

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6581727号
(P6581727)

(45) 発行日 令和1年9月25日(2019.9.25)

(24) 登録日 令和1年9月6日(2019.9.6)

(51) Int. Cl. F I
G06F 9/54 (2006.01) G O 6 F 9/54 B
G06F 9/50 (2006.01) G O 6 F 9/50 1 5 O E

請求項の数 30 (全 34 頁)

(21) 出願番号	特願2018-529031 (P2018-529031)	(73) 特許権者	509123208
(86) (22) 出願日	平成28年12月5日 (2016.12.5)		アビニシオ テクノロジー エルエルシー
(65) 公表番号	特表2019-502989 (P2019-502989A)		アメリカ合衆国 02421 マサチュー
(43) 公表日	平成31年1月31日 (2019.1.31)		セッツ州 レキシントン スプリング ス
(86) 国際出願番号	PCT/US2016/064898		トリート 201
(87) 国際公開番号	W02017/105888	(74) 代理人	100107984
(87) 国際公開日	平成29年6月22日 (2017.6.22)		弁理士 廣田 雅紀
審査請求日	平成30年7月25日 (2018.7.25)	(74) 代理人	100102255
(31) 優先権主張番号	62/268,762		弁理士 小澤 誠次
(32) 優先日	平成27年12月17日 (2015.12.17)	(74) 代理人	100096482
(33) 優先権主張国・地域又は機関	米国 (US)		弁理士 東海 裕作
		(74) 代理人	100188352
			弁理士 松田 一弘
		(74) 代理人	100113860
			弁理士 松橋 泰典

最終頁に続く

(54) 【発明の名称】 動的なパーティショニングを使用してデータを処理すること

(57) 【特許請求の範囲】

【請求項1】

保留キュー、再配置キュー、及び準備完了キューを含むワークユニットを記憶するキューを使用して、コンピューティングシステムにおいてデータを処理するための方法であって、

前記コンピューティングシステムの1又は2以上のノードのストレージに記憶された複数の保留キューにワークユニットを受信するステップと、

前記1又は2以上のノードの各々における少なくとも1つのプロセッサを使用して、前記複数の保留キューから前記ワークユニットを処理するステップとを含み、前記処理するステップが、

前記保留キューごとに、その保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの一つを、パーティショニングルールに従って決定するステップであって、各パーティションが一意的保留キューに関連付けられる、前記決定するステップと、

各保留キューから、ワークユニットを、その保留キューと同じストレージに記憶された少なくとも2つの他のキューの対応する一つに分離するステップであって、前記2つの他のキューが、ワークユニットが記憶されている前記保留キューに関連付けられていないパーティションが割り当てられたワークユニットごとの再配置キューと、ワークユニットが記憶されている前記保留キューに関連付けられているパーティションが割り当てられたワークユニットごとの準備完了キューとを含む、前記分離するステップと、

各再配置キュー内の各ワークユニットを、そのワークユニットに割り当てられた前記パーティションに関連付けられる前記保留キューに再配置するステップと、

前記複数の保留キューが空になる前に、前記パーティショニングルールの新しいパーティショニングルールへの変更に応答して、前記ワークユニットを再パーティショニングするステップと

を含む、前記方法。

【請求項 2】

パーティショニングルールの変更が、複数のパーティション内のパーティションの総数の変更を含む、請求項 1 に記載の方法。

【請求項 3】

コンピューティングシステムの 1 又は 2 以上のノードが複数のノードを含み、前記複数のノードのうちの少なくとも第 1 のノードが、第 1 のパーティションに関連付けられる第 1 の保留キューを記憶し、前記第 1 の保留キューに対応する第 1 の再配置キュー及び第 1 の準備完了キューを記憶するストレージを有し、前記複数のノードのうちの少なくとも第 2 のノードが、第 2 のパーティションに関連付けられる第 2 の保留キューを記憶し、前記第 2 の保留キューに対応する第 2 の再配置キュー及び第 2 の準備完了キューを記憶するストレージを有する、請求項 1 又は 2 に記載の方法。

【請求項 4】

パーティショニングルールの変更が、保留キューのうちのいずれかを記憶するために使用されるストレージを有する複数のノード内のノードの総数の変更を含む、請求項 3 に記載の方法。

【請求項 5】

第 1 の再配置キュー内に 1 又は 2 以上のワークユニットを再配置するステップが、第 1 のノードから第 2 のノードに前記 1 又は 2 以上のワークユニットを含むメッセージを送信するステップを含み、第 2 のパーティションが、前記メッセージに含まれる前記 1 又は 2 以上のワークユニットの各々に割り当てられる、請求項 3 又は 4 に記載の方法。

【請求項 6】

第 1 のノードのストレージがまた、第 3 のパーティションに関連付けられる第 3 の保留キューを記憶し、前記第 3 の保留キューに対応する第 3 の再配置キュー及び第 3 の準備完了キューを記憶する、請求項 3 ~ 5 のいずれかに記載の方法。

【請求項 7】

再パーティショニングするステップが、特定の保留キューに対応する再配置キュー及び準備完了キューに記憶された任意のワークユニットを、その特定の保留キューに戻すステップを含む、請求項 1 に記載の方法。

【請求項 8】

パーティショニングルールが、各ワークユニットに対して複数のパーティションのうちの 1 つのコンテンツに依存する割当てを使用し、特定のワークユニットに割り当てられた前記パーティションが、前記ワークユニットに含まれる少なくとも 1 つの値に基づいて決定される、請求項 7 に記載の方法。

【請求項 9】

各保留キューが、初期パーティション決定がまだ行われていないワークユニットを記憶するための第 1 の部分と、それらのワークユニットをその保留キューに再配置させる初期パーティション決定が行われたワークユニットを記憶するための第 2 の部分とを含む、請求項 7 又は 8 に記載の方法。

【請求項 10】

ワークユニットを各保留キューから対応する再配置キュー又は準備完了キューに分離するステップが、第 2 の部分が空である場合にのみ第 1 の部分に記憶されたワークユニットを分離するステップを含む、請求項 9 に記載の方法。

【請求項 11】

少なくとも第 1 の保留キューが、パーティショニングルールの変更後にその保留キュー

10

20

30

40

50

に戻された任意のワークユニットを記憶するための第3の部分を含む、請求項9又は10に記載の方法。

【請求項12】

再パーティショニングするステップが、

第1の保留キューに対して、新しいパーティショニングルールに従って前記第1の保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの一つを決定するステップと、

ワークユニットを前記第1の保留キューから対応する再配置キュー又は準備完了キューに分離するステップと

をさらに含む、請求項11に記載の方法。

10

【請求項13】

ワークユニットを第1の保留キューから対応する再配置キュー又は準備完了キューに分離するステップが、第3の部分が空である場合にのみ第2の部分に記憶されたワークユニットを分離するステップと、前記第2の部分と前記第3の部分の両方が空である場合にのみ第1の部分に記憶されたワークユニットを分離するステップとを含む、請求項12に記載の方法。

【請求項14】

パーティショニングルールが、複数のパーティションのうちの一つのコンテンツに依存しない割当てを各ワークユニットに使用し、特定のワークユニットに割り当てられた前記パーティションが、前記ワークユニットに含まれる任意の値から独立して決定される、請求項1に記載の方法。

20

【請求項15】

パーティショニングルールが、複数のパーティションのうちの一つの負荷バランス割当てを各ワークユニットに使用し、特定のワークユニットに割り当てられた前記パーティションが、異なるパーティションに対する負荷の定量的尺度に基づいて決定される、請求項14に記載の方法。

【請求項16】

異なるパーティションに対する負荷の定量的尺度が、異なるパーティションに関連付けられる保留キューに再配置されるべき記憶されたワークユニットの異なるセットのバックログの定量的尺度である、請求項15に記載の方法。

