない。むしろ、本明細書で開示された例は、アクセラレータ及びノ又は他の処理要素上の利用可能なデータ及び利用可能なメモリに基づいて、所与の C B B に割り当てられているワークロードノードのどれを行うかを決定する。さらに、各 C B B は、クレジットの第 1 の数で表される、第 1 バッファで利用可能な所与のワークロードに関連づけられたデータの量、及び、クレジットの第 2 の数で表される、第 2 バッファで利用可能なスペースの量を追跡する。これは、所与の C B B でのワークロードノードのダイナミックランタイムスケジューリングを可能にする。

20

#### 【0024】

ワークロードノードごとに、クレジットの第 1 の数が第 1 閾値を満たしかつクレジットの第 2 の数が第 2 閾値を満たす場合に、C B B はワークロードノードを実行し得る。これは、ワークロード全体の所与のグラフから独立した順不同にパイプライン化された実行を可能にする。本明細書で開示された例は、アクセラレータの 1 又は複数の計算ビルディングブロックにワークロードのスタティックマッピングの順不同にパイプライン化された実行を可能にする装置を提供する。例示的な装置は、クレジットの第 1 の数をメモリ内へ読み込むインターフェースと、クレジットの第 1 の数をバッファのメモリ利用可能性に関連付けられたクレジットの閾値数と比較する比較器と、クレジットの第 1 の数がクレジットの閾値数を満たす場合に、1 又は複数の計算ビルディングブロックの最初の一つで実行されるワークロードのワークロードノードを選択するディスパッチャとを含む。

30

#### 【0025】

図 3 は、本開示の教示に従い構築される例示的なコンピューティングシステム 300 を示すブロック図である。図 3 の例において、コンピューティングシステム 300 は、例示的なシステムメモリ 302 及び例示的な異種システム 304 を含む。例示的な異種システム 304 は、例示的なホストプロセッサ 306、例示的な第 1 通信バス 308、例示的な第 1 アクセラレータ 310 a、例示的な第 2 アクセラレータ 310 b 及び例示的な第 3 アクセラレータ 310 c を含む。例示的な第 1 アクセラレータ 310 a、例示的な第 2 アクセラレータ 310 b 及び例示的な第 3 アクセラレータ 310 c の各々は、いくつかはアクセラレータの演算に対して汎用的で、いくつかはそれぞれのアクセラレータの演算に対して特化した様々の C B B を含む。

40

#### 【0026】

図 3 の例において、システムメモリ 302 は異種システム 304 に結合される。システムメモリ 302 はメモリである。図 3 において、システムメモリ 302 は、ホストプロセッサ 306、第 1 アクセラレータ 310 a、第 2 アクセラレータ 310 b 及び第 3 アクセラレータ 310 c のうちの少なくとも 1 つの間での共有ストレージである。図 3 の例において、システムメモリ 302 はコンピューティングシステム 300 に位置する物理ストレージである。しかしながら、他の例において、システムメモリ 302 はコンピューティン

50

グシステム 300 の外部にあってよく及び / 又はそうでなければ離れていてよい。さらなる例において、システムメモリ 302 は仮想記憶装置であってよい。図 3 の例において、システムメモリ 302 は、永続ストレージ（例えば読み出し専用メモリ（ROM）、プログラマブル ROM（PROM）、消去可能 PROM（EPROM）、電氣的消去可能 PROM（EEPROM 等））である。他の例において、システムメモリ 302 は、永続基本入出力システム（BIOS）又はフラッシュストレージであってよい。さらなる例において、システムメモリ 302 は揮発性メモリであってよい。

【0027】

図 3 において、異種システム 304 はシステムメモリ 302 と結合される。図 3 の例において、異種システム 304 は、ホストプロセッサ 306、及び / 又は、第 1 アクセラレータ 310 a、第 2 アクセラレータ 310 b 又は第 3 アクセラレータ 310 c の 1 又は複数でワークロードを実行することによって、ワークロードを処理する。図 3 において、異種システム 304 は SOC である。代替的に、異種システム 304 はいかなるその他のタイプのコンピューティング又はハードウェアシステムであってよい。

【0028】

図 3 の例において、ホストプロセッサ 306 は、コンピュータ又はコンピューティングデバイス（例えばコンピューティングシステム 300）に関連付けられた演算の完了の実行、遂行及び / 又は促進するための命令（例えば機械可読命令）を実行する処理要素である。図 3 の例において、ホストプロセッサ 306 は、異種システム 304 にとって基本の処理要素であり、かつ、少なくとも 1 つのコアを含む。代替的に、ホストプロセッサ 306 は、（例えば 1 つより多い CPU が用いられる例において）共同した一次的な処理要素であってよいが、他の例において、ホストプロセッサ 306 は、二次的な処理要素であってよい。

【0029】

図 3 に図示の例において、第 1 アクセラレータ 310 a、第 2 アクセラレータ 310 b 及び / 又は第 3 アクセラレータ 310 c の 1 又は複数は、ハードウェアアクセラレーションのようなコンピューティングタスクのための異種システム 304 で実行するプログラムによって利用されてよい処理要素である。例えば、第 1 アクセラレータ 310 a は、AI に対するマシンビジョンタスク（例えば VPU）を処理する処理速度及び全体性能を向上するように設計され及び / 又はそうでなければ構成され若しくは構築された処理リソースを含む処理要素である。

【0030】

本明細書で開示された例において、ホストプロセッサ 306、第 1 アクセラレータ 310 a、第 2 アクセラレータ 310 b 及び第 3 アクセラレータ 310 c の各々は、コンピューティングシステム 300 及び / 又はシステムメモリ 302 の他の要素と通信する。例えば、ホストプロセッサ 306、第 1 アクセラレータ 310 a、第 2 アクセラレータ 310 b、第 3 アクセラレータ 310 c 及び / 又はシステムメモリ 302 は第 1 通信バス 308 で通信する。本明細書で開示されたいくつかの例において、ホストプロセッサ 306、第 1 アクセラレータ 310 a、第 2 アクセラレータ 310 b、第 3 アクセラレータ 310 c 及び / 又はシステムメモリ 302 は、任意の適切な有線及び / 又は無線通信システムで通信してよい。更に、本明細書で開示されたいくつかの例において、ホストプロセッサ 306、第 1 アクセラレータ 310 a、第 2 アクセラレータ 310 b、第 3 アクセラレータ 310 c 及び / 又はシステムメモリ 302 の各々は、任意の適切な有線及び / 又は無線通信システムで、コンピューティングシステム 300 の外部の任意のコンポーネントと通信してよい。

【0031】

図 3 の例において、第 1 アクセラレータ 310 a は、例示的な畳み込みエンジン 312、例示的な RNN エンジン 314、例示的なメモリ 316、例示的なメモリ管理ユニット（MMU）318、例示的な DSP 320、例示的なコントローラ 322 及び例示的なダイレクトメモリアクセス（DMA）ユニット 324 を含む。更に、例示的な畳み込みエン

10

20

30

40

50

ジン 3 1 2、例示的な R N N エンジン 3 1 4、例示的な D M A ユニット 3 2 4、例示的な D S P 3 2 0 及び例示的なコントローラ 3 2 2 及びの各々は、例示的な第 1 スケジューラ 3 2 6、例示的な第 2 スケジューラ 3 2 8、例示的な第 3 スケジューラ 3 3 0、例示的な第 4 スケジューラ 3 3 2 及び例示的な第 5 スケジューラ 3 3 4 をそれぞれ含む。例示的な D S P 3 2 0 及び例示的なコントローラ 3 2 2 の各々は更に、例示的な第 1 カーネルライブラリ 3 3 6 及び例示的な第 2 カーネルライブラリ 3 3 8 を含む。

【 0 0 3 2 】

図 3 に図示の例において、畳み込みエンジン 3 1 2 は畳み込みに関連したタスクの処理を向上させるよう構成されたデバイスである。さらに、畳み込みエンジン 3 1 2 は、視覚イメージの解析に関連するタスク及び / 又は C N N に関連する他のタスクの処理を向上させる。図 3 において、R N N エンジン 3 1 4 は R N N に関連するタスクの処理を向上するよう構成されたデバイスである。更に、R N N エンジン 3 1 4 は、セグメント化されていない繋がった手書き認識、音声認識の解析に関連するタスク及び / 又は R N N に関連する他のタスクの処理を向上させる。

10

【 0 0 3 3 】

図 3 の例において、メモリ 3 1 6 は、畳み込みエンジン 3 1 2、R N N エンジン 3 1 4、M M U 3 1 8、D S P 3 2 0、コントローラ 3 2 2 及び D M A ユニット 3 2 4 のうちの少なくとも 1 つの間の共有ストレージである。図 3 の例において、メモリ 3 1 6 は第 1 アクセラレータ 3 1 0 a に位置する物理ストレージである。しかしながら、他の例において、メモリ 3 1 6 は、第 1 アクセラレータ 3 1 0 a の外部にあってよく及び / 又はそうでなければ離れていてよい。さらなる例において、メモリ 3 1 6 は仮想記憶装置であってよい。図 3 の例において、メモリ 3 1 6 は、永続ストレージ（例えば R O M、P R O M、E P R O M、E E P R O M 等）である。他の例において、メモリ 3 1 6 は永続 B I O S 又はフラッシュストレージであってよい。さらなる例において、メモリ 3 1 6 は揮発性メモリであってよい。

20

【 0 0 3 4 】

図 3 に図示の例において、例示的な M M U 3 1 8 は、メモリ 3 1 6 及び / 又はシステムメモリ 3 0 2 のアドレスへの参照を含むデバイスである。M M U 3 1 8 は更に、畳み込みエンジン 3 1 2、R N N エンジン 3 1 4、D S P 3 2 0 及び / 又はコントローラ 3 2 2 の 1 又は複数によって用いられる仮想的メモリアドレスを、メモリ 3 1 6 及び / 又はシステムメモリ 3 0 2 内の物理アドレスへ変換する。

30

【 0 0 3 5 】

図 3 の例において、D S P 3 2 0 は、デジタル信号の処理を向上させるデバイスである。例えば、D S P 3 2 0 は、カメラ及び / 又はコンピュータビジョンに関する他のセンサからのデータのような、連続的な実世界の信号を測定、フィルタ及び / 又は圧縮する処理を促進する。図 3 において、コントローラ 3 2 2 は第 1 アクセラレータ 3 1 0 a の制御ユニットとして実装される。例えば、コントローラ 3 2 2 は、第 1 アクセラレータ 3 1 0 a の演算を管理する。いくつかの例において、コントローラ 3 2 2 はクレジットマネージャーを実装する。さらに、コントローラ 3 2 2 は、畳み込みエンジン 3 1 2、R N N エンジン 3 1 4、メモリ 3 1 6、M M U 3 1 8 及び / 又は D S P 3 2 0 の 1 又は複数に、ホストプロセッサ 3 0 6 から受信した機械可読命令にどのように応答するかを命令し得る。

40

【 0 0 3 6 】

図 3 に図示の例において、D M A ユニット 3 2 4 は、畳み込みエンジン 3 1 2、R N N エンジン 3 1 4、D S P 3 2 0 及びコントローラ 3 2 2 のうちの少なくとも 1 つに、ホストプロセッサ 3 0 6 から独立してシステムメモリ 3 0 2 にアクセスすることを可能にするデバイスである。例えば、D M A ユニット 3 2 4 は、アナログ又はデジタル回路、ロジック回路、プログラマブルプロセッサ、プログラマブルコントローラ、グラフィック処理ユニット（G P U）、デジタルシグナルプロセッサ（D S P）、特定用途向け集積回路（A S I C）、プログラマブル論理デバイス（P L D）及び / 又はフィールドプログラマブル論理デバイス（F P L D）の 1 又は複数によって実装され得る。

50

## 【 0 0 3 7 】

図 3 の例において、第 1 スケジューラ 3 2 6、第 2 スケジューラ 3 2 8、第 3 スケジューラ 3 3 0、第 4 スケジューラ 3 3 2 及び第 5 スケジューラ 3 3 4 の各々は、畳み込みエンジン 3 1 2、RNN エンジン 3 1 4、DMA ユニット 3 2 4、DSP 3 2 0 及びコントローラ 3 2 2 がそれぞれ、オフロードされていた及び / 又はそうでなければ第 1 アクセラレータ 3 1 0 a に送信されているワークロードの一部をいつ実行するかを決定するデバイスである。更に、第 1 カーネルライブラリ 3 3 6 及び第 2 カーネルライブラリ 3 3 8 の各々は、1 又は複数のカーネルを含むデータ構造である。第 1 カーネルライブラリ 3 3 6 及び第 2 カーネルライブラリ 3 3 8 のカーネルは、例えば、DSP 3 2 0 及びコントローラ 3 2 2 のそれぞれで高スループットのためにコンパイルされたルーチンである。カーネルは、例えば、コンピューティングシステム 3 0 0 で行われる実行ファイルの実行可能なサブセクションに対応する。

10

## 【 0 0 3 8 】

本明細書で開示された例において、畳み込みエンジン 3 1 2、RNN エンジン 3 1 4、メモリ 3 1 6、MMU 3 1 8、DSP 3 2 0、コントローラ 3 2 2 及び DMA ユニット 3 2 4 の各々は、第 1 アクセラレータ 3 1 0 a の他の要素と通信する。例えば、畳み込みエンジン 3 1 2、RNN エンジン 3 1 4、メモリ 3 1 6、MMU 3 1 8、DSP 3 2 0、コントローラ 3 2 2 及び DMA ユニット 3 2 4 は、例示的な第 2 通信バス 3 4 0 で通信する。いくつかの例において、第 2 通信バス 3 4 0 は、コンフィギュレーションアンドコントロール (CnC) ファブリック及びデータファブリックにより実装されてよい。本明細書で開示されたいくつかの例において、畳み込みエンジン 3 1 2、RNN エンジン 3 1 4、メモリ 3 1 6、MMU 3 1 8、DSP 3 2 0、コントローラ 3 2 2 及び DMA ユニット 3 2 4 は、任意の適切な有線及び / 又は無線通信システムで通信してよい。更に、本明細書に開示されたいくつかの例において、畳み込みエンジン 3 1 2、RNN エンジン 3 1 4、メモリ 3 1 6、MMU 3 1 8、DSP 3 2 0、コントローラ 3 2 2 及び DMA ユニット 3 2 4 の各々は、任意の適切な有線及び / 又は無線通信システムで第 1 アクセラレータ 3 1 0 a の外部の任意のコンポーネントと通信してよい。

20

## 【 0 0 3 9 】

前に言及したように、例示的な第 1 アクセラレータ 3 1 0 a、例示的な第 2 アクセラレータ 3 1 0 b 及び例示的な第 3 アクセラレータ 3 1 0 c の各々は、いくつかはアクセラレータの演算に対して汎用的で、いくつかはそれぞれのアクセラレータの演算に対して特化した様々の CBB を含む。例えば、第 1 アクセラレータ 3 1 0 a、第 2 アクセラレータ 3 1 0 b 及び第 3 アクセラレータ 3 1 0 c の各々は、メモリ、MMU、コントローラ、及び、CBB の各々に対するそれぞれのスケジューラのような汎用 CBB を含む。

30

## 【 0 0 4 0 】

図 3 の例において、第 1 アクセラレータ 3 1 0 a は VPU を実装し、かつ、畳み込みエンジン 3 1 2、RNN エンジン 3 1 4 及び DSP 3 2 0 (例えば第 1 アクセラレータ 3 1 0 a の演算に特化した CBB) を含み、第 2 アクセラレータ 3 1 0 b 及び第 3 アクセラレータ 3 1 0 c は第 2 アクセラレータ 3 1 0 b 及び / 又は第 3 アクセラレータ 3 1 0 c に特化した追加的又は代替的な CBB を含んでよい。例えば、もし第 2 アクセラレータ 3 1 0 b が GPU を実装していれば、第 2 アクセラレータ 3 1 0 b の演算に特化した CBB は、スレッドディスパッチャ、グラフィックテクノロジーインターフェース及び / 又はコンピュータグラフィック及び / 又は画像処理を処理する処理速度及び全体性能を向上するのに好ましい任意のその他の CBB を含み得る。さらに、もし第 3 アクセラレータ 3 1 0 c が FPGAs を実装していれば、第 3 アクセラレータ 3 1 0 c の演算に特化した CBB は、1 又は複数の算術ロジックユニット (ALU) 及び / 又は汎用の計算を処理する処理速度及び全体性能を向上するのに好ましい任意のその他の CBB を含み得る。