30

【請求項17】

第1のノード上で実行する1又は2以上のワーカーを使用して、第1の準備完了キューからワークユニットを処理するステップをさらに含む、請求項1に記載の方法。

【請求項18】

再パーティショニングするステップが、新しいパーティショニングルールが、(1)各ワークユニットに対する複数のパーティションのうちの一つのコンテンツに依存する割当てであって、特定のワークユニットに割り当てられた前記パーティションが、前記ワークユニットに含まれる少なくとも一つの値に基づいて決定される割当てを使用するか、又は、(2)各ワークユニットに対する前記複数のパーティションのうちの一つのコンテンツに依存しない割当てであって、特定のワークユニットに割り当てられた前記パーティションが、前記ワークユニットに含まれる任意の値から独立して決定される割当てを使用するかを決定するステップを含む、請求項17に記載の方法。

40

【請求項19】

再パーティショニングするステップが、第1の準備完了キューからワークユニットを処理するステップが、前記ワークユニットが第1の保留キューに受信されたのと同じ順序で前記第1の保留キューから発するワークユニットを処理するステップを含むかどうかを決定するステップを含む、請求項18に記載の方法。

【請求項20】

再パーティショニングするステップが、古いパーティショニングルールから新しいパーティショニングルールにパーティショニングルールを変更する前に、及び前記新しいパー

50

パーティショニングルールがコンテンツに依存する割当てを使用すると決定し、第1の準備完了キューからワークユニットを処理するステップが、前記ワークユニットが第1の保留キューに受信されたのと同じ順序で前記第1の保留キューから発するワークユニットを処理するステップを含むと決定した後、すべての再配置キューを排出させるステップを含む、請求項19に記載の方法。

【請求項21】

排出させるステップが、(1)各再配置キュー内の各ワークユニットを、古いパーティショニングルールに従ってそのワークユニットに割り当てられたパーティションに関連付けられる保留キューに再配置するステップと、(2)前記排出中に、第1の準備完了キューからのワークユニットの処理を中断するステップとを含む、請求項20に記載の方法。

10

【請求項22】

パーティショニングルールが、各ワークユニットに対して複数のパーティションのうちの1つのキーベースの割当てを使用し、特定のワークユニットに割り当てられた前記パーティションが、前記ワークユニットに含まれるキー値の関数に基づいて決定される、請求項1に記載の方法。

【請求項23】

コンピューティングシステムの1又は2以上のノードが複数のノードを含み、前記複数のノードのストレージが、複数の保留キューと、キー値に基づいて前記複数のノードの各々にアクセス可能な共有データとを含む、請求項22に記載の方法。

【請求項24】

20

共有データが、第1のノードのストレージに記憶された第1の準備完了キューから1又は2以上のワークユニットの処理中に、複数のノードのうちの前記第1のノード内の少なくとも第1のプロセッサによってアクセスされる、請求項23に記載の方法。

【請求項25】

共有データが、ワークユニットに割り当てられた複数のパーティションに従ってパーティショニングされ、前記共有データにアクセスするために使用される各キー値が、前記複数のパーティションのうちの一意的パーティションに関連付けられる、請求項23又は24に記載の方法。

【請求項26】

パーティショニングルールの変更に応じてワークユニットを再パーティショニングするステップが、新しいパーティショニングルールに従って共有データを再パーティショニングし始めた後に発生する、請求項25に記載の方法。

30

【請求項27】

新しいパーティショニングルールに従って共有データを再パーティショニングするステップが、

前記共有データに記憶された複数の参照ユニットの参照ユニットごとに、前記新しいパーティショニングルールに従ってその参照ユニットに割り当てられた複数のパーティションのうちの1つを決定するステップと、

以前のノードから新しいノードに、参照ユニットが記憶されている前記ノードに関連付けられていないパーティションが割り当てられたその任意の参照ユニットを転送するステップと、

40

前記以前のノードと前記新しいノードの両方において前記共有データを再パーティショニングする前記ステップの間に、新しいノードに転送された前記参照ユニットのうち少なくとも1つにアクセスしようと試みるステップとを含む、請求項26に記載の方法。

【請求項28】

処理するステップがさらに、

各準備完了キューから少なくとも1つの結果を提供するステップであって、前記各準備完了キューからの前記結果が、その準備完了キュー内に記憶された1又は2以上のワークユニットに基づく、前記提供するステップと、

50

前記結果が提供された前記ワークユニットを前記各準備完了キューから除去するステップと

を含む、請求項 1 ~ 27 のいずれかに記載の方法。

【請求項 29】

保留キュー、再配置キュー、及び準備完了キューを含むワークユニットを記憶するキューを使用して、データを処理するために、コンピュータ可読媒体上に非一時的な形式で記憶されるソフトウェアであって、前記ソフトウェアが、コンピューティングシステムに、

前記コンピューティングシステムの 1 又は 2 以上のノードのストレージに記憶された複数の保留キューにワークユニットを受信することと、

前記 1 又は 2 以上のノードの各々における少なくとも 1 つのプロセッサを使用して、前記複数の保留キューから前記ワークユニットを処理することと

を行わせるための命令を含み、前記処理することが、

前記保留キューごとに、その保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの 1 つを、パーティショニングルールに従って決定することであって、各パーティションが一意的保留キューに関連付けられる、前記決定することと、

各保留キューから、ワークユニットを、その保留キューと同じストレージに記憶された少なくとも 2 つの他のキューの対応する 1 つに分離することであって、前記 2 つの他のキューが、ワークユニットが記憶されている前記保留キューに関連付けられていないパーティションが割り当てられたワークユニットごとの再配置キューと、ワークユニットが記憶されている前記保留キューに関連付けられているパーティションが割り当てられたワークユニットごとの準備完了キューとを含む、前記分離することと、

各再配置キュー内の各ワークユニットを、そのワークユニットに割り当てられた前記パーティションに関連付けられる前記保留キューに再配置することと、

前記複数の保留キューが空になる前に、前記パーティショニングルールの新しいパーティショニングルールへの変更に応答して、前記ワークユニットを再パーティショニングすることと

を含む、前記ソフトウェア。

【請求項 30】

保留キュー、再配置キュー、準備完了キューを含むワークユニットを記憶するキューを使用してデータを処理するためのコンピューティングシステムであって、

1 又は 2 以上のノードを含み、それぞれが、

複数の保留キューのうちのストレージに記憶された少なくとも 1 つの保留キューにワークユニットを受信するように構成された入力デバイス又はポートと、

前記保留キューから前記ワークユニットを処理するように構成された少なくとも 1 つのプロセッサとを含み、前記複数の保留キューのために処理するステップが、

前記保留キューごとに、その保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの 1 つを、パーティショニングルールに従って決定するステップであって、各パーティションが一意的保留キューに関連付けられる、前記決定するステップと、

各保留キューから、ワークユニットを、その保留キューと同じストレージに記憶された少なくとも 2 つの他のキューの対応する 1 つに分離するステップであって、前記 2 つの他のキューが、ワークユニットが記憶されている前記保留キューに関連付けられていないパーティションが割り当てられたワークユニットごとの再配置キューと、ワークユニットが記憶されている前記保留キューに関連付けられているパーティションが割り当てられたワークユニットごとの準備完了キューとを含む、前記分離するステップと、

各再配置キュー内の各ワークユニットを、そのワークユニットに割り当てられた前記パーティションに関連付けられる前記保留キューに再配置するステップと、

前記複数の保留キューが空になる前に、前記パーティショニングルールの新しいパーティショニングルールへの変更に応答して、前記ワークユニットを再パーティショニングす

10

20

30

40

50

るステップと
を含む、前記コンピューティングシステム。

【発明の詳細な説明】

【技術分野】

【0001】

関連出願の相互参照

本出願は、2015年12月17日に出願された米国特許出願第62/268,762号に対する優先権を主張する。

【0002】

本明細書は、動的なパーティショニングを使用してデータを処理することに関する。

【背景技術】

【0003】

複数のコンピューティングノード（例えば、サーバクラスタのサーバ、及び/又はマルチコアプロセッサのコア）を使用して並列にデータを処理するコンピューティングシステムは、パーティショナを使用してそれらのコンピューティングノード間でデータを分散することができる。データは、ワークの個々のユニットを表す「ワークユニット」において処理され得、ワークユニットは、任意の数のそれぞれのフィールドの値を有する個々のレコードの形態であり得る。或いは、ワークユニットは、レコードに関連付けられるメタデータ、処理されるべきイベントを記述する情報、又は実施されるべきワークの個々の部分
20
を表す任意の他の情報の形態であり得る。各コンピューティングノードには異なるパーティションが割り当てられ得、そのコンピューティングノード上で動作する1又は2以上のワーカー（例えば、プロセス又はスレッド）は、その割り当てられたパーティション内のワークユニットを処理する。レコードが処理されると、フィールド内の値が変更され得る。データ処理のいくつかの段階において、ワークユニットは、そのレコードに適用されるべき識別されたタスクとともに、個々のレコードの形態であり得る。パーティショナは、特定のワークユニットを、そのワークユニットが属する特定のパーティションに分散するために、パーティショニングルールを使用する。いくつかのコンピューティングシステムでは、処理は段階的に実施され、ワークユニットは、同じパーティション又は異なる数のパーティションの間で後続の段階のために再パーティショニングされる。例えば、各コンピューティングノードは、入力ワークユニットを受信するための入力バッファと、出力ワークユニットを提供するための出力バッファとを有し得る。パーティショナは、出力バッファごとに、前の段階からの出力ワークユニットを、次の段階の入力ワークユニットとして処理されるべき入力バッファ間で分散する。そのような再パーティショニングは、通常、前の段階の計算がすべて完了した後に実施される。そのような再パーティショニングに使用されるパーティショニングルールは変更される可能性があるが、変更は異なるワークユニットのセット間で発生する必要がある。言い換えれば、所与のワークユニットのセットについてパーティショニングルールは静的であり、したがって、次のパーティショニングルールを使用して次のセット内の任意のワークユニットがパーティショニングされ得る前に、所与のセット内のすべてのワークユニットが現在のパーティショニングルールを使用
30
して完全にパーティショニングされる。
40

【発明の概要】

【課題を解決するための手段】

【0004】

一態様では、一般に、保留キュー、再配置キュー、及び準備完了キューを含むワークユニットを記憶するキューを使用して、コンピューティングシステムにおいてデータを処理するための方法は、コンピューティングシステムの1又は2以上のノードのストレージに記憶された複数の保留キューにワークユニットを受信するステップと、1又は2以上のノードの各々における少なくとも1つのプロセッサを使用して、複数の保留キューからワークユニットを処理するステップとを含む。処理するステップは、保留キューごとに、その
50

保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの1つを、パーティショニングルールに従って決定するステップであって、各パーティションが一意的保留キューに関連付けられる、前記決定するステップと、各保留キューから、ワークユニットを、その保留キューと同じストレージに記憶された少なくとも2つの他のキューの対応する1つに分離するステップであって、2つの他のキューが、ワークユニットが記憶されている保留キューに関連付けられていないパーティションが割り当てられたワークユニットごとの再配置キューと、ワークユニットが記憶されている保留キューに関連付けられているパーティションが割り当てられたワークユニットごとの準備完了キューとを含む、前記分離するステップと、各再配置キュー内の各ワークユニットを、そのワークユニットに割り当てられたパーティションに関連付けられる保留キューに再配置するステップと、複数の保留キューが空になる前に、パーティショニングルールの新しいパーティショニングルールへの変更に応答して、ワークユニットを再パーティショニングするステップとを含む。

10

【0005】

態様は、以下の特徴のうちの1又は2以上を含んでもよい。疑念を避けるために、以下の特徴のうちのいずれかは、特徴が相互に排他的でない限り、他の以下の特徴のうちのいずれかと組み合わせて使用され得る（例えば、コンテンツに依存する割当てに関連する特徴と、コンテンツに依存しない割当てに関連する特徴は、代替である）。

【0006】

パーティショニングルールの変更は、複数のパーティション内のパーティションの総数の変更を含む。

20

【0007】

コンピューティングシステムの1又は2以上のノードは複数のノードを含み、複数のノードのうちの少なくとも第1のノードは、第1のパーティションに関連付けられる第1の保留キューを記憶し、第1の保留キューに対応する第1の再配置キュー及び第1の準備完了キューを記憶するストレージを有し、複数のノードのうちの少なくとも第2のノードが、第2のパーティションに関連付けられる第2の保留キューを記憶し、第2の保留キューに対応する第2の再配置キュー及び第2の準備完了キューを記憶するストレージを有する。

【0008】

パーティショニングルールの変更は、保留キューのうちのいずれかを記憶するために使用されるストレージを有する複数のノード内のノードの総数の変更を含む。

30

【0009】

第1の再配置キュー内に1又は2以上のワークユニットを再配置するステップは、第1のノードから第2のノードに1又は2以上のワークユニットを含むメッセージを送信するステップを含み、第2のパーティションは、メッセージに含まれる1又は2以上のワークユニットの各々に割り当てられる。

【0010】

第1のノードのストレージはまた、第3のパーティションに関連付けられる第3の保留キューを記憶し、第3の保留キューに対応する第3の再配置キュー及び第3の準備完了キューを記憶する。

40

【0011】

再パーティショニングするステップは、特定の保留キューに対応する再配置キュー及び準備完了キューに記憶された任意のワークユニットを、その特定の保留キューに戻すステップを含む。

【0012】

パーティショニングルールは、各ワークユニットに対して複数のパーティションのうちの1つのコンテンツに依存する割当てを使用し、特定のワークユニットに割り当てられたパーティションは、ワークユニットに含まれる少なくとも1つの値に基づいて決定される。

50

【 0 0 1 3 】

各保留キューは、初期パーティション決定がまだ行われていないワークユニットを記憶するための第1の部分と、それらのワークユニットをその保留キューに再配置させる初期パーティション決定が行われたワークユニットを記憶するための第2の部分とを含む。

【 0 0 1 4 】

ワークユニットを各保留キューから対応する再配置キュー又は準備完了キューに分離するステップは、第2の部分が空である場合にのみ第1の部分に記憶されたワークユニットを分離するステップを含む。

【 0 0 1 5 】

少なくとも第1の保留キューは、パーティショニングルールの変更後にその保留キューに戻された任意のワークユニットを記憶するための第3の部分を含む。

10

【 0 0 1 6 】

再パーティショニングするステップは、第1の保留キューに対して、新しいパーティショニングルールに従って第1の保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの1つを決定するステップと、ワークユニットを第1の保留キューから対応する再配置キュー又は準備完了キューに分離するステップとをさらに含む。

【 0 0 1 7 】

ワークユニットを第1の保留キューから対応する再配置キュー又は準備完了キューに分離するステップは、第3の部分が空である場合にのみ第2の部分に記憶されたワークユニットを分離するステップと、第2の部分と第3の部分の両方が空である場合にのみ第1の部分に記憶されたワークユニットを分離するステップとを含む。

20

【 0 0 1 8 】

パーティショニングルールは、複数のパーティションのうちの1つのコンテンツに依存しない割当てを各ワークユニットに使用し、特定のワークユニットに割り当てられたパーティションは、ワークユニットに含まれる任意の値から独立して決定される。

【 0 0 1 9 】

パーティショニングルールは、複数のパーティションのうちの1つの負荷バランス割当てを各ワークユニットに使用し、特定のワークユニットに割り当てられたパーティションは、異なるパーティションに対する負荷の定量的尺度に基づいて決定される。