40

## 【 0 0 4 1 】

図 3 の異種システム 3 0 4 は、ホストプロセッサ 3 0 6、第 1 アクセラレータ 3 1 0 a、第 2 アクセラレータ 3 1 0 b 及び第 3 アクセラレータ 3 1 0 c を含むが、いくつかの例

50

において、異種システム 304 は、特定用途向け命令セットプロセッサ (ASIP)、物理演算ユニット (PPU)、指定された DSP、画像プロセッサ、コプロセッサ、浮動小数点ユニット、ネットワークプロセッサ、マルチコア及びフロントエンドプロセッサを含む、任意の数の処理要素 (例えばホストプロセッサ及び / 又はアクセラレータ) を含んでよい。

#### 【0042】

さらに、図 3 の例において、畳み込みエンジン 312、RNN エンジン 314、メモリ 316、MMU 318、DSP 320、コントローラ 322、DMA ユニット 324、第 1 スケジューラ 326、第 2 スケジューラ 328、第 3 スケジューラ 330、第 4 スケジューラ 332、第 5 スケジューラ 334、第 1 カーネルライブラリ 336 及び第 2 カーネルライブラリ 338 は第 1 アクセラレータ 310 a 上に実装されるが、畳み込みエンジン 312、RNN エンジン 314、メモリ 316、MMU 318、DSP 320、コントローラ 322、DMA ユニット 324、第 1 スケジューラ 326、第 2 スケジューラ 328、第 3 スケジューラ 330、第 4 スケジューラ 332、第 5 スケジューラ 334、第 1 カーネルライブラリ 336 及び第 2 カーネルライブラリ 338 の 1 又は複数は、ホストプロセッサ 306、第 2 アクセラレータ 310 b 及び / 又は第 3 アクセラレータ 310 c に実装され得る。

#### 【0043】

図 4 は、例示的な 1 又は複数のスケジューラを含む例示的なコンピューティングシステム 400 を示すブロック図である。いくつかの例において、コンピューティングシステム 400 は、図 3 のコンピューティングシステム 300 に対応し得る。図 4 の例において、コンピューティングシステム 400 は、例示的な入力 402、例示的なコンパイラ 404 及び例示的なアクセラレータ 406 を含む。いくつかの例において、アクセラレータ 406 は、図 3 の第 1 アクセラレータ 310 a に対応し得る。図 4 において、入力 402 はコンパイラ 404 に結合される。入力 402 は、アクセラレータ 406 で実行されるべきワークロードである。いくつかの例において、コンパイラ 404 は、図 3 のホストプロセッサ 306 及び / 又は外部デバイスに対応し得る。

#### 【0044】

図 4 の例において、入力 402 は、例えば、機能、アルゴリズム、プログラム、アプリケーション、及び / 又は、アクセラレータ 406 によって実行される他のコードである。いくつかの例において、入力 402 は、機能、アルゴリズム、プログラム、アプリケーション及び / 又は他のコードのグラフ記述であってよい。追加的又は代替的な例において、入力 402 は深層学習及び / 又はコンピュータビジョンのような AI 処理に関するワークロードである。

#### 【0045】

図 4 に図示の例において、コンパイラ 404 は入力 402 及びアクセラレータ 406 に結合される。コンパイラ 404 は入力 402 を受信し、入力 402 をアクセラレータ 406 によって実行される 1 又は複数の実行ファイル内へコンパイルする。例えば、コンパイラ 404 は、入力 402 を受信し、ワークロード (例えば入力 402) の様々なワークロードノードをアクセラレータ 406 の様々な CBB に割り当てるグラフコンパイラである。更に、コンパイラ 404 は、アクセラレータ 406 のメモリ内の 1 又は複数のバッファに対してメモリを割り振る。

#### 【0046】

図 4 の例において、アクセラレータ 406 はコンパイラ 404 に結合されており、例示的なクレジットマネージャー 408、例示的な CnC ファブリック 410、例示的なデータファブリック 411、例示的な畳み込みエンジン 412、例示的な DMA ユニット 414、例示的な RNN エンジン 416、例示的な DSP 418、例示的なメモリ 420 及び例示的な MMU 422 を含む。更に、例示的な畳み込みエンジン 412、例示的な DMA ユニット 414、例示的な RNN エンジン 416 及び例示的な DSP 418 の各々は、例示的な第 1 スケジューラ 424、例示的な第 2 スケジューラ 426、例示的な第 3 スケジ

10

20

30

40

50

ユーラ 4 2 8 及び例示的な第 4 スケジューラ 4 3 0 をそれぞれ含む。さらに、例示的な D S P 4 1 8 は例示的なカーネルライブラリ 4 3 2 を含む。いくつかの例において、第 1 スケジューラ 4 2 4 は図 3 の第 1 スケジューラ 3 2 6 に対応し得る。追加的又は代替的な例において、第 2 スケジューラ 4 2 6 は図 3 の第 3 スケジューラ 3 3 0 に対応し得る。さらなる例において、第 3 スケジューラ 4 2 8 は図 3 の第 2 スケジューラ 3 2 8 に対応し得る。いくつかの例において、第 4 スケジューラ 4 3 0 は図 4 の第 4 スケジューラ 3 3 2 に対応し得る。

#### 【 0 0 4 7 】

図 4 に図示の例において、クレジットマネージャー 4 0 8 はコンパイラ 4 0 4 及び C n C ファブリック 4 1 0 に結合される。クレジットマネージャー 4 0 8 は、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び / 又は D S P 4 1 8 の 1 又は複数に関連付けられたクレジットを管理するデバイスである。いくつかの例において、クレジットマネージャー 4 0 8 は、クレジットマネージャーコントローラとしてコントローラにより実装され得る。クレジットは、メモリ 4 2 0 内で利用可能なワークロードノードに関連付けられたデータ、及び / 又は、ワークロードノードの出力に対してメモリ 4 2 0 内で利用可能なスペースの量を表す。例えば、クレジットマネージャー 4 0 8 は、コンパイラ 4 0 4 から受信した 1 又は複数の実行ファイルに基づいて、メモリ 4 2 0 を、所与のワークロードのワークロードノードごとに関連付けられた 1 又は複数のバッファに区分し得る。もしワークロードノードがバッファにデータを書き込むよう構成されていれば、ワークロードノードはプロデューサーであり、もしワークロードノードがバッファからデータを読み出すよう構成されていれば、ワークロードノードはコンシューマーである。

#### 【 0 0 4 8 】

図 4 の例において、クレジットマネージャー 4 0 8 は更に、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び / 又は D S P 4 1 8 の 1 又は複数にクレジットを送信及び / 又はクレジット受信するよう構成される。いくつかの例において、クレジットマネージャー 4 0 8 は、アクセラレータ 4 0 6 の制御ユニットとして実装される。例えば、クレジットマネージャー 4 0 8 はアクセラレータ 4 0 6 の演算を管理し得る。さらに、クレジットマネージャー 4 0 8 は、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び / 又は D S P 4 1 8 の 1 又は複数に、実行ファイル及び / 又はコンパイラ 4 0 4 から受信した他の機械可読命令に対してどのように応答するかを命令し得る。

#### 【 0 0 4 9 】

図 4 の例において、C n C ファブリック 4 1 0 は、クレジットマネージャー 4 0 8、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び D S P 4 1 8 に結合される。C n C ファブリック 4 1 0 は、クレジットマネージャー 4 0 8、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び / 又は D S P 4 1 8 の 1 又は複数が、クレジットマネージャー 4 0 8、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び / 又は D S P 4 1 8 の 1 又は複数にクレジットを送信及び / 又はそれらからクレジットを受信することを可能にする少なくとも 1 つのロジック回路と電氣的に相互接続されるネットワークである。いくつかの例において、C n C ファブリック 4 1 0 は、図 3 の第 2 通信バス 3 4 0 に対応し得る。

#### 【 0 0 5 0 】

図 4 の例において、データファブリック 4 1 1 は、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6、D S P 4 1 8、メモリ 4 2 0 及び M M U 4 2 2 に結合される。データファブリック 4 1 1 は、クレジットマネージャー 4 0 8、畳み込みエンジン 4 1 2、R N N エンジン 4 1 6、D S P 4 1 8、メモリ 4 2 0 及び / 又は M M U 4 2 2 の 1 又は複数が、クレジットマネージャー 4 0 8、畳み込みエンジン 4 1 2、R N N エンジン 4 1 6、D S P 4 1 8、メモリ 4 2 0 及び / 又は M M U 4 2 2 の 1 又は複数にデータを送信及び / 又はそれらからデータを受信することを可能にする少なくとも 1 つのロジック回路と電氣的に相互接続するネットワークである。いくつかの例において、データ

ファブリック 4 1 1 は図 3 の第 2 通信バス 3 4 0 に対応し得る。

【 0 0 5 1 】

図 4 に図示の例において、畳み込みエンジン 4 1 2 は C n C ファブリック 4 1 0 及びデータファブリック 4 1 1 に結合される。畳み込みエンジン 4 1 2 は、畳み込みに関連するタスクの処理を向上するよう構成されたデバイスである。さらに、畳み込みエンジン 4 1 2 は、視覚イメージの解析に関連付けられたタスク及び / 又は C N N に関連付けられた他のタスクの処理を向上させる。いくつかの例において、畳み込みエンジン 4 1 2 は図 3 の畳み込みエンジン 3 1 2 に対応し得る。

【 0 0 5 2 】

図 4 に図示の例において、DMA ユニット 4 1 4 は C n C ファブリック 4 1 0 及びデータファブリック 4 1 1 に結合される。DMA ユニット 4 1 4 は、畳み込みエンジン 4 1 2、R N N エンジン 4 1 6 又は D S P 4 1 8 のうちの少なくとも 1 つが、対応するプロセッサ (例えばホストプロセッサ 3 0 6) から独立して、アクセラレータ 4 0 6 から離れたメモリ (例えばシステムメモリ 3 0 2) にアクセスすることを可能にするデバイスである。いくつかの例において、DMA ユニット 4 1 4 は図 3 の DMA ユニット 3 2 4 に対応し得る。例えば、DMA ユニット 4 1 4 は、アナログ又はデジタル回路、ロジック回路、プログラマブルプロセッサ、プログラマブルコントローラ、GPU、DSP、ASIC、PLD 及び / 又は FPLD の 1 又は複数によって実装され得る。

【 0 0 5 3 】

図 4 において、R N N エンジン 4 1 6 は C n C ファブリック 4 1 0 及びデータファブリック 4 1 1 に結合される。R N N エンジン 4 1 6 は、R N N に関連するタスクの処理を向上するよう構成されたデバイスである。更に、R N N エンジン 4 1 6 は、セグメント化されていない繋がった手書き認識、音声認識の解析に関連付けられたタスク及び / 又は R N N に関連付けられた他のタスクの処理を向上させる。いくつかの例において、R N N エンジン 4 1 6 は図 3 の R N N エンジン 3 1 4 に対応し得る。

【 0 0 5 4 】

図 4 の例において、D S P 4 1 8 は C n C ファブリック 4 1 0 及びデータファブリック 4 1 1 に結合される。D S P 4 1 8 はデジタル信号の処理を向上させるデバイスである。例えば、D S P 4 1 8 は、カメラ及び / 又はコンピュータビジョンに関する他のセンサからのデータのような、連続的な実世界の信号を測定、フィルタ及び / 又は圧縮する処理を促進する。いくつかの例において、D S P 4 1 8 は図 3 の D S P 3 2 0 に対応し得る。

【 0 0 5 5 】

図 4 の例において、メモリ 4 2 0 はデータファブリック 4 1 1 に結合される。メモリ 4 2 0 は、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び D S P 4 1 8 のうちの少なくとも 1 つの間での共有ストレージである。いくつかの例において、メモリ 4 2 0 は図 3 のメモリ 3 1 6 に対応し得る。メモリ 4 2 0 は、クレジットマネージャ 4 0 8 から受信した実行ファイルに関連付けられたワークロードの 1 又は複数のワークロードノードに関連付けられた 1 又は複数のバッファに区分化され得る。図 4 の例において、メモリ 4 2 0 はアクセラレータ 4 0 6 に位置する物理ストレージである。しかしながら、他の例において、メモリ 4 2 0 はアクセラレータ 4 0 6 の外部にあってよく及び / 又はそうでなければ離れていてよい。さらなる例において、メモリ 4 2 0 は仮想記憶装置であってよい。図 4 の例において、メモリ 4 2 0 は永続ストレージ (例えば ROM、PROM、EPROM、EEPROM 等) である。他の例において、メモリ 4 2 0 は永続 BIOS 又はフラッシュストレージであってよい。さらなる例において、メモリ 4 2 0 は揮発性メモリであってよい。

【 0 0 5 6 】

図 4 に図示の例において、例示的な MMU 4 2 2 はデータファブリック 4 1 1 に結合される。MMU 4 2 2 は、メモリ 4 2 0 及び / 又はアクセラレータ 4 0 6 から離れたメモリのアドレスへの参照を含むデバイスである。MMU 4 2 2 は更に、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び / 又は D S P 4 1 8 の 1 又は複数

10

20

30

40

50

に利用される仮想的メモリアドレスを、メモリ 4 2 0 及び / 又はアクセラレータ 4 0 6 から離れたメモリ内の物理アドレスに変換する。いくつかの例において、MMU 4 2 2 は図 3 の MMU 3 1 8 に対応し得る。

#### 【 0 0 5 7 】

図 4 の例において、第 1 スケジューラ 4 2 4、第 2 スケジューラ 4 2 6、第 3 スケジューラ 4 2 8 及び第 4 スケジューラ 4 3 0 の各々は、クレジットマネージャー 4 0 8 及び / 又はアクセラレータ 4 0 6 の追加的な C B B によって、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び D S P 4 1 8 にそれぞれ割り当てられているワークロードの一部（例えばワークロードノード）を、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6 及び D S P 4 1 8 がそれぞれいつ実行するかを決定するデバイスである。タスク及び / 又は所与のワークロードノードの他の演算に応じて、ワークロードノードはプロデューサー又はコンシューマーであり得る。プロデューサーワークロードノードは他のワークロードノードで利用されるデータを生成し、他方、コンシューマーワークロードノードは他のワークロードノードで生成されたデータを消費及び / 又はそうでなければ処理する。

10

#### 【 0 0 5 8 】

図 4 に図示の例において、カーネルライブラリ 4 3 2 は 1 又は複数のカーネルを含むデータ構造である。いくつかの例において、カーネルライブラリ 4 3 2 は図 3 の第 1 カーネルライブラリ 3 3 6 に対応し得る。カーネルライブラリ 4 3 2 のカーネルは、例えば、D S P 4 1 8 で高スループットとなるためにコンパイルされたルーチンである。カーネルは、例えば、アクセラレータ 4 0 6 上で動作する実行ファイルの実行可能なサブセクションに対応する。図 4 の例において、アクセラレータ 4 0 6 は V P U を実装し、クレジットマネージャー 4 0 8、C n C ファブリック 4 1 0、データファブリック 4 1 1、畳み込みエンジン 4 1 2、DMA ユニット 4 1 4、R N N エンジン 4 1 6、D S P 4 1 8、メモリ 4 2 0 及び MMU 4 2 2 を含むが、アクセラレータ 4 0 6 は図 4 に図示されたこれらに追加的又は代替的な C B B を含んでよい。

20

#### 【 0 0 5 9 】

図 4 の例で、演算において、第 1 スケジューラ 4 2 4 は、畳み込みエンジン 4 1 2 に割り当てられたワークロードノードに対するワークロードノードへの入力バッファ及びワークロードノードからの出力バッファに対応したクレジットを読み込む。例えば、入力バッファはワークロードノードがそこからデータを読み出すよう構成されたバッファである一方、出力バッファはワークロードノードがそこからデータを書き込むよう構成されたバッファである。いくつかの例において、第 1 ワークロードノードの入力バッファは第 2 ワークロードノードの出力バッファであり得る。さらに、第 1 スケジューラ 4 2 4 はクレジットマネージャー 4 0 8 からクレジットを受信及び / 又はそうでなければ取得する。

30

#### 【 0 0 6 0 】

図 4 の例で、演算において、第 1 スケジューラ 4 2 4 は、畳み込みエンジン 4 1 2 に割り当てられたワークロードノードを選択し、選択されたワークロードノードへの入力バッファに格納されているデータを演算するために、第 1 スケジューラ 4 2 4 がクレジットの閾値量を受信しているか否かを決定する。例えば、第 1 スケジューラ 4 2 4 は、入力バッファに対するプロデューサーワークロードノードから受信したクレジット数を、入力バッファに対するクレジットの閾値数と比較する。もし第 1 スケジューラ 4 2 4 はクレジットの閾値量を受信していなければ、第 1 スケジューラ 4 2 4 は畳み込みエンジン 4 1 2 に割り当てられた他のワークロードノードの処理を繰り返す。