30

【 0 0 2 0 】

異なるパーティションに対する負荷の定量的尺度は、異なるパーティションに関連付けられる保留キューに再配置されるべき記憶されたワークユニットの異なるセットのバックログの定量的尺度である。

【 0 0 2 1 】

本方法は、第1のノード上で実行する1又は2以上のワーカーを使用して、第1の準備完了キューからワークユニットを処理するステップをさらに含む。

【 0 0 2 2 】

再パーティショニングするステップは、新しいパーティショニングルールが、(1)各ワークユニットに対する複数のパーティションのうちの1つのコンテンツに依存する割当てであって、特定のワークユニットに割り当てられたパーティションが、ワークユニットに含まれる少なくとも1つの値に基づいて決定される割当てを使用するか、又は、(2)各ワークユニットに対する複数のパーティションのうちの1つのコンテンツに依存しない割当てであって、特定のワークユニットに割り当てられたパーティションが、ワークユニットに含まれる任意の値から独立して決定される割当てを使用するかを決定するステップを含む。

40

【 0 0 2 3 】

再パーティショニングするステップは、第1の準備完了キューからワークユニットを処理するステップが、ワークユニットが第1の保留キューに受信されたのと同じ順序で第1の保留キューから発するワークユニットを処理するステップを含むかどうかを決定するス

50

テップを含む。

【0024】

再パーティショニングするステップは、古いパーティショニングルールから新しいパーティショニングルールにパーティショニングルールを変更する前に、及び新しいパーティショニングルールがコンテンツに依存する割当てを使用すると決定し、第1の準備完了キューからワークユニットを処理するステップが、ワークユニットが第1の保留キューに受信されたのと同じ順序で第1の保留キューから発するワークユニットを処理するステップを含むと決定した後に、すべての再配置キューを排出させるステップを含む。

【0025】

排出させるステップは、(1)各再配置キュー内の各ワークユニットを、古いパーティショニングルールに従ってそのワークユニットに割り当てられたパーティションに関連付けられる保留キューに再配置するステップと、(2)排出中に、第1の準備完了キューからのワークユニットの処理を中断するステップとを含む。

10

【0026】

パーティショニングルールは、各ワークユニットに対して複数のパーティションのうちの1つのキーベースの割当てを使用し、特定のワークユニットに割り当てられたパーティションは、ワークユニットに含まれるキー値の関数に基づいて決定される。

【0027】

コンピューティングシステムの1又は2以上のノードは複数のノードを含み、複数のノードのストレージは、複数の保留キューと、キー値に基づいて複数のノードの各々にアクセス可能な共有データとを含む。

20

【0028】

共有データは、第1のノードのストレージに記憶された第1の準備完了キューから1又は2以上のワークユニットの処理中に、複数のノードのうちの第1のノード内の少なくとも第1のプロセッサによってアクセスされる。

【0029】

共有データは、ワークユニットに割り当てられた複数のパーティションに従ってパーティショニングされ、共有データにアクセスするために使用される各キー値は、複数のパーティションのうちの一意的パーティションに関連付けられる。

【0030】

パーティショニングルールの変更に応じてワークユニットを再パーティショニングするステップは、新しいパーティショニングルールに従って共有データを再パーティショニングし始めた後に発生する。

30

【0031】

新しいパーティショニングルールに従って共有データを再パーティショニングするステップは、共有データに記憶された複数の参照ユニットの参照ユニットごとに、新しいパーティショニングルールに従ってその参照ユニットに割り当てられた複数のパーティションのうちの1つを決定するステップと、以前のノードから新しいノードに、参照ユニットが記憶されているノードに関連付けられていないパーティションが割り当てられたその任意の参照ユニットを転送するステップと、以前のノードと新しいノードの両方において共有データを再パーティショニングするステップの間に、新しいノードに転送された参照ユニットのうち少なくとも1つにアクセスしようと試みるステップとを含む。

40

【0032】

処理するステップはさらに、各準備完了キューから少なくとも1つの結果を提供するステップであって、各準備完了キューからの結果が、その準備完了キュー内に記憶された1又は2以上のワークユニットに基づく、ステップと、結果が提供されたワークユニットを各準備完了キューから除去するステップとを含む。

【0033】

別の態様では、一般に、ソフトウェアは、保留キュー、再配置キュー、及び準備完了キューを含むワークユニットを記憶するキューを使用して、データを処理するために、コン

50

コンピュータ可読媒体上に非一時的な形式で記憶される。ソフトウェアは、コンピューティングシステムに、コンピューティングシステムの1又は2以上のノードのストレージに記憶された複数の保留キューにワークユニットを受信することと、1又は2以上のノードの各々における少なくとも1つのプロセッサを使用して、複数の保留キューからワークユニットを処理することとを行わせるための命令を含む。処理することは、保留キューごとに、その保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの一つを、パーティショニングルールに従って決定することであって、各パーティションが一意的保留キューに関連付けられる、前記決定することと、各保留キューから、ワークユニットを、その保留キューと同じストレージに記憶された少なくとも2つの他のキューの対応する1つに分離することであって、2つの他のキューが、ワークユニットが記憶されている保留キューに関連付けられていないパーティションが割り当てられたワークユニットごとの再配置キューと、ワークユニットが記憶されている保留キューに関連付けられているパーティションが割り当てられたワークユニットごとの準備完了キューとを含む、前記分離することと、各再配置キュー内の各ワークユニットを、そのワークユニットに割り当てられたパーティションに関連付けられる保留キューに再配置することと、複数の保留キューが空になる前に、パーティショニングルールの新しいパーティショニングルールへの変更に応答して、ワークユニットを再パーティショニングすることとを含む。

10

【0034】

別の態様では、一般に、保留キュー、再配置キュー、準備完了キューを含むワークユニットを記憶するキューを使用して、データを処理するためのコンピューティングシステムは、1又は2以上のノードを含む。各ノードは、複数の保留キューのうちストレージに記憶された少なくとも1つの保留キューにワークユニットを受信するように構成された入力デバイス又はポートと、保留キューからワークユニットを処理するように構成された少なくとも1つのプロセッサとを含む。複数の保留キューのために処理するステップは、保留キューごとに、その保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの一つを、パーティショニングルールに従って決定するステップであって、各パーティションが一意的保留キューに関連付けられる、前記決定するステップと、各保留キューから、ワークユニットを、その保留キューと同じストレージに記憶された少なくとも2つの他のキューの対応する1つに分離するステップであって、2つの他のキューが、ワークユニットが記憶されている保留キューに関連付けられていないパーティションが割り当てられたワークユニットごとの再配置キューと、ワークユニットが記憶されている保留キューに関連付けられているパーティションが割り当てられたワークユニットごとの準備完了キューとを含む、前記分離するステップと、各再配置キュー内の各ワークユニットを、そのワークユニットに割り当てられたパーティションに関連付けられる保留キューに再配置するステップと、複数の保留キューが空になる前に、パーティショニングルールの新しいパーティショニングルールへの変更に応答して、ワークユニットを再パーティショニングするステップとを含む。

20

30

【0035】

別の態様では、一般に、保留キュー、再配置キュー、及び準備完了キューを含むワークユニットを記憶するキューを使用して、データを処理するためのコンピューティングシステムは、コンピューティングシステムの1又は2以上のノードのストレージに記憶された複数の保留キューにワークユニットを受信するための手段と、複数の保留キューからワークユニットを処理するための手段とを含む。処理するステップは、保留キューごとに、その保留キューに記憶された各ワークユニットに割り当てられた複数のパーティションのうちの一つを、パーティショニングルールに従って決定するステップであって、各パーティションが一意的保留キューに関連付けられる、前記決定するステップと、各保留キューから、ワークユニットを、その保留キューと同じストレージに記憶された少なくとも2つの他のキューの対応する1つに分離するステップであって、2つの他のキューが、ワークユニットが記憶されている保留キューに関連付けられていないパーティションが割り当てられたワークユニットごとの再配置キューと、ワークユニットが記憶されている保留キュー