40

#### 【 0 0 6 1 】

図 4 に図示した例において、演算において、もし第 1 スケジューラ 4 2 4 が、選択されたワークロードノードへの入力バッファに格納されているデータを演算するためにクレジットの閾値量を受信していれば、第 1 スケジューラ 4 2 4 は、第 1 スケジューラ 4 2 4 が選択されたワークロードノードに対して出力バッファにデータを書き込むためにクレジットの閾値量を受信しているか否かを判断する。例えば、第 1 スケジューラ 4 2 4 は、出力

50

バッファに対するコンシューマーワークロードノードから受信したクレジット数を、選択されたワークロードノードのための出力バッファに対するクレジットの閾値数と比較する。もし第1スケジューラ424がクレジットの閾値量を受信していなければ、第1スケジューラ424は畳み込みエンジン412に割り当てられた他のワークロードノードの処理を繰り返す。もし第1スケジューラ424が出力バッファにデータを書き込むためにクレジットの閾値量を受信していれば、第1スケジューラ424は選択されたワークロードノードの実行が準備できたことを示す。次に、第1スケジューラ424は畳み込みエンジン412に割り当てられた追加的なワークロードノードに対してこの処理を繰り返す。

#### 【0062】

図4の例で、演算において、畳み込みエンジン412に割り当てられたワークロードノードが解析された後に、第1スケジューラ424は実行の準備ができたワークロードノードをスケジューリングする。第1スケジューラ424は次に、スケジュールに従ってワークロードノードをディスパッチする。ディスパッチされたワークロードノードが畳み込みエンジン412によって実行された後に、第1スケジューラ424は、入力バッファ及び/又は出力バッファに対応するクレジットをクレジットマネージャー408に送信する。第1スケジューラ424は実行されるスケジュール内に追加的なワークロードノードがあるかどうかを判断する。もしスケジュール内に追加的なワークロードノードがあるなら、第1スケジューラ424は、畳み込みエンジン412で実行されるスケジュール内の次のワークロードノードを生じさせる。

#### 【0063】

図5は図3及び4の1又は複数のスケジューラを実装し得る例示的なスケジューラ500のブロック図である。例えば、スケジューラ500は、図3の第1スケジューラ326、第2スケジューラ328、第3スケジューラ330、第4スケジューラ332及び/又は第5スケジューラ334、及び/又は、図4の第1スケジューラ424、第2スケジューラ426、第3スケジューラ428及び/又は第4スケジューラ430、及び/又は、図6のスケジューラ600、及び/又は、図7の第1スケジューラ722、第2スケジューラ724、第3スケジューラ726及び/又は第4スケジューラ728の例示的な実装である。

#### 【0064】

図5の例において、スケジューラ500は、例示的なワークロードインターフェース502、例示的なバッファクレジット格納装置504、例示的なクレジット比較器506、例示的なワークロードノードディスパッチャ508及び例示的な通信バス510を含む。スケジューラ500は、スケジューラ500が関連付けられるところCBBに割り当てられているワークロードの一部(例えばワークロードノード)を、スケジューラ500が関連付けられるところのCBBがいつ実行するかを決定するデバイスである。

#### 【0065】

図5に図示の例において、ワークロードインターフェース502は、スケジューラ500、バッファクレジット記録装置504、クレジット比較器506及び/又はワークロードノードディスパッチャ508の外部の他のデバイスと通信するよう構成されたデバイスである。

例えば、ワークロードインターフェース502は、スケジューラ500が関連付けられるところのCBBによって実行されるワークロードノードを受信及び/又はそうでなければ取得し得る。更に又は代替的に、ワークロードインターフェース502は他のスケジューラ、他のCBB及び/又は他のデバイスにクレジットを送信及び/又はそれらから受信し得る。さらに、ワークロードインターフェース502は、ワークロードノードへの入力バッファ及び/又はワークロードノードからの出力バッファに対応するクレジットを、バッファクレジット記録装置504内へ及び/又はそこから読み込み得る。

#### 【0066】

いくつかの例において、例示的なワークロードインターフェース502は例示的なインターフェースする手段を実装する。インターフェース手段は、図8の少なくともブロック

10

20

30

40

50

802、818及び822によって実装されるような実行可能命令によって実装される。例えば、図8のブロック802、818及び822の実行可能命令は、図9の例に示される例示的なプロセッサ910及び/又は例示的なアクセラレータ912のような少なくとも1つのプロセッサで実行されてよい。他の例において、インターフェース手段は、ハードウェアロジック、ハードウェア実装ステートマシン、論理回路、及び/又は、ハードウェア、ソフトウェア及び/又はファームウェアの他の任意の組み合わせによって実装される。

#### 【0067】

図5に図示される例において、バッファクレジット記録装置504は、ワークロードインターフェース502、クレジット比較器506及び/又はワークロードノードディスパッチャ508のうちの少なくとも1つの間での共有ストレージである。バッファクレジット記録装置504はスケジューラ500に位置する物理ストレージである。しかしながら、他の例において、バッファクレジット記録装置504はスケジューラ500の外部にあってよく及び/又はそうでなければそれから離れていてよい。さらなる例において、バッファクレジット記録装置504は仮想記憶装置であってよい。図5の例において、バッファクレジット記録装置504は永続ストレージ（例えばROM、PROM、EPROM、EEPROM等）である。他の例において、バッファクレジット記録装置504は永続BIOS又はフラッシュストレージであってよい。さらなる例において、バッファクレジット記録装置504は揮発性メモリであってよい。

#### 【0068】

図5の例において、バッファクレジット記録装置504は、スケジューラ500が関連付けられるところのCBBに割り当てられたワークロードノードに関連付けられたワークロードノードへの入力バッファ及び/又はワークロードノードからの出力バッファに対応したクレジットを格納することに関連付けられたメモリである。例えば、バッファクレジット記録装置504は、スケジューラ500が関連付けられるところのCBBに割り当てられたワークロードノードごとに対するフィールド、及び、スケジューラ500が関連付けられるところのCBBに割り当てられたワークロードノードに関連付けられたワークロードノードへの各入力バッファ及び/又はワークロードノードからの各出力バッファに対するフィールド、を含むデータ構造として実装され得る。

#### 【0069】

図5の図示の例において、バッファクレジット記録装置504は更に又は代替的に、スケジューラ500が関連付けられるところのCBBに割り当てられているワークロードノード、及び/又は、ワークロードノードへの入力バッファ及び/又はワークロードノードからの出力バッファに対応するクレジットの閾値量を格納し得る。さらに、バッファクレジット記録装置504は、各ワークロードノードへの入力バッファ及び/又は各ワークロードノードからの出力バッファに対するクレジットの閾値数に関連付けられたフィールドを含む。

#### 【0070】

図5の例において、ワークロードノードがプロデューサーである（例えばワークロードノードが他のワークロードノードによって利用されるデータを生成する）場合に、クレジットの閾値数は、スケジューラ500が関連付けられるところのCBBがプロデューサーワークロードノードを実行し得る前に満たされるべき出力バッファのスペースの閾値量（例えばメモリ420の区分化されたスペース）に対応する。更に、ワークロードノードがコンシューマーである（例えばワークロードノードが他のワークロードノードによって生成されたデータを処理する）場合に、クレジットの閾値数は、スケジューラ500が関連付けられるところのCBBがコンシューマーワークロードノードを実行し得る前に満たされるべき入力バッファのデータの閾値量（例えばメモリ420の区分化されたスペース）に対応する。

#### 【0071】

いくつかの例において、例示的なバッファクレジット格納装置504は例示的な格納す

10

20

30

40

50

る手段を実装する。格納手段は、図 8 において実装されたもののような実行可能命令によって実装され得る。例えば、実行可能命令は図 9 の例に示された例示的なプロセッサ 9 1 0 及び / 又は例示的なアクセラレータ 9 1 2 のような少なくとも 1 つのプロセッサで実行されてよい。他の例において、格納手段は、ハードウェアロジック、ハードウェア実装ステートマシン、論理回路、及び / 又は、ハードウェア、ソフトウェア及び / 又はファームウェアの他の任意の組み合わせによって実装される。

#### 【 0 0 7 2 】

図 5 に図示された例において、クレジット比較器 5 0 6 は、スケジューラ 5 0 0 が関連付けられるところの C B B に割り当てられたワークロードノードへの入力バッファ及び / 又はワークロードノードからの出力バッファに対応したクレジットの閾値数が受信されているか否かを判断するよう構成されたデバイスである。クレジット比較器 5 0 6 は、スケジューラ 5 0 0 が関連付けられるところの C B B に割り当てられたワークロードノードを選択するよう構成される。

10

#### 【 0 0 7 3 】

図 5 の例において、クレジット比較器 5 0 6 は更に、選択されたワークロードノードに対して入力バッファに格納されたデータを演算するために、スケジューラ 5 0 0 がクレジットの閾値量を受信したか否かを判断するよう構成される。例えば、クレジット比較器 5 0 6 は、外部デバイス（例えばクレジットマネージャー 4 0 8、コントローラ 3 2 2 等）から受信したクレジット数に関連付けられたバッファクレジット記録装置 5 0 4 内のフィールドを、選択されたワークロードノードへの入力バッファに対するクレジットの閾値数に関連付けられたバッファクレジット記録装置 5 0 4 内のフィールドと比較する。もしスケジューラ 5 0 0 がクレジットの閾値量を受信していないならば、クレジット比較器 5 0 6 はスケジューラ 5 0 0 が関連付けられるところの C B B に割り当てられた他のワークロードノードの処理を繰り返す。

20

#### 【 0 0 7 4 】

図 5 に図示した例において、もしスケジューラ 5 0 0 が入力バッファに格納されたデータを演算するためにクレジットの閾値量を受信しているならば、クレジット比較器 5 0 6 は選択されたワークロードノードに対する出力バッファにデータを書き込むためにクレジットの閾値量をスケジューラ 5 0 0 が受信したか否かを判断する。例えば、クレジット比較器 5 0 6 は、選択されたワークロードノードに対する出力バッファに対する外部デバイス（例えばクレジットマネージャー 4 0 8、コントローラ 3 2 2 等）から受信したクレジット数に関連付けられたバッファクレジット記録装置 5 0 4 内のフィールドを、出力バッファに対するクレジットの閾値数に関連付けられたバッファクレジット記録装置 5 0 4 内のフィールドと比較する。

30

#### 【 0 0 7 5 】

図 5 の例において、もしスケジューラ 5 0 0 がクレジットの閾値量を受信していないならば、クレジット比較器 5 0 6 はスケジューラ 5 0 0 に関連付けられるところの C B B に割り当てられた他のワークロードノードの処理を繰り返す。もしスケジューラ 5 0 0 が出力バッファにデータを書き込むためにクレジットの閾値量を受信しているならば、クレジット比較器 5 0 6 は選択されたワークロードノードが実行する準備があることを示す。次に、クレジット比較器 5 0 6 は、スケジューラ 5 0 0 が関連付けられるところの C B B に割り当てられた追加的なワークロードノードに対するこの処理を繰り返す。

40

#### 【 0 0 7 6 】

いくつかの例において、例示的なクレジット比較器 5 0 6 は例示的な比較する手段を実装する。比較手段は、少なくとも図 8 のブロック 8 0 4、8 0 6、8 0 8、8 1 0 及び 8 1 2 によって実装されるような実行可能命令によって実装される。例えば、図 8 のブロック 8 0 4、8 0 6、8 0 8、8 1 0 及び 8 1 2 の実行可能命令は、図 9 の例に示された例示的なプロセッサ 9 1 0 及び / 又は例示的なアクセラレータ 9 1 2 のような少なくとも 1 つのプロセッサで実行されてよい。他の例において、比較手段は、ハードウェアロジック、ハードウェア実装ステートマシン、論理回路、及び / 又は、ハードウェア、ソフトウェ

50

ア及び／又はファームウェアの他の任意の組み合わせによって実装される。

【 0 0 7 7 】

図 5 の例において、ワークロードノードディスパッチャ 5 0 8 は、スケジューラ 5 0 0 に関連付けられるところの C B B で実行されるべくスケジューラ 5 0 0 に関連付けられるところの C B B に割り当てられた 1 又は複数のワークロードノードをスケジューリングするデバイスである。例えば、スケジューラ 5 0 0 に関連付けられるところの C B B に割り当てられたワークロードノードが解析された後に、ワークロードノードディスパッチャ 5 0 8 は実行する準備ができたワークロードノードをスケジューリングする。例えば、ワークロードノードディスパッチャ 5 0 8 は、実行する準備ができたワークロードノードをラウンドロビンスケジュールのようなスケジューリングアルゴリズムに基づいてスケジューリングする。ワークロードノードディスパッチャ 5 0 8 は次に、スケジュールに従ってワークロードノードをディスパッチする。他の例において、ワークロードノードディスパッチャ 5 0 8 は、実行する準備ができたワークロードノードをスケジューリングする任意の他の適切な任意のアルゴリズムを利用し得る。

10

【 0 0 7 8 】

図 5 に図示された例において、ディスパッチされたワークロードノードがスケジューラ 5 0 0 が関連付けられるところの C B B によって実行されるにつれ、ワークロードインターフェース 5 0 2 は、ワークロードインターフェース 5 0 2 がクレジットをそこから受信するところの外部デバイス（例えばクレジットマネージャ 4 0 8、コントローラ 3 2 2 等）への入力バッファに関連付けられたクレジットを送信する。ワークロードノードディスパッチャ 5 0 8 は更に、実行されるスケジュール内に追加的なワークロードノードがあるかどうかを判断する。もしスケジュール内に追加的なワークロードノードがあるならば、ワークロードノードディスパッチャ 5 0 8 はスケジュール内の次のワークロードノードをディスパッチする。

20

【 0 0 7 9 】

いくつかの例において、例示的なワークロードノードディスパッチャ 5 0 8 は、例示的なディスパッチする手段を実装する。ディスパッチ手段は、少なくとも図 8 のブロック 8 1 4、8 1 6 及び 8 2 0 によって実装されるような実行可能命令によって実装される。例えば、図 8 のブロック 8 1 4、8 1 6 及び 8 2 0 の実行可能命令は、図 9 の例に示される例示的なプロセッサ 9 1 0 及び／又は例示的なアクセラレータ 9 1 2 のような少なくとも 1 つのプロセッサで実行されてよい。他の例において、ディスパッチ手段は、ハードウェアロジック、ハードウェア実装ステートマシン、論理回路、及び／又は、ハードウェア、ソフトウェア及び／又はファームウェアの他の任意の組み合わせによって実装される。

30

【 0 0 8 0 】

本明細書で開示された例において、ワークロードインターフェース 5 0 2、バッファクレジット記録装置 5 0 4、クレジット比較器 5 0 6 及びワークロードノードディスパッチャ 5 0 8 の各々は、スケジューラ 5 0 0 の他の要素と通信する。例えば、ワークロードインターフェース 5 0 2、バッファクレジット記録装置 5 0 4、クレジット比較器 5 0 6 及びワークロードノードディスパッチャ 5 0 8 は例示的な通信バス 5 1 0 で通信する。本明細書で開示されるいくつかの例において、ワークロードインターフェース 5 0 2、バッファクレジット記録装置 5 0 4、クレジット比較器 5 0 6 及びワークロードノードディスパッチャ 5 0 8 は、任意の適切な有線及び／又は無線通信システムで通信してよい。更に、本明細書で開示されたいくつかの例において、ワークロードインターフェース 5 0 2、バッファクレジット記録装置 5 0 4、クレジット比較器 5 0 6 及びワークロードノードディスパッチャ 5 0 8 の各々は、任意の適切な有線及び／又は無線通信システムでスケジューラ 5 0 0 の外部の任意のコンポーネントと通信してよい。

40

【 0 0 8 1 】

図 6 は、図 5 のバッファクレジット記録装置 5 0 4 のさらなる詳細を示す例示的なスケジューラ 6 0 0 のブロック図である。スケジューラ 6 0 0 は、図 3 の第 1 スケジューラ 3 2 6、第 2 スケジューラ 3 2 8、第 3 スケジューラ 3 3 0、第 4 スケジューラ 3 3 2 及び

50

／又は第5スケジューラ334、及び／又は、図4の第1スケジューラ424、第2スケジューラ426、第3スケジューラ428及び／又は第4スケジューラ430、及び／又は、図5のスケジューラ500、及び／又は、図7の第1スケジューラ722、第2スケジューラ724、第3スケジューラ726及び／又は第4スケジューラ728の例示的な実装である。