40

50

に関連付けられているパーティションが割り当てられたワークユニットごとの準備完了キューを含む、前記分離するステップと、各再配置キュー内の各ワークユニットを、そのワークユニットに割り当てられたパーティションに関連付けられる保留キューに再配置するステップと、複数の保留キューが空になる前に、パーティショニングルールの新しいパーティショニングルールへの変更に応答して、ワークユニットを再パーティショニングするステップを含む。

【0036】

異なるキューの名前（例えば、「保留キュー」、「再配置キュー」、「準備完了キュー」）、又は異なるサブキュー（例えば、「保留中のサブキュー」、「再循環サブキュー」）は、純粋にラベル付け目的にすぎず、それらの名前には技術輸入が全くない点に留意されたい。異なる名前を有する異なるキュー又はサブキュー（本明細書では、異なる「キューのタイプ」とも呼ばれる）は、同じ技術を使用して実装されてもよく、本明細書に明示的に記載されているように、それらが使用される方法に関してのみ異なってもよい。異なる方法で使用されるキュー又はサブキューは、「第1の」、「第2の」、「第3の」（又は、同様のもの）を使用してラベル付けされてもよいが、純粋に読者の便宜のために、より記述的な言葉を使用してラベル付けされていてもよい。

10

【0037】

態様は、以下の利点のうちの1又は2以上を含むことができる。

【0038】

いくつかのデータ処理システムは、「ワークスティーリング」などの技法を使用して、サーバクラスタの異なるサーバ間で計算のワークユニットを分散する。例えば、各サーバ上のワーカーは、対応するワークキューからワークユニットを消費し得、空のワークキューを有するアイドル状態のサーバは、そのワークキュー内のワークユニットのバックログが大きい別のサーバからワークユニットを取得することができる。これは、負荷バランスを容易にする一種の動的な並列処理の形式を提供するが、部分的にはワークユニットを要求するメッセージを通信するためのコストのために、及び部分的にはデータのローカル性を管理することの複雑さ（すなわち、一部のワークユニットは別のサーバによって盗まれるべきではなく、ローカルサーバに残すべきである）のために、ワークスティーリング(work stealing)のスケラビリティ(scalability)が制限されることがある。

20

【0039】

ローカル性を実現する1つの方法は、特定のワークユニットを処理する必要があるタスクを、そのタスクが生成されたサーバから、その特定のワークユニットがローカルに記憶されたサーバに移行させることである。しかし、そのような移行は、動的ではあるが、タスクに関連付けられる状態をシリアル化して転送する必要があるため、コストが高くなる可能性がある。場合によっては、特定のワークユニットが適切なサーバに送信され、ワーカーがその特定のワークユニットで実施されるべきタスクを生成する方が効率的な場合もある。さらに、短命のタスクの場合、アクティブなタスクの移行は、そのタスクの全寿命のかなりの部分を占める可能性がある。

30

【0040】

異なるサーバ（又は、他のタイプのコンピューティングノード）間でのワークユニットの明示的なパーティショニングは、適切なパーティショニングルールを使用してローカル性の管理を容易にすることができる。場合によっては、ローカル性が重要ではない場合もあり、その場合、サーバ間でワークユニットを均等に分散するラウンドロビン(round-robin)パーティショニングルールなどの、コンテンツに依存しないパーティショニングが使用され得る。しかし、他の場合、ワークユニットは、レコードのキーフィールド内のキー値に基づいてパーティションを割り当てるパーティショニングルールなどのコンテンツに依存する方法でサーバに分散される必要がある場合がある。

40

【0041】

本明細書に記載される技法は、パーティショニングルールを動的に（すなわち、ワークユニットのグループのパーティショニング中に）変更することを可能にする。したがって

50

、データ処理システム内の複雑さ又は効率の損失に過度の負担をかけることなしに、ワークユニットの動的なパーティショニングの利点が達成され得る。例えば、以下でより詳細に説明するように、パーティショニングの異なる態様に複数のキューがあるため、新しいパーティショニングルールに従って再度処理する必要があり得る任意のワークユニットを決定する方がより簡単である。また、必要に応じて特定のワークユニット間の相対的順序付けの保存を容易にするパーティショニング手順の態様もある。また、データ処理システムは、様々なコンピューティングノードにおいてパーティショニングを追加又は除去するためにパーティショニングルールを変更することによって、必要に応じてコンピューティングリソースを増減させることもできる。ワークスティーリングに関連付けられる負荷バランスの利点の一部は、パーティショニングルールの適切な管理によっても達成され得る。

10

【 0 0 4 2 】

本発明の他の特徴及び利点は、以下の説明及び特許請求の範囲から明らかになるであろう。

【 図面の簡単な説明 】

【 0 0 4 3 】

【 図 1 】 データ処理システムのブロック図である。

【 図 2 】 データ処理グラフの図である。

【 図 3 A 】 複数のキュー構造のセットの図である。

【 図 3 B 】 複数のキュー構造の代替セットの図である。

20

【 図 4 A 】 ワークユニットの接続とパスの例を示す、複数のキュー構造のセットの図である。

【 図 4 B 】 ワークユニットの接続とパスの例を示す、複数のキュー構造のセットの図である。

【 図 5 A 】 データ処理グラフの図である。

【 図 5 B 】 データ処理グラフの図である。

【 発明を実施するための形態 】

【 0 0 4 4 】

図 1 は、動的なパーティショニング技法が使用され得るデータ処理システム 1 0 0 の例を示す。システム 1 0 0 は、ネットワーク 1 0 4 を介して通信する複数のコンピューティングノード 1 0 2 (例えば、ローカルエリアネットワーク (LAN, local area network) 及び/又はワイドエリアネットワーク (WAN, wide-area network) を介して通信するサーバ、或いはマルチコアプロセッサ内のネットワークを介して通信するコア、或いはそれらの組合せ) を含む。システム 1 0 0 はまた、コンピューティングノード 1 0 2 に加えて、データベース又は他のストレージシステムなどのデータのソース又はリポジトリ、或いはオンラインデータストリームへの接続を表す 1 又は 2 以上のデータストレージノード 1 0 6 を含む、他のタイプのノードを含むことができる。コンピューティングノード 1 0 2 は、例えば、プロセッサ 1 1 0 (例えば、中央処理装置 (CPU, central processing unit) 又はマルチコアプロセッサ)、ワーキングメモリ 1 1 2 (例えば、揮発性ストレージ媒体)、及び永続ストレージ 1 1 4 (例えば、不揮発性ストレージ媒体) を有するサーバコンピュータとして実装され得る。各コンピューティングノード 1 0 2 は、ノード 1 0 2 に記憶されたワークユニットを処理するためにワーカーを実行するために、UNIX オペレーティングシステムのバージョンなどの適切なオペレーティングシステムの制御下で実行環境をホストするように構成され得る。システム 1 0 0 内の別のノードは、プログラミングのための、或いはコンピューティングノード 1 0 2 の一部又は全部及びネットワーク 1 0 4 の少なくとも一部を備えるサーバクラスタと相互作用するためのユーザインターフェースを提供するアクセスノード 1 1 8 であり得る。いくつかのコンピューティングノード 1 0 2 は、サーバクラスタの一部でなくてもよいが、他のコンピューティングノード 1 0 2 の機能のサブセット (例えば、データを記憶するが、ワーカーを実行しないノード) を提供してもよい。

30

40

50

【 0 0 4 5 】

アクセスノード 1 1 8 は、開発者が、データ処理システム 1 0 0 によって実行されるべきデータ処理プログラムを特定することができる開発環境を提供し得る。いくつかの実施形態では、データ処理プログラムは、頂点間の有向リンク（ワークエレメント、すなわちデータの流れを表す）によって接続された頂点（データ処理コンポーネント又はデータセットを表す）を含むデータフローグラフである。例えば、そのような環境は、参照により本明細書に組み込まれる、「Managing Parameters for Graph-Based Data processing graphs」（グラフに基づくデータ処理グラフのためのパラメータ管理）と題された米国特許出願公開第 2 0 0 7 / 0 0 1 1 6 6 8 号明細書により詳細に記載されている。そのようなグラフベースの計算を実行するためのシステムは、参照により本明細書に組み込まれる、
10 「EXECUTING COMPUTATIONS EXPRESSED AS GRAPHS」と題された米国特許第 5 9 6 6 0 7 2 号明細書に記載されている。このシステムに従って作成されたデータフローグラフは、グラフコンポーネントによって表される個々のプロセスに出入りする情報を取得することと、プロセス間で情報を移動することと、プロセスの実行順序を定義することとを行うための方法を提供する。このシステムは、あらゆる使用可能な方法からプロセス間通信方法を選択するアルゴリズムを含む（例えば、グラフのリンクによる通信パスが、TCP/IP 又は UNIX ドメインソケットを使用することもでき、プロセス間でデータを渡すために共有メモリを使用することもできる）。

【 0 0 4 6 】

いくつかの実施形態では、データ処理プログラムは、データフローに加えて、参照により本明細書に組み込まれる、「COMPILING GRAPH-BASED PROGRAM SPECIFICATIONS」（グラフに基づくプログラム仕様のコンパILING）と題された米国特許出願公開第 2 0 1 6 / 0 0 7 0 7 2 9 号明細書に記載される特徴を含む、データ処理グラフである。例えば、この公開は、データ処理グラフが、データの転送又は制御の転送のいずれか、或いはその両方を表すノード間のリンクを含み得ることを説明している。リンクの特性を示す 1 つの方法は、コンポーネントに異なるタイプのポートを提供することである。リンクは、上流コンポーネントの出力ポートから下流コンポーネントの入力ポートに結合された有向リンクである。ポートには、データエレメントがリンクからどのように書き込まれ、及び読み取られるか、並びに / 又は、データを処理するためにコンポーネントがどのように制御されるかの特性を表すインジケータがある。
20

【 0 0 4 7 】

これらのポートは、いくつかの異なる特性を有し得、そのうちのいくつかは、コンパイラ及び / 又は実行時(runtime) コンピューティングシステムが、並列化を達成するために、どのようにワークユニットをパーティショニングするかに影響を及ぼす。ポートの 1 つの特徴は、入力ポート又は出力ポートとしてのその方向性である。有向リンクは、上流コンポーネントの出力ポートから下流コンポーネントの入力ポートに伝達されるデータ及び / 又は制御を表す。ポートのもう 1 つの特徴は、そのタイプである。開発者は、異なるタイプのポートをリンクすることが許容されている。データ処理グラフのデータ処理特性のいくつかは、異なるタイプのポートがどのようにリンクされているかに依存する。例えば、異なるタイプのポート間のリンクは、各実行が、その実行セット内のコンポーネントの複数のインスタンスにパーティショニングされ得るワークユニットの上流集合によって供給される階層型の並列性を提供する、異なる「実行セット」内のコンポーネントのネストされたサブセットにつながることを可能にする。コンポーネントが有し得る異なるタイプのポートは、次のものを含む。
30

・コンポーネントのインスタンスが、それぞれ、ポートに接続されたリンクを通過する収集のすべてのワークユニットを読み書きすることを意味する、収集の入力ポート又は出力ポート。それらの収集ポート間に単一のリンクを有するコンポーネントのペアの場合、下流コンポーネントは、上流コンポーネントによって書き込まれるときにワークユニットを読み取ることが一般に許容され、上流コンポーネントと下流コンポーネントとの間のパイプラインの並列性が可能になる。また、ワークユニットは並べ替えが行われてもよく、
40
50

これにより並列化の効率化が可能になる。いくつかのグラフ表示では、例えば図 2、図 5 A、及び図 5 B に示されるデータ処理グラフにおいて、そのような収集ポートは、一般に、コンポーネントにおいて四角形のコネクタ記号で示されている。

・コンポーネントのインスタンスが、ポートに接続されたリンクから、又はそこへ、最大で 1 つのワークユニットにおいてそれぞれ読み書きを行うことを意味する、スカラ入力ポート又は出力ポート。それらのスカラポート間に単一のリンクを有するコンポーネントのペアの場合、上流コンポーネントの実行が終了した後の下流コンポーネントのシリアル実行は、単一のワークユニットの転送を制御の転送として使用して強制される。いくつかのグラフ表示では、例えば図 2、図 5 A、及び図 5 B に示されるデータ処理グラフにおいて、そのようなスカラポートは、一般に、コンポーネントにおいて三角形のコネクタ記号

10

で示されている。
・スカラ入力又は出力と類似しているが、ワークユニットは送信される必要がなく、コンポーネント間で制御の転送を通信するために使用される、制御入力又は出力ポート。それらの制御ポート間のリンクを有するコンポーネントのペアの場合、上流コンポーネントの実行が終了した後の下流コンポーネントのシリアル実行が強制される（たとえ、それらのコンポーネントが収集ポート間のリンクも有する場合であっても）。いくつかの図形表現では、そのような制御ポートは、一般に、コンポーネントにおいて円形のコネクタ記号で示されている。

【 0 0 4 8 】

これらの異なるタイプのポートは、データ処理グラフの柔軟な設計を可能にし、ポートタイプの重複するプロパティを使用して、データと制御フローとの強力な組合せを可能にする。特に、データを何らかの形で伝達する収集ポートとスカラポートの 2 つのタイプのポートがあり（「データポート」と呼ばれる）、シリアル実行を強制するスカラポートと制御ポートの 2 つのタイプのポートがある（「シリアルポート」と呼ばれる）。データ処理グラフは、一般に、「ソースコンポーネント」（接続された入力データポートを有さない場合がある）である 1 又は 2 以上のコンポーネント、及び「シンクコンポーネント」（接続された出力データポートを有さない場合がある）である 1 又は 2 以上のコンポーネントを有する。

20

【 0 0 4 9 】

異なるタイプのポート間の接続も、異なる実行セットを示す。特に、収集タイプ出力ポートからスカラタイプ入力ポートへのリンクは、実行セットへの入口点を意味し、スカラタイプ出力ポートから収集タイプ入力ポートへのリンクは、実行セットの出口点を意味する。コンパイラは、コンポーネントを実行セットに自動的に割り当てるために、発見アルゴリズムを実施することができ、コンピューティングシステムは、本明細書に記載される動的なパーティショニング技法を使用してワークユニットの集合からワークユニットを処理するために、実行セットを動的に並列化することができる。

30

【 0 0 5 0 】

図 2 は、プログラミンググラフィカルインターフェース内のデータ処理グラフ 2 0 0 の例を示す。データ処理グラフ 2 0 0 は、第 1 のデータセット（D 1）2 0 2、第 1 のコンポーネント（C 1）2 0 4、第 2 のコンポーネント（C 2）2 0 6、第 3 のコンポーネント（C 3）2 0 8、第 4 のコンポーネント（C 4）2 1 0、第 5 のコンポーネント（C 5）2 1 2、第 6 のコンポーネント（C 6）2 1 4、及び第 2 のデータセット（D 2）2 1 6 を含む。コンパイラは、（ルート、「レベル 0」実行セット 2 2 2 に加えて）「レベル 1」実行セット 2 1 8 及び「レベル 2」実行セット 2 2 0 の発見をもたらす実行セットディスカバリ手順を使用して、データ処理グラフ 2 0 0 の個々のコンポーネント及びデータセットを実行セットに割り当てることができる。ルート実行セット 2 2 2 以外の各実行セットは、ワークユニットの上流集合によって供給され、各入来ワークユニット上で動作する実行セットの異なるインスタンスに対応するワークユニットは、実行時にコンピューティングノード 1 0 2 の間でパーティショニングされ得る。そのようなパーティショニングを実施する能力は、コンパイラが実行可能コードに挿入する「for - all」プリミテ