#### 【0082】

図6の例において、スケジューラ600は、例示的なワークロードインターフェース502、例示的なバッファクレジット格納装置504、例示的なクレジット比較器506及び例示的なワークロードノードディスパッチャ508を含む。スケジューラ600は、スケジューラ600が関連付けられるところのCBBが、スケジューラ600が関連付けられるところのCBBに割り当てられているワークロードの一部（例えばワークロードノード）をいつ実行するかを決定するデバイスである。

10

#### 【0083】

図6に図示の例において、ワークロードインターフェース502は、スケジューラ600の外部の1又は複数のデバイス、バッファクレジット記録装置504及びワークロードノードディスパッチャ508に結合される。ワークロードインターフェース502は、スケジューラ600の外部の他のデバイス、バッファクレジット記録装置504及び／又はワークロードノードディスパッチャ508と通信するように構成されているデバイスである。例えば、ワークロードインターフェース502は、スケジューラ600が関連付けられるところのCBBによって実行されるワークロードノードを受信及び／又はそうでなければ取得し得る。更に又は代替的に、ワークロードインターフェース502は、クレジットをスケジューラ600の外部の1又は複数のデバイスへ送信及び／又はそこから受信し得る。さらに、ワークロードインターフェース502は、ワークロードノードへの入力バッファ及び／又はワークロードノードからの出力バッファに対応するクレジットを、バッファクレジット記録装置504内へ及び／又はから読み込み得る。

20

#### 【0084】

図6に図示した例において、バッファクレジット記録装置504は、ワークロードインターフェース502、クレジット比較器506及び／又はワークロードノードディスパッチャ508のうちの少なくとも1つの間での共有ストレージである。バッファクレジット記録装置504はスケジューラ500に位置する物理ストレージである。しかしながら、他の例において、バッファクレジット記録装置504はスケジューラ500の外部にあってよく及び／又はそうでなければそれから離れていてもよい。さらなる例において、バッファクレジット記録装置504は仮想記憶装置であってよい。図5の例において、バッファクレジット記録装置504は、永続ストレージ（例えばROM、PROM、EPROM、EEPROM等）である。他の例において、バッファクレジット記録装置504は永続BIOS又はフラッシュストレージであってよい。さらなる例において、バッファクレジット記録装置504は揮発性メモリであってよい。

30

#### 【0085】

図6の例において、バッファクレジット記録装置504は、第1ワークロードノードWN[0]、第2ワークロードノードWN[1]及び第nワークロードノードWN[n]に対応する行を含むデータ構造である。バッファクレジット記録装置504は更に、第1コンシューマー（例えばコンシューマー[0]）に対する入力バッファ、第1コンシューマー（例えばコンシューマー[1]）に対する入力バッファ、第1プロデューサー（例えばプロデューサー[0]）に対する出力バッファ、及び、第mプロデューサー（例えばプロデューサー[m]）に対する出力バッファに対応する列を含む。バッファクレジット記録装置504はさらに、各ワークロードノードへの入力バッファ及び／又は各ワークロードノードからの出力バッファのクレジットの閾値数に対応する列を含む。

40

#### 【0086】

図6に図示の例において、第1ワークロードノードWN[0]、第2ワークロードノードWN[1]及び第nワークロードノードWN[n]の各々は、スケジューラ600が関

50

連付けられるところの C B B に割り当てられる。バッファクレジット記録装置 5 0 4 において、第 1 ワークロードノード W N [ 0 ]、第 2 ワークロードノード W N [ 1 ] 及び第 n ワークロードノード W N [ n ] に対応する行と、第 1 コンシューマー（例えばコンシューマー [ 0 ]）に対する入力バッファ、第 1 コンシューマー（例えばコンシューマー [ 1 ]）に対する入力バッファ、第 1 プロデューサー（例えばプロデューサー [ 0 ]）に対する出力バッファ、第 m プロデューサー（例えばプロデューサー [ m ]）に対する出力バッファに対応する列との間の交差は、そのバッファに対する 1 又は複数の外部デバイスから受信するクレジット数に対応するフィールドを示す。さらに、各ワークロードノードへの入力バッファ及び / 又は各ワークロードノードからの出力バッファに対するクレジットの閾値数に対応する列は、スケジューラ 6 0 0 が関連付けられるところの C B B がそれぞれのワークロードノードを演算し得る前にバッファに対して満たされるべきクレジットの閾値数を示す。

10

#### 【 0 0 8 7 】

図 6 の例において、バッファクレジット記録装置 5 0 4 における、第 1 ワークロードノード W N [ 0 ]、第 2 ワークロードノード W N [ 1 ] 及び第 n ワークロードノード W N [ n ] に対応する行と、第 1 コンシューマー（例えばコンシューマー [ 0 ]）に対する入力バッファ及び第 1 コンシューマー（例えばコンシューマー [ 1 ]）に対する入力バッファに対応する列との間の交差のフィールドは、外部デバイス（例えばクレジットマネージャー 4 0 8、コントローラ 3 2 2 等）によってゼロの値に初期化される。更に、バッファクレジット記録装置 5 0 4 における、第 1 ワークロードノード W N [ 0 ]、第 2 ワークロードノード W N [ 1 ] 及び第 n ワークロードノード W N [ n ] に対応する行と、第 1 プロデューサー（例えばプロデューサー [ 0 ]）に対する出力バッファ及び第 m プロデューサー（例えばプロデューサー [ m ]）に対する出力バッファに対応する列との間の交差のフィールドは、外部デバイス（例えばクレジットマネージャー 4 0 8、コントローラ 3 2 2 等）によって関連するバッファ内に区分化されたメモリの量に対応した値に初期化される。さらに、入力バッファ及び / 又は出力バッファに対するクレジットの閾値数に対応する列は、外部デバイス（例えばクレジットマネージャー 4 0 8、コントローラ 3 2 2、ホストプロセッサ 3 0 6 で実行するソフトウェア等）によって初期化される。

20

#### 【 0 0 8 8 】

図 6 に図示の例において、クレジット比較器 5 0 6 は、バッファクレジット記録装置 5 0 4 及びワークロードノードディスパッチャ 5 0 8 に結合される。クレジット比較器 5 0 6 は、スケジューラ 6 0 0 が関連付けられるところの C B B に割り当てられたワークロードノードへの入力バッファ及び / 又はワークロードノードからの出力バッファに対応するクレジットの閾値数を受信しているか否かを判断するよう構成されたデバイスである。図 6 の例において、ワークロードノードディスパッチャ 5 0 8 は、ワークロードインターフェース 5 0 2、バッファクレジット記録装置 5 0 4、クレジット比較器 5 0 6 及びスケジューラ 6 0 0 の外部の 1 又は複数のデバイスに結合される。ワークロードノードディスパッチャ 5 0 8 は、例えば、スケジューラ 6 0 0 が関連付けられるところの C B B で実行されるべく、スケジューラ 6 0 0 が関連付けられるところの C B B に割り当てられた 1 又は複数のワークロードノードをスケジューリングするデバイスである。

30

40

#### 【 0 0 8 9 】

図 6 の例で、演算において、ワークロードインターフェース 5 0 2 が外部デバイス（例えばクレジットマネージャー 4 0 8、コントローラ 3 2 2 等）からワークロードノードを受信及び / 又はそうでなければ取得した場合に、ワークロードインターフェース 5 0 2 はワークロードノードをワークロードノードに対応するバッファクレジット記録装置 5 0 4 のそれぞれのフィールド内へ読み込む。さらに、クレジット比較器 5 0 6 は、スケジューラ 6 0 0 が関連付けられるところの C B B に割り当てられたワークロードノードを選択する。

#### 【 0 0 9 0 】

図 6 に図示の例において、クレジット比較器 5 0 6 は、選択されたワークロードノード

50

に対する入力バッファに格納されたデータを演算するためにクレジットの閾値量をスケジューラ 600 が受信しているか否かを判断する。例えば、クレジット比較器 506 は、外部デバイス（例えばクレジットマネージャ 408、コントローラ 322 等）から受信したクレジット数に関連付けられたバッファクレジット記録装置 504 内のフィールドを、選択されたワークロードノードへの入力バッファに対するクレジットの閾値数に関連付けられたバッファクレジット記録装置 504 内のフィールドと比較する。クレジットの閾値数は、スケジューラ 600 が関連付けられるところの CBB がコンシューマワークロードノードを実行し得る前に満たされるべき入力バッファのデータの閾値量（例えばメモリ 420 の区分化されたスペース）に対応する。もしスケジューラ 600 がクレジットの閾値量を受信していないならば、クレジット比較器 506 は、スケジューラ 600 が関連付けられるところの CBB に割り当てられた他のワークロードノードの処理を繰り返す。

10

#### 【0091】

図 6 に図示した例において、もしスケジューラ 600 が入力バッファに格納されたデータを演算するためにクレジットの閾値量を受信しているならば、クレジット比較器 506 は、スケジューラ 600 が選択されたワークロードノードに対する出力バッファへデータを書き込むためにクレジットの閾値量を受信しているか否かを判断する。例えば、クレジット比較器 506 は、選択されたワークロードノードに対する出力バッファに対する外部デバイス（例えばクレジットマネージャ 408、コントローラ 322 等）から受信したクレジット数に関連付けられたバッファクレジット記録装置 504 内のフィールドを、出力バッファに対するクレジットの閾値数に関連付けられたバッファクレジット記録装置 504 内のフィールドと比較する。クレジットの閾値数は、スケジューラ 600 が関連付けられるところの CBB がプロデューサワークロードノードを実行し得る前に満たされるべき出力バッファのスペースの閾値量（例えばメモリの区分化されたスペース）に対応し得る。

20

#### 【0092】

図 6 の例において、もしスケジューラ 600 がクレジットの閾値量を受信していないならば、クレジット比較器 506 はスケジューラ 600 が関連付けられるところの CBB に割り当てられた他のワークロードノードの処理を繰り返す。もしスケジューラ 600 が出力バッファへデータを書き込むためにクレジットの閾値量を受信しているならば、クレジット比較器 506 は、選択されたワークロードノードが実行する準備ができていることを示す。次に、クレジット比較器 506 は、スケジューラ 600 が関連付けられるところの CBB に割り当てられた追加的なワークロードノードに対するこの処理を繰り返す。

30

#### 【0093】

図 6 の例において、ワークロードノードディスパッチャ 508 は、スケジューラ 600 が関連付けられるところの CBB で実行されるべく、スケジューラ 600 が関連付けられるところの CBB に割り当てられた 1 又は複数のワークロードノードをスケジューリングするデバイスである。例えば、スケジューラ 600 が関連付けられるところの CBB に割り当てられたワークロードノードが解析された後に、ワークロードノードディスパッチャ 508 は、実行の準備ができたワークロードノードをスケジューリングする。例えば、ワークロードノードディスパッチャ 508 は、実行の準備ができたワークロードノードを、ラウンドロビンスケジュールのようなスケジューリングアルゴリズムに基づいてスケジューリングする。ワークロードノードディスパッチャ 508 は次に、スケジュールに従ってワークロードノードをディスパッチする。他の例において、ワークロードノードディスパッチャ 508 は実行の準備ができたワークロードノードをスケジューリングする任意の他の適切な任意のアルゴリズムを利用し得る。

40

#### 【0094】

図 6 に図示した例において、ディスパッチされたワークロードノードがスケジューラ 600 が関連付けられるところの CBB によって実行されるにつれ、ワークロードインターフェース 502 は、入力バッファに関連付けられたクレジットをワークロードインターフェース 502 がそこからクレジットを受信したところの外部デバイス（例えばクレジット

50

マネージャー 408、コントローラ 322 等)へ送信する。ワークロードノードディスパッチャ 508 は更に、実行されるべきスケジュール内に追加的なワークロードノードがあるかどうかを判断する。もしスケジュール内に追加的なワークロードノードがあるならば、ワークロードノードディスパッチャ 508 はスケジュール内の次のワークロードノードをディスパッチする。

#### 【0095】

図 7 は、パイプライン及びバッファを実装している異種システムのアクセラレータで実行するワークロードを表す例示的なグラフ 700 の図である。例えば、アクセラレータは、図 3 の第 1 アクセラレータ 310 a であり、異種システムは異種システム 304 である。例示的なグラフ 700 は、例示的な第 1 ワークロードノード 702 (WN[0])、例示的な第 2 ワークロードノード 704 (WN[1])、例示的な第 3 ワークロードノード 706 (WN[2])、例示的な第 4 ワークロードノード 708 (WN[3]) 及び例示的な第 5 ワークロードノード 710 (WN[4]) を含む。図 7 の例において、アクセラレータは、ワークロードノードを様々な CBB に割り当てる例示的なクレジットマネージャー 712 からのスケジュールに基づいたグラフ 700 によって表されるワークロードを実行するよう構成されている。例えば、クレジットマネージャー 712 及び/又は他のコントローラは、第 1 ワークロードノード 702 (WN[0]) を例示的な第 1 CBB 714 に、第 2 ワークロードノード 704 (WN[1]) を例示的な第 2 CBB 716 に、第 3 ワークロードノード 706 (WN[2]) を例示的な第 3 CBB 718 に、第 4 ワークロードノード 708 (WN[3]) を例示的な第 4 CBB 720 に、及び、第 5 ワークロードノード 710 (WN[4]) を例示的な第 2 CBB 716 に割り当てる。

#### 【0096】

図 7 の例において、例示的な第 1 CBB 714、例示的な第 2 CBB 716、例示的な第 3 CBB 718 及び例示的な第 4 CBB 720 の各々は、例示的な第 1 スケジューラ 722、例示的な第 2 スケジューラ 724、例示的な第 3 スケジューラ 726 及び例示的な第 4 スケジューラ 728 を含む。第 1 スケジューラ 722、第 2 スケジューラ 724、第 3 スケジューラ 726 及び第 4 スケジューラ 728 の各々は、図 5 のスケジューラ 500 及び/又は図 6 のスケジューラ 600 によって実装され得る。

#### 【0097】

図 7 に図示の例において、第 1 ワークロードノード 702 (WN[0]) 及び第 2 ワークロードノード 704 (WN[1]) は例示的な第 1 バッファ 730 に関連付けられる。第 1 バッファ 730 は、第 1 ワークロードノード 702 (WN[0]) の出力バッファ及び第 2 ワークロードノード 704 (WN[1]) の入力バッファである。第 2 ワークロードノード 704 (WN[1]) 及び第 3 ワークロードノード 706 (WN[2]) は例示的な第 2 バッファ 732 に関連付けられる。第 2 バッファ 732 は第 2 ワークロードノード 704 (WN[1]) の出力バッファ及び第 3 ワークロードノード 706 (WN[2]) の入力バッファである。第 4 ワークロードノード 708 (WN[3]) 及び第 5 ワークロードノード 710 (WN[4]) は例示的な第 3 バッファ 734 に関連付けられる。第 3 バッファ 734 は第 4 ワークロードノード 708 (WN[3]) の出力バッファ及び第 5 ワークロードノード 710 (WN[4]) の入力バッファである。第 1 バッファ 730、第 2 バッファ 732 及び第 3 バッファ 734 の各々は循環バッファによって実装され得る。図 7 の例において、第 1 バッファ 730、第 2 バッファ 732 及び第 3 バッファ 734 の各々は、アクセラレータのメモリの 5 つの区分を含み、それらの各々はデータのタイルを格納し得る。

#### 【0098】

図 7 に図示した例において、第 1 ワークロードノード 702 (WN[0]) はプロデューサーワークロードノードであり、クレジットマネージャー 712 は第 1 バッファ 730 に対する 5 つのクレジットで第 1 スケジューラ 722 を初期化する。同様に、第 2 ワークロードノード 704 (WN[1]) はプロデューサーワークロードノードなので、クレジットマネージャー 712 は第 2 バッファ 732 に対する 5 つのクレジットで第 2 スケジューラ 724 を初期化する。

ーラ 724 を初期化する。更に、第 4 ワークロード ノード 708 (WN[3]) はプロデューサー ワークロード ノードであり、クレジット マネージャー 712 は第 3 バッファ 734 に対する 5 つのクレジットで第 4 スケジューラ 728 を初期化する。

【 0 0 9 9 】

第1スケジューラ722、第2スケジューラ724及び第4スケジューラ728の各々に提供された5つのクレジットは、第1バッファ730、第2バッファ732及び第3バッファ734のサイズの表現である。更に、第2ワークロードノード704（WN[1]）はまたコンシューマーワークロードノードでもあり、クレジットマネージャー712は第1バッファ730に対するゼロクレジットで第2スケジューラ724を初期化する。さらに、第3ワークロードノード706（WN[2]）はコンシューマーワークロードノードなので、クレジットマネージャー712は第2バッファ732に対するゼロクレジットで第3スケジューラ726を初期化する。さらに、第5ワークロードノード710（WN[4]）はコンシューマーワークロードノードであり、クレジットマネージャー712は第3バッファ734に対するゼロクレジットで第2スケジューラ724を初期化する。

【 0 1 0 0 】

図 7 の例において、第 1 スケジューラ 7 2 2 は第 1 ワークロード ノード 7 0 2 ( W N [ 0 ] ) への入力バッファ及びそこから出力バッファの両方に対するクレジットの閾値数を受信しているので、第 1 スケジューラ 7 2 2 は第 1 C B B 7 1 4 で実行する第 1 ワークロード ノード 7 0 2 ( W N [ 0 ] ) をディスパッチする。更に、第 4 スケジューラ 7 2 8 は、第 4 ワークロード ノード 7 0 8 ( W N [ 3 ] ) への入力バッファ及びそこから出力バッファの両方に対するクレジットの閾値数を受信しているので、第 4 スケジューラ 7 2 8 は第 4 C B B 7 2 0 で実行する第 4 ワークロード ノード 7 0 8 ( W N [ 3 ] ) をディスパッチする。第 1 ワークロード ノード 7 0 2 ( W N [ 0 ] ) が第 1 C B B 7 1 4 で実行されるにつれ、第 1 C B B 7 1 4 はデータを第 1 バッファ 7 3 0 に送信する。同様に、第 4 ワークロード ノード 7 0 8 ( W N [ 3 ] ) が第 4 C B B 7 2 0 で実行されるにつれ、第 4 C B B 7 2 0 はデータを第 3 バッファ 7 3 4 に送信する。

【 0 1 0 1 】

図 7 に図示された例において、第 1 C B B 7 1 4 及び第 4 C B B 7 2 0 の各々がそれぞれ第 1 ワークロードノード 7 0 2 ( W N [ 0 ] ) 及び第 4 ワークロードノード 7 0 8 ( W N [ 3 ] ) に関連付けられたデータのタイルを送信するにつれ、第 1 スケジューラ 7 2 2 及び第 4 スケジューラ 7 2 8 は、それぞれ、第 1 C B B 7 1 4 及び第 4 C B B 7 2 0 から第 1 バッファ 7 3 0 及び第 3 バッファ 7 3 4 へ送信されたデータのタイルごとにクレジットマネージャー 7 1 2 へクレジットを送信する。クレジットマネージャー 7 1 2 は、第 1 スケジューラ 7 2 2 から受信したクレジットを第 2 スケジューラ 7 2 4 へ、第 4 スケジューラ 7 2 8 から受信したクレジットを第 2 スケジューラ 7 2 4 へ送信する。第 4 C B B 7 2 0 が第 4 ワークロードノード 7 0 8 ( W N [ 3 ] ) を実行するにつれ、第 4 C B B 7 2 0 は第 3 バッファ 7 3 4 に格納するデータの 2 つのタイルを生成する。同様に、第 1 C B B 7 1 4 が第 1 ワークロードノード 7 0 2 ( W N [ 0 ] ) を実行するにつれ、第 1 C B B 7 1 4 は第 1 バッファ 7 3 0 に格納するデータの 5 つのタイルを生成する。

【 0 1 0 2 】

図 7 の例において、第 1 C B B 7 1 4 が第 1 ワークロードノード 7 0 2 ( W N [ 0 ] ) を実行するよりも迅速に、第 4 C B B 7 2 0 は第 4 ワークロードノード 7 0 8 ( W N [ 3 ] ) を実行する。第 2 バッファ 7 3 2 には利用可能なメモリはあるが、第 2 ワークロードノード 7 0 4 ( W N [ 1 ] ) が依存しているデータが準備される前に第 5 ワークロードノード 7 1 0 ( W N [ 4 ] ) が依存しているデータが準備されるので、第 2 スケジューラ 7 2 4 は、第 2 ワークロードノード 7 0 4 ( W N [ 1 ] ) とは対照的に、第 5 ワークロードノード 7 1 0 ( W N [ 4 ] ) を第 2 C B B 7 1 6 で実行するものとして選択する。

【 0 1 0 3 】

図 7 に図示の例において、第 5 ワークロードノード 7 1 0 (WN[4]) が第 2 CBB 7 1 6 で実行され、第 2 CBB 7 1 6 が第 3 バッファ 7 3 4 に格納されたデータのタイル

を消費するにつれ、第2スケジューラ724は、第3バッファ734に関連付けられたクレジットを、第3バッファ734からの第2CBB716で消費されたデータのタイルごとに、クレジットマネージャー712へ返送する。次に、第1バッファ730及び第2バッファ732に対するクレジットの閾値量が満たされると、第2スケジューラ724は第2CBB716で実行する第2ワークロードノード704(WN[1])をディスパッチする。第2CBB716が第2ワークロードノード704(WN[1])に関連付けられたデータのタイルを生成し、第2バッファ732にデータを出力するにつれ、第2スケジューラ724は、第2CBB716から第2バッファ732へ送信されたデータのタイルごとに、第2バッファ732に関連付けられたクレジットをクレジットマネージャー712に送信する。

10

#### 【0104】

図7の例において、第2スケジューラ724から第2バッファ732に関連付けられたクレジットを受信すると、クレジットマネージャー712は、第2バッファ732に関連付けられたクレジットを第3スケジューラ726に送信する。第3スケジューラ726が第2バッファ732に関連付けられたクレジットの閾値量を受信した場合に、第3スケジューラ726は、第3CBB718で実行する第3ワークロードノード706(WN[2])をディスパッチする。第3CBB718が第3ワークロードノード706(WN[2])を実行し、第3CBB718が第2バッファ732に格納されたデータのタイルを消費するにつれ、第3スケジューラ726は、第2バッファ732に関連付けられたクレジットを、第2バッファ732からの第3CBB718で消費されたデータのタイルごとに、クレジットマネージャー712へ返送する。

20

#### 【0105】

追加的又は代替的な例において、第1CBB714は図4の畳み込みエンジン412に対応し得、第1スケジューラ722は図4の第1スケジューラ424に対応し得る。いくつかの例において、第2CBB716は図4のRNNエンジン416に対応し得、第2スケジューラ724は図4の第3スケジューラ428に対応し得る。さらなる例において、第3CBB718は図4のDMAユニット414に対応し得、第3スケジューラ726は図4の第2スケジューラ426に対応し得る。いくつかの例において、第4CBB720は図4のDSP418に対応し得、第4スケジューラ728は図4の第4スケジューラ430に対応し得る。

30

#### 【0106】

図3の第1スケジューラ326、第2スケジューラ328、第3スケジューラ330、第4スケジューラ332及び/又は第5スケジューラ334、及び/又は、図4の第1スケジューラ424、第2スケジューラ426、第3スケジューラ428及び/又は第4スケジューラ430、及び/又は、図7の第1スケジューラ722、第2スケジューラ724、第3スケジューラ726及び/又は第4スケジューラ728の実装の例示的な態様は図5及び/又は図6に図示されているが、図5及び/又は図6に図示された要素、処理及び/又はデバイスの1又は複数は、組み合わせたり、分割されたり、再構成されたり、省略されたり、除去されたり、及び/又は、その他の方式によって実装されてよい。さらに、図5の例示的なワークロードインターフェース502、例示的なバッファクレジット格納装置504、例示的なクレジット比較器506、例示的なワークロードノードディスパッチャ508、例示的な通信バス510及び/又は、より一般的に、例示的なスケジューラ500及び/又は図6の例示的なスケジューラ600は、ハードウェア、ソフトウェア、ファームウェア、及び/又は、ハードウェア、ソフトウェア及び/又はファームウェアの任意の組み合わせによって実装されてよい。従って、例えば、図5の例示的なワークロードインターフェース502、例示的なバッファクレジット格納装置504、例示的なクレジット比較器506、例示的なワークロードノードディスパッチャ508、例示的な通信バス510及び/又は、より一般的に、例示的なスケジューラ500及び/又は図6の例示的なスケジューラ600のいずれも、アナログ又はデジタル回路、ロジック回路、プログラマブルプロセッサ、プログラマブルコントローラ、グラフィック処理ユニット(G

40

50

P U )、デジタルシグナルプロセッサ ( D S P )、特定用途向け集積回路 ( A S I C )、プログラマブル論理デバイス ( P L D ) 及び / 又はフィールドプログラマブル論理デバイス ( F P L D ) の 1 又は複数によって実装され得る。本特許の任意の装置又はシステムの請求項を純粋なソフトウェア及び / 又はファームウェア実装を包含するように読む場合に、図 5 の例示的なワークロードインターフェース 5 0 2、例示的なバッファクレジット格納装置 5 0 4、例示的なクレジット比較器 5 0 6、例示的なワークロードノードディスパッチャ 5 0 8、例示的な通信バス 5 1 0、及び / 又は、より一般的に、例示的なスケジューラ 5 0 0 及び / 又は図 6 の例示的なスケジューラ 6 0 0 のうちの少なくとも 1 つは、ソフトウェア及び / 又はファームウェアを含む、メモリ、デジタルバーサタイルディスク ( D V D )、コンパクトディスク ( C D )、ブルーレイディスク等のような、非一時的コンピュータ可読記憶装置デバイス、又は、ストレージディスクを含むように本明細書では明示的に定義される。更にまた、図 5 の例示的なスケジューラ 5 0 0 及び / 又は図 6 の例示的なスケジューラ 6 0 0 は、図 5 及び / 又は図 6 に図示されている、要素、処理及び / 又はデバイスに追加し、又は、代わりに、1 又は複数のそれらを含んでよく、及び / 又は、図示された要素、処理及びデバイスのいずれか一つより多く又は全てを含んでよい。本明細書で用いられるように、「通信」及びそれらの変形を含む語句は、直接的な通信及び / 又は 1 又は複数の中間媒介コンポーネントを介した間接的な通信を包含し、直接物理的な (例えば有線) 通信及び / 又は一定の通信を必要としておらず、むしろ、周期的な間隔、スケジュールされた間隔、非周期的な間隔及び / 又は一回的なイベントでの選択的な通信を更に含む。

#### 【 0 1 0 7 】

図 5 のスケジューラ 5 0 0 及び / 又は図 6 のスケジューラ 6 0 0 を実装するための、例示的なハードウェアロジック、機械可読命令、ハードウェア実装ステートマシン及び / 又はそれらの任意の組み合わせを表すフローチャートが図 8 に示される。機械可読命令は、図 9 に関して以下で説明される例示的なプロセッサプラットフォーム 9 0 0 に示されるプロセッサ 9 1 0 及び / 又はアクセラレータ 9 1 2 のようなコンピュータプロセッサによる実行のための実行可能プログラム又は実行可能プログラムの一部の 1 又は複数であってよい。プログラムは、プロセッサ 9 1 0 及び / 又はアクセラレータ 9 1 2 に関連付けられた、C D - R O M、フロッピーディスク、ハードドライブ、D V D、ブルーレイディスクのような非一時的コンピュータ可読記憶媒体又はメモリに格納されたソフトウェア内に具現されてよいが、プログラム全体及び / 又はそれらの一部は代替的にプロセッサ 9 1 0 及び / 又はアクセラレータ 9 1 2 以外のデバイスで実行され得、及び / 又は、ファームウェア又は専用のハードウェアに具現され得る。さらに、例示的なプログラムが図 8 に図示されたフローチャートに言及して記載されるが、図 5 の例示的なスケジューラ 5 0 0 及び / 又は図 6 スケジューラ 6 0 0 を実装する多くの他の方法が代替的に用いられてよい。例えば、ブロックの実行の順序が変更されてよく、及び / 又は、記載されたいくつかのブロックは変更、除去又は組み合わせされてよい。更に又は代替的に、ブロックのいずれか又は全部は、ソフトウェア又はファームウェアを実行することなく対応する演算を実行するよう構築された 1 又は複数のハードウェア回路 (例えば別個及び / 又は統合されたアナログ及び / 又はデジタル回路、F P G A、A S I C、比較器、オペレーショナルアンプ (オペアンプ)、ロジック回路等) によって実装されてよい。

#### 【 0 1 0 8 】

本明細書に記載された機械可読命令は、圧縮されたフォーマット、暗号化されたフォーマット、細分化されたフォーマット、コンパイルされたフォーマット、実行可能なフォーマット、パッケージされたフォーマット等の 1 又は複数で格納されてよい。本明細書で記載される機械可読命令は、機械実行可能命令を創造、製造及び / 又は生成するのに利用されてよいデータ (例えば命令の一部、コード、コードの表現等) として格納されてよい。例えば、機械可読命令は細分化されて、ストレージデバイス及び / 又はコンピューティングデバイス (例えばサーバ) の 1 又は複数に格納されてよい。機械可読命令は、コンピューティングデバイス及び / 又は、他の機械でそれらを直接的に可読、変換可能及び / 又は

実行可能にするために、インストール、変更、適合、更新、組み合わせ、補完、構成、復号化、解凍、アンパッキング、分配、再割り当て、コンパイル等の1又は複数が必要であってよい。例えば、機械可読命令は、複数部分に格納されてよく、それらは、個別に圧縮され、暗号化されかつ別個のコンピューティングデバイスに格納され、それらの部分は、復号化され、解凍され及び組み合わせされた場合に本明細書で記載されるようなプログラムに実装される実行可能命令のセットを形成する。

#### 【0109】

別の例において、機械可読命令は、コンピュータによってそれらが読み出しえるが、特定のコンピューティングデバイス又は他のデバイスで命令を実行するために、ライブラリ（例えばダイナミックリンクライブラリ（DLL））、ソフトウェア開発キット（SDK）  
10、アプリケーションプログラミングインターフェース（API）等の追加を必要としてよい状態で格納されてよい。別の例において、機械可読命令は、機械可読命令及び/又は対応するプログラムが全部又は一部実行され得る前に構成されること（例えば、格納された設定、データ入力、記録されたネットワークアドレス等）が必要でもよい。従って、開示された機械可読命令及び/又は対応するプログラムは、格納され又はそうでなければ残り若しくは送信中の場合の、機械可読命令及び/又はプログラムの特定のフォーマット又は状態に関わらず、そのような機械可読命令及び/又はプログラムを包含するよう意図される。

#### 【0110】

本明細書で記載される機械可読命令は、過去、存在又は未来命令言語、スクリプト言語、プログラミング言語等のいずれかで表され得る。例えば、機械可読命令は以下の言語の任意のものを使用して表され得る、すなわち、C、C++、Java（登録商標）、C#、  
20、PERL、PYTHON、JavaScript（登録商標）、ハイパーテキストマークアップ言語（HTML）、構造化照会言語（SQL）、Swift等。

#### 【0111】

上で言及したように、図8の例示的な処理は、ハードディスクドライブ、フラッシュメモリ、読み出し専用メモリ、コンパクトディスク、デジタルバーサタイルディスク、キャッシュ、ランダムアクセスメモリ、及び/又は、任意の期間（例えば拡張された時間の間、恒久的に、短い期間の間、一時的なバッファリングの間及び/又は情報のキャッシングの間）情報が格納されるその他のストレージデバイス及び/又はストレージディスクのよ  
30うな、非一時的コンピュータ及び/又は機械可読媒体に格納された実行可能命令（例えばコンピュータ及び/又は機械可読命令）を使用することを実装されてよい。本明細書で用いられるように、非一時的コンピュータ可読媒体という用語は、コンピュータ可読記憶装置デバイス及び/又はストレージディスクの任意のタイプを含むように、かつ、伝播する信号を除外しかつ送信媒体を除外するよう明示的に定義される。

#### 【0112】

「含み」及び「備え」（及びそれらの全ての型及び時制）は本明細書において非限定的な用語に用いられている。従って、請求項がプリアンブルとして又は請求項内の任意の種類の記述において「含む」又は「備える」の任意の型（例えば、備える、含む、備え、含み、有し等）が用いられるときにはいつでも、対応する請求項又は記述の範囲の外部に入ることなく、追加的な要素、用語等が存在してよいと理解されるべきである。本明細書で用い  
40られている、「少なくとも」という語句は、例えば、請求項のプリアンブルの遷移用語として用いられる場合、それは「備え」及び「含み」が非限定的な用語であるのと同じ態様において非限定的である。「及び/又は」という用語は例えば、A、B及び/又はCのような型で用いられる場合、（1）A単独、（2）B単独、（3）C単独、（4）AとB、（5）AとC、（6）BとC及び（7）AとBとCのような、A、B、Cの任意の組み合わせ又はサブセットを指す。本明細書において、構造、コンポーネント、項目、オブジェクト及び/又は物を説明する文脈の中で用いられると、「A及びBのうちの少なくとも1つ」という語句は、（1）少なくとも1つのA、（2）少なくとも1つのB及び（3）少なくとも1つA及び少なくとも1つのBのいずれかを  
50含む実装を指すことを意図している