40

50

ィブを使用してコンパイラによって示されてよく、これは実行時に複数のキュー構造 (MQS, multiple queue structure) と呼ばれる複数のキューを有するバッファデータ構造の作成をトリガし、以下でより詳細に説明される。

【0051】

実行時に、for-allプリミティブは、「サーバグループ」と呼ばれるコンピューティングノード102のグループに分散されるワーカーのプールによって実施されるfor-all操作をトリガする(各コンピューティングノード102は、ワーカーのプールを管理する「サーバ」である)。実行セットによって処理されるべきワークユニットの集合内のすべてのワークユニットの場合、その実行セット内のコンポーネントのグループのインスタンスが、サーバグループ内のサーバのうちの1つの上で実行される。サーバグループ内の各サーバには、対応するパーティショナによって管理されるMQSがあり、そのサーバグループ用のMQSのグループをともに形成する。サーバグループ内のサーバは、for-all操作の期間にわたって必ずしも一定ではない(例えば、サーバはサーバグループに動的に追加、又はそこから除去され得る)。for-all操作は、サーバグループ内の変更を認識する。サーバが追加されると、割り振られたMQSからワークを消費し、MQSを管理するすべてのパーティショナによって使用されるパーティショニング機能を変更するために、新しいサーバ上のワーカープールが始動する。サーバが除去された場合、本明細書に記載される方法に従って、そのワーカープールは削除され、そのワークは再配置されるか排出される。

10

【0052】

MQSは、for-all操作以外の他のタイプの操作のバッファとしても使用され得る。例えば、「ロールアップ」や「結合」などのデータベース操作は、操作の結果を累積するために追加のステップを含むことによって、for-all操作として実装され得る。また、上流とは異なる数のコンピューティングノード102の下流上で実行するデータ処理プログラムの上流部分と下流部分との間でMQSが使用され得る。例えば、コンピューティングノード102がマルチコアプロセッサ上のコアであり、下流プロセッサが上流プロセッサの2倍のコアを有する場合、1つの上流スレッドから2つの下流スレッドにワークユニットを渡すためにMQSのグループが使用され得る。メモリ階層に従ってインテリジェントにパーティショニングされている場合、1つのコアで処理されるすべてのデータは、同じローカルキャッシュを共有する他のコアに移動することができる。

20

30

【0053】

各MQSは、(実行前にコンパイラによって)静的に決定され、(実行時に)動的に決定され、及び/又はワークユニットが処理されるときに動的に変更され得るパーティショニングルールで構成されたパーティショナによって管理される。例えば、コンパイラは、ほとんどのパーティショニングを使用しないデフォルトのパーティショニングルール(又は、1つのローカルパーティショニングのみを使用する、この場合、パーティショニングが実施される必要はない)を示す場合があり、次いで、実行時に実際のワークロードがあるしきい値より大きくなると、パーティションの数を増やすためにある時点で動的な決定が行われ得る。そのようなパーティショニングルールの変更後、以下でより詳細に説明されるように、MQSからワークユニットを消費するワーカーによって処理されるまで、任意のワークユニットが再パーティショニングされ得る(例えば、パーティショニングを追加又は除去するために、或いは別の理由のために)。

40

【0054】

所与のワークユニットのパーティションを計算するために、様々なパーティショニングルールのいずれかが使用され得る。場合によっては、パーティションはワークユニットのキー値に基づいており、ハッシュ関数などの関数がキー値に適用される。パーティションの数を法とするハッシュ関数の結果が使用され得る。他のオプションは範囲パーティショニングを含み、キー値の範囲が異なるパーティションに割り当てられる。場合によっては、並列化の効率を高めるために、固有ではないキーに対して統計的な範囲パーティショニングが使用され得、より多数のワークユニットに属する可能性が高いことが知られている

50

「重い」キーは、単独のパーティションに置かれるか、又はより少ない数のワークユニットに属する可能性が高いと知られている「軽い」キーのみに配置される。例えば、キーが市外局番に対応する場合、大都市の市外局番が異なるパーティションに割り当てられてもよく、最も大きな都市の市外局番が、最も小さい町の市外局番と同じパーティションに意図的に割り当てられてもよい。

【0055】

データ処理プログラムによって指定されたタスクを実行するワーカーによって消費又は生産されるワークユニット（おそらくはワークユニットの一部）に加えて、場合によっては、データ処理プログラムの実行時間を超えて持続する可能性のある他のより長命のデータもある。そのようなより長命のデータの「参照ユニット」は、例えば、永続ストレージ 114 からコンピューティングノード 102 上で実行するワーカー、又はデータストレージノード 106 からアクセスされ得る。いくつかの実施形態では、参照ユニットは、共有変数ごとに一意のキー値を使用してアクセス（すなわち、読取り又は書込み）され得るキー付き共有変数の形式で、すべてのコンピューティングノード 102 にアクセス可能な共有データとして提供される。これらのキー付き共有変数にアクセスする際のデータ競合を回避するために、ロックなどの適切なアクセス制御が使用され得る。参照ユニットはまた、コンピューティングノード 102 にアクセス可能なデータセット（例えば、ファイル又はデータベーステーブル）内に記憶され得る。場合によっては、参照ユニットは、システム 100 の異なるノード上に記憶された複数のパーティションにわたってパーティショニングされ得る。

【0056】

データ処理システム 100 は、パーティションのうちの 1 つ及び対応する MQS をそれぞれ管理する複数のパーティションを使用してワークユニットを（例えば、for-all プリミティブを実装するために）パーティショニングするように構成され得、ワークユニットは、コンピューティングノード 102 のワーキングメモリ 112 又は永続ストレージ 114、或いは両方の組合せに記憶され得る（例えば、キャッシングスキームを使用して）。図 3A を参照すると、MQS 300A、300B、300C、300D のセットの例は、各 MQS を構成する複数のキューが互いにどのように関係しているか、及び他の MQS のキューにどのように関係しているかを示す。MQS 内の各キューは「先入れ先出し（FIFO, first in, first out）」プロパティを有し、したがって、ワークユニットは、一方の端（バックエンド又は「末尾」エンド）に追加され、また他方の端（フロントエンド又は「先頭」エンド）から除去される。パーティションの数 N はパーティショニングルールによって決定され、パーティションのそれぞれについて同じである。この例では N = 4 であり、A、B、C、及び D とラベル付けされた異なるパーティションを有する。典型的には、コンピューティングノード 102 ごとに 1 つのパーティションがあり、その場合、各 MQS は異なるコンピューティングノード 102 に記憶される。しかしながら、単一のコンピューティングノード 102 に、例えばそのコンピューティングノード 102 上の異なるプロセス及び/又は異なるコアに関連付けられ得る複数のパーティションを割り当てることも可能である。その場合、コンピューティングノード 102 は、複数の MQS を記憶することができる。

【0057】

各 MQS は、パーティショニングされるべき入来ワークユニットを受信する保留キューを含む。したがって、この例では、MQS 300A、300B、300C、300D ごとに、それぞれ 1 つの保留キュー 302A、302B、302C、302D（したがって、各パーティション A、B、C、D に関連付けられる一意の保留キュー）が存在する。これらのワークユニットは、コンピューティングノード 102 のサブセット上にある 1 又は 2 以上のソースから到着している可能性があり、したがって、ワークユニットの保留キューへの最初の到着は、保留キューのすべてより少ない数のデータを取り込み、及び/又は保留キューを不均等に満たす可能性がある。パーティションによって使用されるパーティショニングルールの目標は、保留キューに到着するワークユニットを特定の方法で再分散す