。同様に、本明細書において、構造、コンポーネント、項目、オブジェクト及び／又は物を説明する文脈の中で用いられると、「A又はBのうちの少なくとも1つ」という語句は、(1)少なくとも1つのA、(2)少なくとも1つのB及び(3)少なくとも1つA及び少なくとも1つのBのいずれかを含む実装を指すことを意図している。本明細書において、処理、命令、動作、活動及び／又は段階の遂行又は実行を説明する文脈の中で用いられると、「A及びBのうちの少なくとも1つ」という語句は、(1)少なくとも1つのA、(2)少なくとも1つのB及び(3)少なくとも1つA及び少なくとも1つのBのいずれかを含む実装を指すことを意図している。同様に、本明細書において、処理、命令、動作、活動及び／又は段階の遂行又は実行を説明する文脈の中で用いられると、「A又はBのうちの少なくとも1つ」という語句は、(1)少なくとも1つのA、(2)少なくとも1つのB及び(3)少なくとも1つA及び少なくとも1つのBのいずれかを含む実装を指すことを意図している。

10

**【0113】**

本明細書で用いられると、単数の参照(例えば「a」、「an」、「第1」、「第2」等)は複数を除外していない。「a」又は「an」エンティティという用語は、本明細書で用いられると、そのエンティティの1又は複数を指す。「a」(又は「an」)、「1又は複数」、「少なくとも1つの」という用語は本明細書では同じ意味で用いられ得る。さらに、個別に列挙されているが、複数の手段、要素又は方法動作は、例えば単一のユニット又はプロセッサによって実装されてよい。更に、別個の特徴が異なる例又はクレイムに含まれていてよいが、これらはおそらく組み合わせられてよく、異なる例又はクレイムに含まれていることは特徴の組み合わせが実現可能でない及び／又は有利でないことを暗示するものではない。

20

**【0114】**

図8は、図5のスケジューラ500及び／又は図6のスケジューラ600を実装するために実行されてよい機械可読命令によって実装され得る処理800を表すフローチャートである。処理800はブロック802で開始し、そこにおいて、ワークロードインターフェース502は、スケジューラ500及び／又はスケジューラ600が関連付けられるところのCBBに割り当てられたワークロードノードへの入力バッファ及び／又はそこからの出力バッファに対応するクレジットを、バッファクレジット記録装置504内へ読み込む。

30

**【0115】**

図8に図示される例において、処理800はブロック804で継続し、そこにおいて、クレジット比較器506はスケジューラ500及び／又はスケジューラ600が関連付けられるところのCBBに割り当てられたワークロードノードを選択する。ブロック806において、クレジット比較器506は、選択されたワークロードノードに対する入力バッファに格納されたデータを演算するためにスケジューラ500及び／又はスケジューラ600がクレジットの閾値量を受信しているか否かを判断する。例えば、クレジット比較器506は、外部デバイス(例えばクレジットマネージャ408、コントローラ322等)から受信したクレジット数に関連付けられた配列又は他のデータ構造内のフィールドを、選択されたワークロードノードへの入力バッファに対するクレジットの閾値数に関連付けられた配列又は他のデータ構造内のフィールドと比較する。もしスケジューラ500及び／又はスケジューラ600が選択されたワークロードノードに対する入力バッファに格納されているデータを演算するためにクレジットの閾値量を受信していない、とクレジット比較器506が判断したら(ブロック806:NO)、処理800はブロック812に進む。

40

**【0116】**

図8の例において、もしスケジューラ500及び／又はスケジューラ600が入力バッファに格納されているデータを演算するためにクレジットの閾値量を受信している、とクレジット比較器506が判断したら(ブロック806:YES)、処理800はブロック808に進む。ブロック808において、クレジット比較器506は、選択されたワーク

50

ロードノードに対する出力バッファにデータを書き込むためにスケジューラ 500 及び / 又はスケジューラ 600 がクレジットの閾値量を受信しているか否かを判断する。例えば、クレジット比較器 506 は、選択されたワークロードノードに対する出力バッファに対する外部デバイス（例えばクレジットマネージャー 408、コントローラ 322 等）から受信したクレジット数に関連付けられた配列又は他のデータ構造内のフィールドを、出力バッファに対するクレジットの閾値数に関連付けられた配列又は他のデータ構造内のフィールドと比較する。もしスケジューラ 500 及び / 又はスケジューラ 600 がクレジットの閾値量を受信していない、とクレジット比較器 506 が判断したら（ブロック 808：NO）、処理 800 はブロック 812 に進む。もし出力バッファにデータを書き込むためにスケジューラ 500 及び / 又はスケジューラ 600 がクレジットの閾値量を受信している、とクレジット比較器 506 が判断したら（ブロック 808：YES）、クレジット比較器 506 は、選択されたワークロードノードの実行の準備ができたことをブロック 810 において示す。

10

**【0117】**

図 8 に図示した例において、ブロック 812 で、クレジット比較器 506 は処理すべき追加的なワークロードノードがあるか否かを判断する。もしクレジット比較器 506 が処理すべき追加的なワークロードノードがあると判断したなら（ブロック 812：YES）、クレジット比較器 506 は追加的なワークロードノードを選択して、処理 800 はブロック 806 に進む。もしクレジット比較器 506 が処理すべき追加的なワークロードノードがないと判断したなら（ブロック 812：NO）、処理 800 はブロック 814 に進む。

20

**【0118】**

図 8 に図示の例において、ブロック 814 で、ワークロードノードディスパッチャ 508 は、実行の準備ができたワークロードノードをスケジューリングする。ブロック 816 において、ワークロードノードディスパッチャ 508 はスケジュールに従ってワークロードノードをディスパッチする。ブロック 818 において、ディスパッチされたワークロードノードがスケジューラ 500 及び / 又はスケジューラ 600 が関連付けられるところの CBB によって実行されるにつれ、ワークロードインターフェース 502 は、入力バッファに関連付けられたクレジットを、そこからワークロードインターフェース 502 がクレジットを受信したところの外部デバイス（例えばクレジットマネージャー 408、コントローラ 322 等）へ送信する。

30

**【0119】**

図 8 に図示した例において、ブロック 820 で、ワークロードノードディスパッチャ 508 は実行すべきスケジュール内に追加的なワークロードノードがあるかどうかを判断する。もしワークロードノードディスパッチャ 508 がスケジュール内に追加的なワークロードノードがあると判断したなら（ブロック 820：YES）、処理 800 はブロック 816 に進む。もしワークロードノードディスパッチャ 508 がスケジュールに実行すべき追加的なワークロードノードがないと判断したなら（ブロック 820：NO）、処理 800 はブロック 822 に進む。

**【0120】**

図 8 の例において、ブロック 822 で、ワークロードインターフェース 502 は演算を継続するか否かを判断する。例えば、ワークロードインターフェース 502 で演算を継続するとの判断が生じるであろう条件は、追加的なワークロードノードを受信することを含む。もしワークロードインターフェース 502 が演算を継続することを決定したら（ブロック 822：YES）、処理 800 はブロック 802 に進む。もしワークロードインターフェース 502 が演算を継続しないことを決定したら（ブロック 822：NO）、処理 800 は終了する。

40

**【0121】**

図 9 は、図 5 のスケジューラ 500 及び / 又は図 6 のスケジューラ 600 の 1 又は複数のインスタス化を実装するために図 8 の命令を実行するよう構築された例示的なプロセスプラットフォーム 900 のブロック図である。プロセスプラットフォーム 900 は

50

、例えば、サーバ、パーソナルコンピュータ、ワークステーション、自己学習機械（例えばニューラルネットワーク）、モバイルデバイス（例えばセルフォン、スマートフォン、i P a d（登録商標）のようなタブレット）、パーソナルデジタルアシスタント（P D A）、インターネット機器、D V Dプレイヤー、C Dプレイヤー、デジタルビデオレコーダ、ブルーレイプレイヤー、ゲームコンソール、パーソナルビデオレコーダ、セットトップボックス、ヘッドセット若しくは他のウェアラブルデバイス、又は、任意のその他のタイプのコンピューティングデバイスであり得る。

#### 【0122】

図示の例のプロセッサプラットフォーム900は、プロセッサ910及びアクセラレータ912を含む。図示の例のプロセッサ910はハードウェアである。例えば、プロセッサ910は、任意の所望なファミリー若しくは製造者からの集積回路、ロジック回路、マイクロプロセッサ、GPU、DSP又はコントローラの1又は複数によって実装され得る。ハードウェアプロセッサは、半導体ベース（例えばシリコンベース）デバイスであってよい。更に、アクセラレータ912は、例えば、集積回路、ロジック回路、マイクロプロセッサ、GPU、DSP、FPGA、VPU、コントローラ、及び/又は、任意の所望なファミリー若しくは製造者からの他のCBBの1又は複数により実装され得る。図示の例のアクセラレータ912はハードウェアである。ハードウェアアクセラレータは、半導体ベース（例えばシリコンベース）デバイスであってよい。この例において、アクセラレータ912は、例示的な畳み込みエンジン312、例示的なRNNエンジン314、例示的なメモリ316、例示的なMMU318、例示的なDSP320、例示的なコントローラ322及び例示的なDMAユニット324を実装する。さらに、例示的な畳み込みエンジン312、例示的なRNNエンジン314、例示的なDMAユニット324、例示的なDSP320及び例示的なコントローラ322の各々は、例示的な第1スケジューラ326、例示的な第2スケジューラ328、例示的な第3スケジューラ330、例示的な第4スケジューラ332及び例示的な第5スケジューラ334をそれぞれ含む。図9の例において、例示的な第1スケジューラ326、例示的な第2スケジューラ328、例示的な第3スケジューラ330、例示的な第4スケジューラ332及び例示的な第5スケジューラ334の各々は、例示的なワークロードインターフェース502、例示的なバッファクレジット格納装置504、例示的なクレジット比較器506、例示的なワークロードノードディスパッチャ508、及び/又は、より一般的に、スケジューラ500を含む。

#### 【0123】

追加的又は代替的な例において、プロセッサ910は、例示的な畳み込みエンジン312、例示的なRNNエンジン314、例示的なメモリ316、例示的なMMU318、例示的なDSP320、例示的なコントローラ322及び例示的なDMAユニット324を実装する。さらにこのような追加的又は代替的な例において、例示的な畳み込みエンジン312、例示的なRNNエンジン314、例示的なDMAユニット324、例示的なDSP320及び例示的なコントローラ322の各々は、例示的な第1スケジューラ326、例示的な第2スケジューラ328、例示的な第3スケジューラ330、例示的な第4スケジューラ332及び例示的な第5スケジューラ334をそれぞれ含む。このような追加的又は代替的な例において、例示的な第1スケジューラ326、例示的な第2スケジューラ328、例示的な第3スケジューラ330、例示的な第4スケジューラ332及び例示的な第5スケジューラ334の各々は、例示的なワークロードインターフェース502、例示的なバッファクレジット格納装置504、例示的なクレジット比較器506、例示的なワークロードノードディスパッチャ508、及び/又は、より一般的に、スケジューラ500を含む。

#### 【0124】

図示の例のプロセッサ910はローカルメモリ911（例えばキャッシュ）を含む。図示の例のプロセッサ910は、バス918で揮発性メモリ914及び不揮発性メモリ916を含むメインメモリと通信する。さらには図示の例のアクセラレータ912は、ローカルメモリ913（例えばキャッシュ）を含む。図示の例のアクセラレータ912は、バス

918で揮発性メモリ914及び不揮発性メモリ916を含むメインメモリと通信する。揮発性メモリ914は、同期ダイナミックランダムアクセスメモリ（SDRAM）、ダイナミックランダムアクセスメモリ（DRAM）、RAMBUS（登録商標）ダイナミックランダムアクセスメモリ（RDRAM（登録商標））及び／又は任意のその他のタイプのアクセスメモリデバイスによって実装されてよい。不揮発性メモリ916は、フラッシュメモリ及び／又はその他の所望の任意のタイプのメモリデバイスによって実装されてよい。メインメモリ914、916へのアクセスはメモリコントローラによって制御される。

#### 【0125】

図示の例のプロセッサプラットフォーム900はまたインターフェース回路920をも含む。インターフェース回路920は、Ethernet（登録商標）インターフェース、ユニバーサルシリアルバス（USB）、Bluetooth（登録商標）インターフェース、近距離無線通信（NFC）インターフェース、及び／又は、PCIエクスプレスインターフェースのような任意のタイプのインターフェース規格によって実装されてよい。

#### 【0126】

図示の例において、1又は複数入力デバイス922はインターフェース回路920に接続される。入力デバイス922は、ユーザにデータ及び／又はコマンドをプロセッサ910及び／又はアクセラレータ912内へ入力させるのを可能にする。入力デバイスは、例えば、オーディオセンサ、マイク、カメラ（静止画又は動画）、キーボード、ボタン、マウス、タッチスクリーン、トラックパッド、トラックボール、isopoint及び／又は音声認識システムによって実装され得る。

#### 【0127】

1又は複数出力デバイス924もまた、図示の例のインターフェース回路920に接続される。出力デバイス924は、例えば、ディスプレイデバイス（例えば発光ダイオード（LED）、有機発光ダイオード（OLED）、液晶ディスプレイ（LCD）、カソードレイ管ディスプレイ（CRT）、面内スイッチング（IPS）ディスプレイ、タッチスクリーン等）、触覚出力デバイス、プリンタ及び／又はスピーカによって実装され得る。図示の例のインターフェース回路920は従って、典型的には、グラフィックドライバカード、グラフィックドライバチップ及び／又はグラフィックドライバプロセッサを含む。

#### 【0128】

図示の例のインターフェース回路920はまた、送信機、受信機、トランシーバ、モデム、住宅ゲートウェイ、無線アクセスポイント、及び／又は、ネットワーク926を介した外部の機械（例えば任意の種類のコンピューティングデバイス）とのデータの交換を促進するネットワークインターフェースのような通信デバイスを含む。通信は、例えば、Ethernet（登録商標）接続、デジタル加入者ライン（DSL）接続、電話線接続、同軸ケーブルシステム、衛星システム、ラインオブサイト無線システム、セルラ電話システム等を介し得る。

#### 【0129】

図示の例のプロセッサプラットフォーム900はまた、ソフトウェア及び／又はデータを格納するための1又は複数の大容量ストレージデバイス928を含む。このような大容量ストレージデバイス928の例は、フロッピーディスクドライブ、ハードドライブディスク、コンパクトディスクドライブ、ブルーレイディスクドライブ、独立ディスクの冗長アレイ（RAID）システム、デジタルバーサタイルディスク（DVD）ドライブを含む。

#### 【0130】

図8の機械実行可能命令932は、大容量ストレージデバイス928内、揮発性メモリ914内、不揮発性メモリ916内、及び／又は、CD又はDVDのようなリムーバブル非一時的コンピュータ可読記憶媒体上に格納されてよい。

#### 【0131】

上記から、ワークロードのスタティックマッピングの順不同にパイプライン化された実行が可能に、例示的な方法、装置及び製造物が開示されていることが理解されるであろう。さらに、ワークロードノードが依存するところのデータが利用可能であり、かつ、ワー

10

20

30

40

50

クロードノードの実行によって生成された出力を格納するのに利用可能な十分なメモリがある場合に、計算ビルディングブロックがワークロードノードを実行することが可能であるように、例示的な方法、装置及び製造物が開示されている。更に、本明細書で開示された例は、スケジュール及び/又は他の順序から独立してワークロードノードが割り当てられるところの計算ビルディングブロックによってワークロードノードが実行されることを可能にする。開示された方法、装置及び製造物は、処理デバイスの利用を増加することによってコンピューティングデバイスの使用の効率性を向上する。さらに、本明細書で開示された例示的な方法、装置及び製造物は、ワークロードを処理及び/又はそうでなければ実行するために処理デバイスによって用いられる計算サイクルの数を減少させる。従って、開示された方法、装置及び製造物は、コンピュータの機能を1又は複数改善するよう方向付けられている。

10

**【0132】**

ワークロードのスタティックマッピングの順不同にパイプライン化された実行を可能にする例示的な方法、装置、システム及び製造物が本明細書に開示されている。さらなる例及びそれらの組み合わせは以下のものを含む：例1は、クレジットの第1の数をメモリ内へ読み込むインターフェースと、クレジットの第1の数をバッファ内のメモリ利用可能性に関連付けられたクレジットの閾値数と比較する比較器と、クレジットの第1の数がクレジットの閾値数を満たす場合に、1又は複数の計算ビルディングブロックの最初の一つで実行されるワークロードのワークロードノードを選択するディスパッチャとを備える装置を含む。

20

**【0133】**

例2は、例1の装置を含み、インターフェースは、インターフェースがクレジットマネージャーからクレジットの第1の数を受信した場合にクレジットの第1の数をメモリ内へ読み込み、ワークロードノードに関連付けられたデータの1又は複数のタイルが1又は複数の計算ビルディングブロックの最初の一つからバッファへ送信されるにつれ、バッファに送信された各タイルに対してクレジットをクレジットマネージャーに送信するものである。

**【0134】**

例3は例1の装置を含み、バッファはワークロードノードに関連付けられた出力バッファであり、クレジットの第1の数は出力バッファに対応しており、クレジットの閾値数は出力バッファ内のメモリの閾値量に対応する。

30

**【0135】**

例4は例1の装置を含み、バッファはワークロードノードに関連付けられた入力バッファであり、クレジットの第1の数は入力バッファに対応しており、クレジットの閾値数は入力バッファ内のデータの閾値量に対応する。

**【0136】**

例5は例1の装置を含み、バッファは第1バッファであり、クレジットの閾値数はクレジットの第1閾値数であり、比較器はクレジットの第2の数を第2バッファ内のメモリ利用可能性に関連付けられたクレジットの第2閾値数と比較するものであり、ディスパッチャは、(1)クレジットの第1の数がクレジットの第1閾値数を満たし、(2)クレジットの第2の数がクレジットの第2閾値数を満たす場合に、1又は複数の計算ビルディングブロックの最初の一つで実行されるワークロードノードを選択するものである。

40

**【0137】**

例6は例5の装置を含み、第2バッファはワークロードノードに関連付けられた入力バッファであり、クレジットの第2の数は入力バッファに対応しており、クレジットの第2閾値数は入力バッファ内のデータの閾値量に対応する。

**【0138】**

例7は例1の装置を含み、クレジットの閾値数はクレジットの第1閾値数であり、ワークロードノードは第1ワークロードノードであり、(1)クレジットの第1の数がクレジットの第1閾値数と満たし、(2)クレジットの第2の数がクレジットの第2閾値数を満

50

たす場合に、ディスパッチャは、1又は複数の計算ビルディングブロックの最初の一つで実行される第1ワークロードノード及び第2ワークロードノードをスケジューリングするものである。

【0139】

例8は命令を備える非一時的コンピュータ可読記憶媒体を含み、命令は実行された場合に、少なくとも1つのプロセッサに、クレジットの第1の数をメモリ内へ読み込むこと、クレジットの第1の数をバッファ内のメモリ利用可能性に関連付けられたクレジットの閾値数と比較すること、及び、クレジットの第1の数がクレジットの閾値数を満たす場合に、計算ビルディングブロックで実行されるワークロードのワークロードノードを選択することを少なくとも生じせしめる。

10

【0140】

例9は例8の非一時的コンピュータ可読記憶媒体を含み、命令は実行された場合に、少なくとも1つのプロセッサに、クレジットの第1の数がクレジットマネージャーから受信された場合にクレジットの第1の数をメモリ内へ読み込むこと、及び、ワークロードノードに関連付けられたデータの1又は複数のタイルが計算ビルディングブロックからバッファへ送信されるにつれ、バッファに送信された各タイルに対してクレジットをクレジットマネージャーに送信することを生じせしめる。

【0141】

例10は例8の非一時的コンピュータ可読記憶媒体を含み、バッファはワークロードノードに関連付けられた出力バッファであり、クレジットの第1の数は出力バッファに対応しており、クレジットの閾値数は出力バッファ内のメモリの閾値量に対応する。

20

【0142】

例11は例8の非一時的コンピュータ可読記憶媒体を含み、バッファはワークロードノードに関連付けられた入力バッファであり、クレジットの第1の数は入力バッファに対応しており、クレジットの閾値数は入力バッファ内のデータの閾値量に対応する。

【0143】

例12は例8の非一時的コンピュータ可読記憶媒体を含み、バッファは第1バッファであり、クレジットの閾値数はクレジットの第1閾値数であり、かつ、命令は実行された場合に、少なくとも1つのプロセッサに、クレジットの第2の数を第2バッファ内のメモリ利用可能性に関連付けられたクレジットの第2閾値数と比較すること、(1)クレジットの第1の数がクレジットの第1閾値数を満たし、(2)クレジットの第2の数がクレジットの第2閾値数を満たす場合に、計算ビルディングブロックで実行されるワークロードノードを選択することを生じせしめる。

30

【0144】

例13は例12の非一時的コンピュータ可読記憶媒体を含み、第2バッファはワークロードノードに関連付けられた入力バッファであり、クレジットの第2の数は第2バッファに対応しており、クレジットの第2閾値数は入力バッファ内のデータの閾値量に対応する。

【0145】

例14は例8の非一時的コンピュータ可読記憶媒体を含み、クレジットの閾値数はクレジットの第1閾値数であり、ワークロードノードは第1ワークロードノードであり、命令は実行された場合に、少なくとも1つのプロセッサに、(1)クレジットの第1の数がクレジットの第1閾値数を満たし、(2)クレジットの第2の数がクレジットの第2閾値数を満たす場合に、計算ビルディングブロックで実行される第1ワークロードノード及び第2ワークロードノードをスケジューリングすることを生じせしめる。

40

【0146】

例15は、インターフェースする手段であって、クレジットの第1の数をメモリ内へ読み込むためのインターフェースする手段と、比較する手段であって、クレジットの第1の数をバッファ内のメモリ利用可能性に関連付けられたクレジットの閾値数と比較するための比較する手段と、ディスパッチする手段であって、クレジットの第1の数がクレジットの閾値数を満たす場合に、1又は複数の計算ビルディングブロックの最初の一つで実行さ

50

れるワークロードのワークロードノードを選択するためのディスパッチする手段とを備える装置を含む。

【 0 1 4 7 】

例 1 6 は例 1 5 の装置を含み、インターフェースする手段は、インターフェースする手段がクレジットマネージャーからクレジットの第 1 の数を受信した場合にクレジットの第 1 の数をメモリ内へ読み込み、ワークロードノードに関連付けられたデータの 1 又は複数のタイルが 1 又は複数の計算ビルディングブロックの最初の一つからバッファへ送信されるにつれ、バッファに送信された各タイルに対してクレジットをクレジットマネージャーに送信するものである。

【 0 1 4 8 】

例 1 7 は例 1 5 の装置を含み、バッファはワークロードノードに関連付けられた出力バッファであり、クレジットの第 1 の数は出力バッファに対応しており、クレジットの閾値数は出力バッファ内のメモリの閾値量に対応する。

【 0 1 4 9 】

例 1 8 は例 1 5 の装置を含み、バッファはワークロードノードに関連付けられた入力バッファであり、クレジットの第 1 の数は入力バッファに対応しており、クレジットの閾値数は入力バッファ内のデータの閾値量に対応する。

【 0 1 5 0 】

例 1 9 は例 1 5 の装置を含み、バッファは第 1 バッファであり、クレジットの閾値数はクレジットの第 1 閾値数であり、比較する手段はクレジットの第 2 の数を第 2 バッファ内のメモリ利用可能性に関連付けられたクレジットの第 2 閾値数と比較するものであり、ディスパッチする手段は、( 1 ) クレジットの第 1 の数がクレジットの第 1 閾値数を満たし、( 2 ) クレジットの第 2 の数がクレジットの第 2 閾値数を満たす場合に、1 又は複数の計算ビルディングブロックの最初の一つで実行されるワークロードノードを選択するものである。

【 0 1 5 1 】

例 2 0 は例 1 9 の装置を含み、第 2 バッファはワークロードノードに関連付けられた入力バッファであり、クレジットの第 2 の数は入力バッファに対応しており、クレジットの第 2 閾値数は入力バッファ内のデータの閾値量に対応する。

【 0 1 5 2 】

例 2 1 は例 1 5 の装置を含み、クレジットの閾値数はクレジットの第 1 閾値数であり、ワークロードノードは第 1 ワークロードノードであり、( 1 ) クレジットの第 1 の数がクレジットの第 1 閾値数を満たし、( 2 ) クレジットの第 2 の数がクレジットの第 2 閾値数を満たす場合に、ディスパッチする手段は、1 又は複数の計算ビルディングブロックの最初の一つで実行される第 1 ワークロードノード及び第 2 ワークロードノードをスケジューリングするものである。

【 0 1 5 3 】

例 2 2 は、クレジットの第 1 の数をメモリ内へ読み込むことと、クレジットの第 1 の数をバッファ内のメモリ利用可能性に関連付けられたクレジットの閾値数と比較することと、クレジットの第 1 の数がクレジットの閾値数を満たす場合に、1 又は複数の計算ビルディングブロックの最初の一つで実行されるワークロードのワークロードノードを選択することとを備える方法を含む。

【 0 1 5 4 】

例 2 3 は、例 2 2 の方法を含み、クレジットマネージャーからクレジットの第 1 の数を受信した場合にクレジットの第 1 の数をメモリ内へ読み込むことと、ワークロードノードに関連付けられたデータの 1 又は複数のタイルが 1 又は複数の計算ビルディングブロックの最初の一つからバッファへ送信されるにつれ、バッファに送信された各タイルに対してクレジットをクレジットマネージャーに送信することとをさらに含む。

【 0 1 5 5 】

例 2 4 は例 2 2 の方法を含み、バッファはワークロードノードに関連付けられた出力バ

10

20

30

40

50

ッファであり、クレジットの第 1 の数は出力バッファに対応しており、クレジットの閾値数は出力バッファ内のメモリの閾値量に対応する。

【 0 1 5 6 】

例 2 5 は例 2 2 の方法を含み、バッファはワークロードノードに関連付けられた入力バッファであり、クレジットの第 1 の数は入力バッファに対応しており、クレジットの閾値数は入力バッファ内のデータの閾値量に対応する。

【 0 1 5 7 】

特定の例示的な方法、装置及び製造物が本明細書で開示されているが、この特許のカバレッジの範囲はそれに限定されない。反対に、この特許は、この特許の請求項の範囲内に適正に含まれる全ての方法、装置及び製造物を包含する。

【 0 1 5 8 】

以下の請求項は本明細書においてこの参照により本詳細な説明に組み込まれ、各請求項は本開示の別個の実施形態を独自に代表する。

他の可能な項目

[ 項目 1 ]

クレジットの第 1 の数をメモリ内へ読み込むインターフェースと、  
クレジットの上記第 1 の数をバッファ内のメモリ利用可能性に関連付けられたクレジットの閾値数と比較する比較器と、  
クレジットの上記第 1 の数がクレジットの上記閾値数を満たす場合に、上記 1 又は複数の計算ビルディングブロックの最初の一つで実行される上記ワークロードのワークロードノードを選択するディスパッチャと  
を備える装置。

[ 項目 2 ]

上記インターフェースは、  
上記インターフェースがクレジットマネージャーからクレジットの上記第 1 の数を受信した場合にクレジットの上記第 1 の数をメモリ内へ読み込み、  
上記ワークロードノードに関連付けられたデータの 1 又は複数のタイルが上記 1 又は複数の計算ビルディングブロックの上記最初の一つから上記バッファへ送信されるにつれ、  
上記バッファに送信された各タイルに対してクレジットを上記クレジットマネージャーに送信するものである  
項目 1 の装置。

[ 項目 3 ]

上記バッファは上記ワークロードノードに関連付けられた出力バッファであり、クレジットの上記第 1 の数は上記出力バッファに対応しており、クレジットの上記閾値数は上記出力バッファ内のメモリの閾値量に対応する項目 1 の装置。

[ 項目 4 ]

上記バッファは上記ワークロードノードに関連付けられた入力バッファであり、クレジットの上記第 1 の数は上記入力バッファに対応しており、クレジットの上記閾値数は上記入力バッファ内のデータの閾値量に対応する項目 1 の装置。

[ 項目 5 ]

上記バッファは第 1 バッファであり、クレジットの上記閾値数はクレジットの第 1 閾値数であり、上記比較器はクレジットの第 2 の数を第 2 バッファ内のメモリ利用可能性に関連付けられたクレジットの第 2 閾値数と比較するものであり、上記ディスパッチャは、( 1 ) クレジットの上記第 1 の数がクレジットの上記第 1 閾値数を満たし、( 2 ) クレジットの上記第 2 の数がクレジットの上記第 2 閾値数を満たす場合に、上記 1 又は複数の計算ビルディングブロックの上記最初の一つで実行される上記ワークロードノードを選択するものである項目 1 の装置。

[ 項目 6 ]

上記第 2 バッファは上記ワークロードノードに関連付けられた入力バッファであり、クレジットの上記第 2 の数は上記入力バッファに対応しており、クレジットの上記第 2 閾値

10

20

30

40

50

数は上記入力バッファ内のデータの閾値量に対応する項目 5 の装置。

[ 項目 7 ]

クレジットの上記閾値数はクレジットの第 1 閾値数であり、上記ワークロードノードは第 1 ワークロードノードであり、( 1 ) クレジットの上記第 1 の数がクレジットの上記第 1 閾値数を満たし、( 2 ) クレジットの第 2 の数がクレジットの第 2 閾値数を満たす場合に、上記ディスパッチャは、上記 1 又は複数の計算ビルディングブロックの上記最初の一つで実行される上記第 1 ワークロードノード及び第 2 ワークロードノードをスケジューリングするものである項目 1 の装置。

[ 項目 8 ]

実行された場合に、少なくとも 1 つのプロセッサに、

クレジットの第 1 の数をメモリ内へ読み込むこと、

クレジットの上記第 1 の数をバッファ内のメモリ利用可能性に関連付けられたクレジットの閾値数と比較すること、及び、

クレジットの上記第 1 の数がクレジットの上記閾値数を満たす場合に、計算ビルディングブロックで実行される上記ワークロードのワークロードノードを選択すること  
を少なくとも生じせしめる命令を備える非一時的コンピュータ可読記憶媒体。

[ 項目 9 ]

上記命令は実行された場合に、上記少なくとも 1 つのプロセッサに、

クレジットの上記第 1 の数がクレジットマネージャーから受信された場合にクレジットの上記第 1 の数をメモリ内へ読み込むこと、及び、

上記ワークロードノードに関連付けられたデータの 1 又は複数のタイルが上記計算ビルディングブロックから上記バッファへ送信されるにつれ、上記バッファに送信された各タイルに対してクレジットを上記クレジットマネージャーに送信することを生じせしめる項目 8 の非一時的コンピュータ可読記憶媒体。

[ 項目 10 ]

上記バッファは上記ワークロードノードに関連付けられた出力バッファであり、クレジットの上記第 1 の数は上記出力バッファに対応しており、クレジットの上記閾値数は上記出力バッファ内のメモリの閾値量に対応する項目 8 の非一時的コンピュータ可読記憶媒体。

[ 項目 11 ]

上記バッファは上記ワークロードノードに関連付けられた入力バッファであり、クレジットの上記第 1 の数は上記入力バッファに対応しており、クレジットの上記閾値数は上記入力バッファ内のデータの閾値量に対応する項目 8 の非一時的コンピュータ可読記憶媒体。

[ 項目 12 ]

上記バッファは第 1 バッファであり、クレジットの上記閾値数はクレジットの第 1 閾値数であり、かつ、上記命令は実行された場合に、上記少なくとも 1 つのプロセッサに、

クレジットの第 2 の数を第 2 バッファ内のメモリ利用可能性に関連付けられたクレジットの第 2 閾値数と比較すること、

( 1 ) クレジットの上記第 1 の数がクレジットの上記第 1 閾値数を満たし、( 2 ) クレジットの上記第 2 の数がクレジットの上記第 2 閾値数を満たす場合に、上記計算ビルディングブロックで実行される上記ワークロードノードを選択すること  
を生じせしめる項目 8 の非一時的コンピュータ可読記憶媒体。

[ 項目 13 ]

上記第 2 バッファは上記ワークロードノードに関連付けられた入力バッファであり、クレジットの上記第 2 の数は上記第 2 バッファに対応しており、クレジットの第 2 閾値数は上記入力バッファ内のデータの閾値量に対応する項目 12 の非一時的コンピュータ可読記憶媒体。

[ 項目 14 ]

クレジットの上記閾値数はクレジットの第 1 閾値数であり、上記ワークロードノードは第 1 ワークロードノードであり、命令は実行された場合に、上記少なくとも 1 つのプロセッサに、( 1 ) クレジットの上記第 1 の数がクレジットの上記第 1 閾値数を満たし、( 2

10

20

30

40

50

）クレジットの第 2 の数がクレジットの第 2 閾値数を満たす場合に、上記計算ビルディングブロックで実行される上記第 1 ワークロードノード及び第 2 ワークロードノードをスケジューリングすることを生じせしめる項目 8 の非一時的コンピュータ可読記憶媒体。

〔項目 15〕

インターフェースする手段であって、クレジットの第 1 の数をメモリ内へ読み込むための上記インターフェースする手段と、

比較する手段であって、クレジットの上記第 1 の数をバッファ内のメモリ利用可能性に関連付けられたクレジットの閾値数と比較するための上記比較する手段と、

ディスパッチする手段であって、クレジットの上記第 1 の数がクレジットの上記閾値数に一致する場合に、上記 1 又は複数の計算ビルディングブロックの最初の一つで実行される上記ワークロードのワークロードノードを選択するための上記ディスパッチする手段とを備える装置。

10

〔項目 16〕

上記インターフェースする手段は、

上記インターフェースする手段がクレジットマネージャーからクレジットの上記第 1 の数を受信した場合にクレジットの上記第 1 の数をメモリ内へ読み込み、

上記ワークロードノードに関連付けられたデータの 1 又は複数のタイルが上記 1 又は複数の計算ビルディングブロックの上記最初の一つから上記バッファへ送信されるにつれ、上記バッファに送信された各タイルに対してクレジットを上記クレジットマネージャーに送信する

20

ものである項目 15 の装置。

〔項目 17〕

上記バッファは上記ワークロードノードに関連付けられた出力バッファであり、クレジットの上記第 1 の数は上記出力バッファに対応しており、クレジットの上記閾値数は上記出力バッファ内のメモリの閾値量に対応する項目 15 の装置。〔項目 18〕上記バッファは上記ワークロードノードに関連付けられた入力バッファであり、クレジットの上記第 1 の数は上記入力バッファに対応しており、クレジットの上記閾値数は上記入力バッファ内のデータの閾値量に対応する項目 15 の装置。

〔項目 19〕

上記バッファは第 1 バッファであり、クレジットの上記閾値数はクレジットの第 1 閾値数であり、上記比較する手段はクレジットの第 2 の数を第 2 バッファ内のメモリ利用可能性に関連付けられたクレジットの第 2 閾値数と比較するものであり、上記ディスパッチする手段は、（1）クレジットの上記第 1 の数がクレジットの上記第 1 閾値数を満たし、（2）クレジットの上記第 2 の数がクレジットの上記第 2 閾値数を満たす場合に、上記 1 又は複数の計算ビルディングブロックの上記最初の一つで実行される上記ワークロードノードを選択するものである項目 15 の装置。

30

〔項目 20〕

上記第 2 バッファは上記ワークロードノードに関連付けられた入力バッファであり、クレジットの上記第 2 の数は上記入力バッファに対応しており、クレジットの上記第 2 閾値数は上記入力バッファ内のデータの閾値量に対応する項目 19 の装置。

40

〔項目 21〕

クレジットの上記閾値数はクレジットの第 1 閾値数であり、上記ワークロードノードは第 1 ワークロードノードであり、（1）クレジットの上記第 1 の数がクレジットの上記第 1 閾値数を満たし、（2）クレジットの第 2 の数がクレジットの第 2 閾値数を満たす場合に、上記ディスパッチする手段は、上記 1 又は複数の計算ビルディングブロックの上記最初の一つで実行される上記第 1 ワークロードノード及び第 2 ワークロードノードをスケジューリングするものである項目 15 の装置。

〔項目 22〕

クレジットの第 1 の数をメモリ内へ読み込むことと、

クレジットの上記第 1 の数をバッファ内のメモリ利用可能性に関連付けられたクレジッ

50

トの閾値数と比較することと、

クレジットの上記第 1 の数がクレジットの上記閾値数に一致する場合に、上記 1 又は複数の計算ビルディングブロックの最初の一つで実行される上記ワークロードのワークロードノードを選択することと

を備える方法。

〔項目 2 3〕

クレジットマネージャーからクレジットの上記第 1 の数を受信する場合にクレジットの上記第 1 の数をメモリ内へ読み込むことと、

上記ワークロードノードに関連付けられたデータの 1 又は複数のタイルが上記 1 又は複数の計算ビルディングブロックの上記最初の一つから上記バッファへ送信されるにつれ、上記バッファに送信された各タイルに対してクレジットを上記クレジットマネージャーに送信することと

をさらに含む項目 2 2 の方法。

〔項目 2 4〕

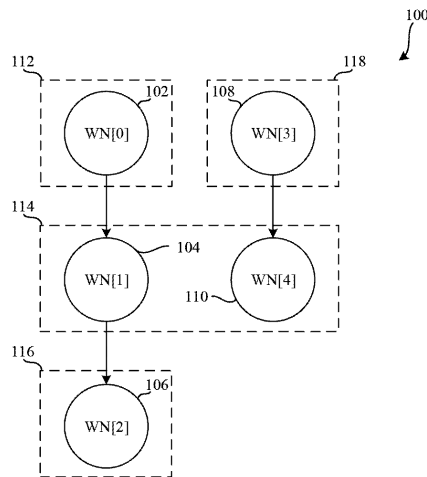
上記バッファは上記ワークロードノードに関連付けられた出力バッファであり、クレジットの上記第 1 の数は上記出力バッファに対応しており、クレジットの上記閾値数は上記出力バッファ内のメモリの閾値量に対応する項目 2 2 の方法。

〔項目 2 5〕

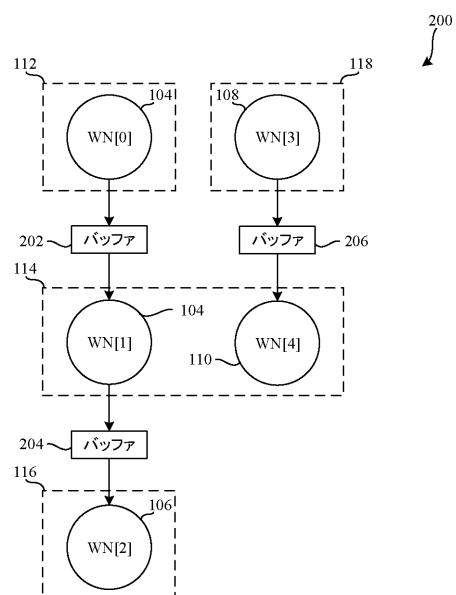
上記バッファは上記ワークロードノードに関連付けられた入力バッファであり、クレジットの上記第 1 の数は上記入力バッファに対応しており、クレジットの上記閾値数は上記入力バッファ内のデータの閾値量に対応する項目 2 2 の方法。

【図面】

【図 1】



【図 2】



10

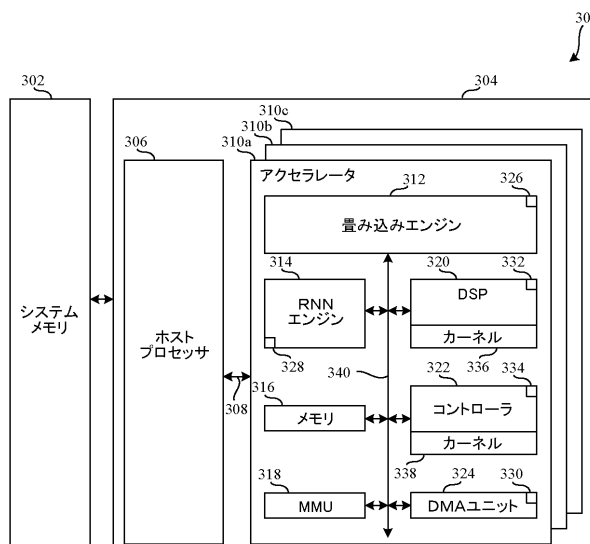
20

30

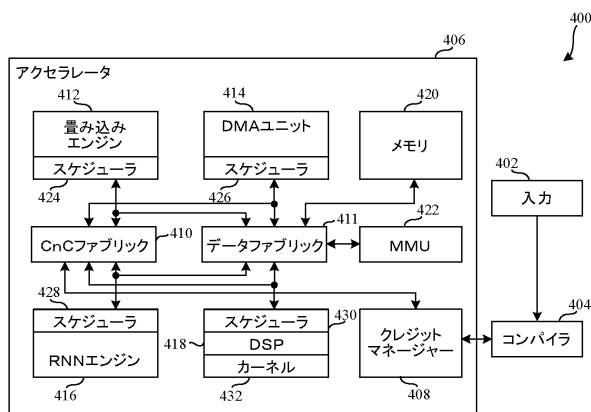
40

50

【 図 3 】

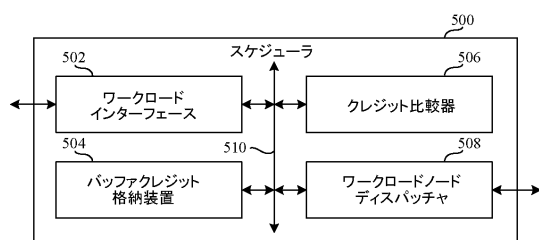


【圖 4】

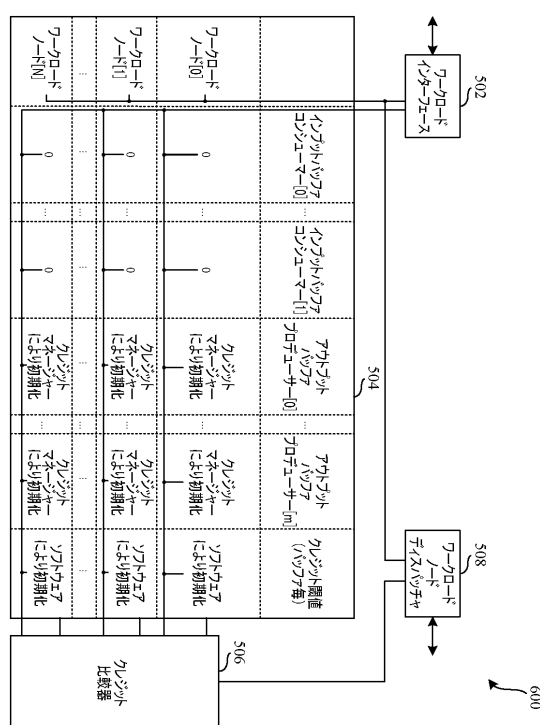


10

【圖 5】



【 図 6 】



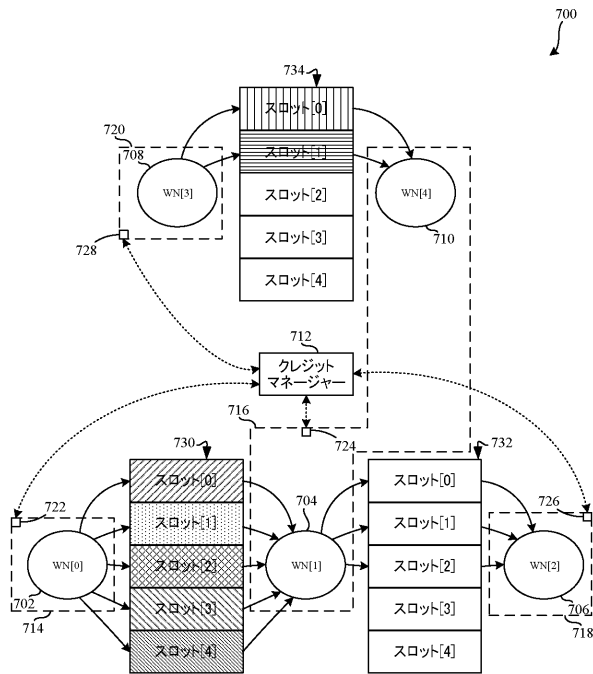
20

30

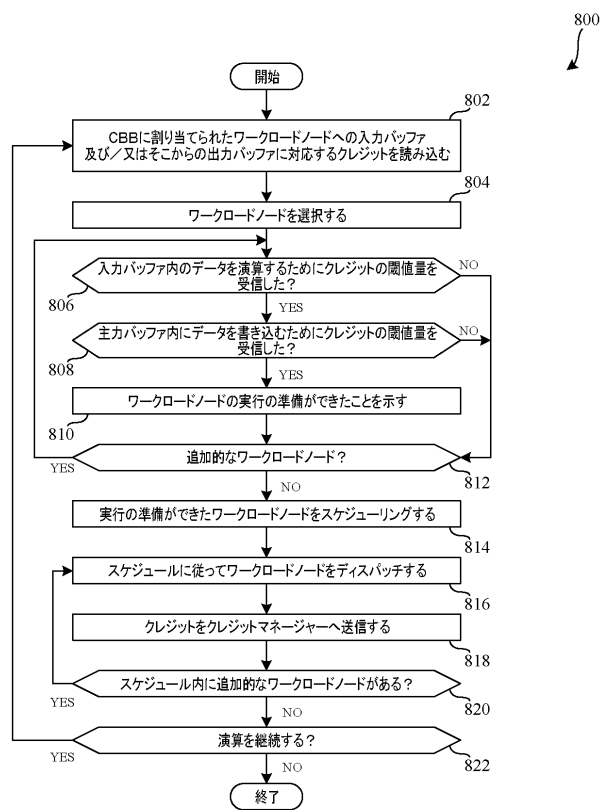
40

50

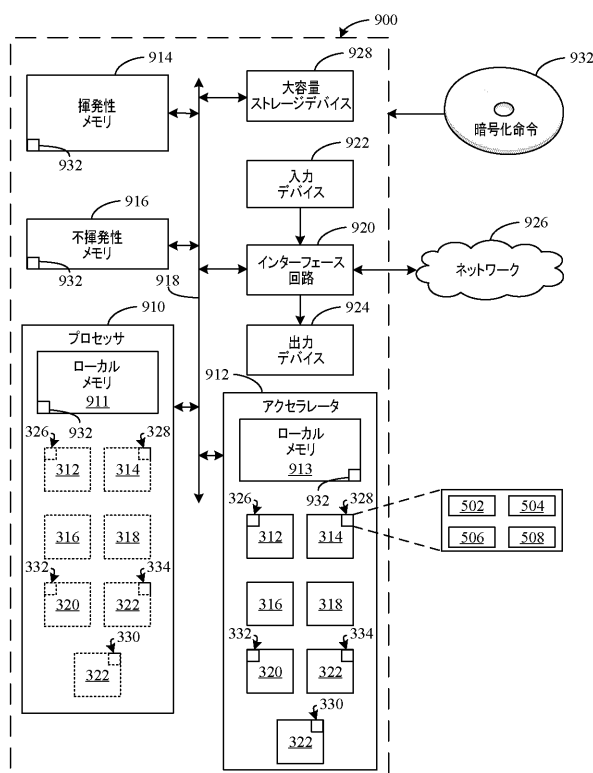
【図 7】



【図 8】



【図 9】



10

20

30

40

50

## フロントページの続き

- レッジ ブーレバード・２２００ インテル・コーポレーション内  
(72)発明者 ロネン ガバイ  
アメリカ合衆国 ９５０５４ カリフォルニア州・サンタクララ・ミッション カレッジ ブーレバ  
ード・２２００ インテル・コーポレーション内  
(72)発明者 ロニ ロスナー  
アメリカ合衆国 ９５０５４ カリフォルニア州・サンタクララ・ミッション カレッジ ブーレバ  
ード・２２００ インテル・コーポレーション内  
(72)発明者 ジギ ウォルター  
アメリカ合衆国 ９５０５４ カリフォルニア州・サンタクララ・ミッション カレッジ ブーレバ  
ード・２２００ インテル・コーポレーション内  
(72)発明者 オレン アガム  
アメリカ合衆国 ９５０５４ カリフォルニア州・サンタクララ・ミッション カレッジ ブーレバ  
ード・２２００ インテル・コーポレーション内  
審査官 坂東 博司  
(56)参考文献 特開２００５－０１８６２０（ＪＰ，Ａ）  
米国特許出願公開第２００４／０２６８０８３（ＵＳ，Ａ１）  
特表２０１７－５２５０４７（ＪＰ，Ａ）  
特開平０７－０２１１４４（ＪＰ，Ａ）  
特開２０１２－１９０４１５（ＪＰ，Ａ）  
米国特許出願公開第２０１２／０２３９８３３（ＵＳ，Ａ１）  
米国特許出願公開第２０１６／０１４００７１（ＵＳ，Ａ１）  
米国特許出願公開第２０１９／００５０２６１（ＵＳ，Ａ１）  
(58)調査した分野 (Int.Cl.，ＤＢ名)  
Ｇ０６Ｆ ９／５０  
Ｇ０６Ｆ ９／３８  
Ｇ０６Ｆ ９／４８