10

20

30

40

50

ることである。

【0058】

プロセッサ110上で実行されるパーティショナは、所与のMQSのワークユニットを処理する。具体的には、MQSごとのパーティショナは、パーティショニングルールに従って、N個のパーティショニングのうちいずれがそのMQSの保留キューに記憶された各ワークユニットに割り当てられているかを決定する。パーティショニングルールは、「コンテンツに依存する」方法であって、特定のワークユニットに割り当てられたパーティションが、ワークユニットに含まれる少なくとも1つの値に基づいて決定される方法で、又は、「コンテンツに依存しない」方法であって、特定のワークユニットに割り当てられたパーティションが、ワークユニットに含まれる任意の値から独立して決定される方法で、ワークユニットにパーティションを割り当ててよい。コンテンツに依存するパーティショニングルールは、例えば、ワークユニット内のキー値（又は、キー値に適用されるハッシュ関数）を使用し得る。これらのキー値は、ワークユニットのプライマリキー、ワークユニットに割り当てられた明示的なパーティショニングキー、同じシリアル化キーが割り当てられたワークユニットをそれらの到着順にシリアルに処理されなければならないシリアル化キー、特定の共有変数にアクセスするために使用される共有変数キー、或いはロールアップキー又は結合キーなどのデータベース操作のためのキーなどの、様々なタイプのキー値のいずれかであり得る。コンテンツに依存しないパーティショニングルールは、例えば、ラウンドロビン分布、擬似ランダム分布、又は、MQS内のワークユニットのバックログがより小さいパーティショニングを優先する負荷バランス分布、又はローカルコンピューティングノード102上のパーティショニングにワークユニットを維持することを優先する分布などの、ワークユニットコンテンツ以外の要因に依存する分布を使用し得る。

10

20

【0059】

特定のワークユニットに割り当てられたパーティションの決定に基づいて、パーティショナは、各保留キューからワークユニットを、それぞれのMQS内の2つの他のタイプのキューのうちの対応する1つに分離する。キューの1つのタイプは、その対応する保留キューと同じストレージにそれぞれ記憶された再配置キュー304A、304B、304C、304Dのセット（それぞれ保留キュー302A、302B、302C、302Dに対応する）である。ワークユニットが記憶されている保留キューに関連付けられていないパーティションが割り当てられたワークユニットについては、パーティショナは、ワークユニットをその保留キューから対応する再配置キューに移動する。別のタイプのキューは、対応する保留キューと同じストレージにそれぞれ記憶された準備完了キューのセット306A、306B、306C、306D（それぞれ保留キュー302A、302B、302C、302Dに対応する）である。ワークユニットが記憶されている保留キューに関連付けられているパーティションが割り当てられたワークユニットについては、パーティショナは、ワークユニットをその保留キューから対応する準備完了キューに移動する。準備完了キュー306A、306B、306C、306Dごとに、その準備完了キューからワークユニットを消費する、対応するワーカプール308A、308B、308C、308Dが存在する。

30

40

【0060】

所与のMQSの再配置キューも、そのMQSのパーティショナによって処理される。再配置キューがいっぱいになると、パーティショナは、各再配置キュー内の各ワークユニットを、そのワークユニットに割り当てられたパーティションに関連付けられる保留キューに再配置する。したがって、図3Aに示されるように、再配置キュー304A内のワークユニットは保留キュー302B、302C、302Dのいずれかに送られ、再配置キュー304B内のワークユニットは保留キュー302A、302C、302Dのいずれかに送られ、再配置キュー304C内のワークユニットは、保留キュー302A、302B、302Dのいずれかに送られ、再配置キュー304D内のワークユニットは、保留キュー302A、302B、302Cのいずれかに送られる。MQSが異なるコンピューティング

50

ノード102に記憶されている場合、再配置キュー（例えば、再配置キューA）から特定のパーティション（例えば、パーティションB）に割り当てられた1又は2以上のワークユニットを再配置することは、適切なコンピューティングノード102の間のネットワーク104を介して、1又は2以上のワークユニットを含むメッセージを送信することを含む。ワークユニットは個別に送信されてよく、又は効率のために、パーティショナ（又はパーティショナとは別のバックグラウンドプロセス）が、同じパーティション/MQSに移動するワークユニットのバッチで、同じメッセージにバンドルされたワークユニットを定期的に（例えば、100ミリ秒ごとに）送信する（つまり、メッセージの総数を減らす）ことができる。或いは、MQSが同じコンピューティングノード102に記憶されている場合、1又は2以上のワークユニットを再配置することは、それらのワークユニット（又は、それらのワークユニットへのポインタ）をワーキングメモリ112のある部分から別の部分に単に移動することを伴い得る。また、いくつかの実施形態では、各MQS内の単一の再配置キューの代わりに、潜在的な宛先パーティション（すなわち、N-1再配置キュー）ごとに1つずつ、複数の別個の再配置キューが存在し得る。

10

【0061】

準備完了キューは、MQSを記憶しているコンピューティングノードのワーカーによって処理される。準備完了キューがいっぱいになると、その準備完了キューが記憶されている同じコンピューティングノード102上の1又は2以上のワーカーは、その準備完了キューからワークユニットを消費する。準備完了キューからワークユニットを除去した後にワークユニットを「消費」するためにワーカーが実施するタスクは、MQSの割振りを開始したデータ処理プログラムに依存し、場合によってはワークユニット自体にも依存する。タスクの1つの例は、データ処理グラフから発見された実行セットによって表される部分グラフである。タスクの別の例は、消費されたワークユニットを他のワークユニットと結合する結合操作などのデータベース操作、又は各ワークユニットが消費されると累積された結果を更新することによってワークユニットを集計するロールアップ操作である。ロールアップ操作の場合、累積された結果は部分集合の配列などの参照ユニットとして記憶され得、ワークユニットが消費されると、配列内の新しい項目として追加されるか（キーが以前に表示されていない場合）、部分集計のうちの1つを更新するために使用される。

20

【0062】

いくつかの実施形態では、各保留キューは2つの部分に分割されるか、又は均等に、2つの別個のキューとして維持される。「保留中のサブキュー」と呼ばれる保留キューの第1の部分は、初期パーティション決定がまだ行われていない（すなわち、異なるMQSから再配置されていないワークユニットに対して）ワークユニットを記憶する。「再循環サブキュー」と呼ばれる保留キューの第2の部分は、それらのワークユニットをその保留キューに再配置させる初期パーティション決定が行われたワークユニットを記憶する。ワークユニットが各保留キューから対応する再配置キュー又は準備完了キューに分離されると、ワークユニットは最初に再循環サブキューに記憶されたものから取り出され、再循環サブキューが空である場合にのみ、保留中のサブキューに記憶されたサブキューから取り出される。

30

【0063】

いくつかの実施形態では、各保留キューは3つの部分に分割されるか、又は均等に、3つの別個のキューとして維持される。保留中のサブキュー及び再循環サブキューに加えて、「戻りサブキュー」と呼ばれる保留キューの第3の部分は、パーティショニングルールの変更後にその保留キューに戻された任意のワークユニットを記憶する。パーティショニングルールが変更されると、ワークユニットは、次の順序で、新しいパーティショニングルールに従って保留キューのこれらの部分から（準備完了キュー又は再配置キューに）除去される：最初に戻りサブキュー、次いで再循環サブキュー、次いで保留中のサブキュー。したがって、ワークユニットは、戻りサブキューが空の場合にのみ再循環サブキューから取り出され、再循環サブキューと戻りサブキューの両方が空の場合にのみ、ワークユニットは保留中のサブキューから取り出される。

40

50

【 0 0 6 4 】

上述のように、保留キューを2つ又は3つの部分に分割することは、ワークユニットがそれらのコンテンツに従って「正しく」再パーティショニングされることを保証するために役立つ。しかしながら、コンテンツに依存しないパーティショニングの場合、再パーティショニングのためにワークユニットを戻りサブキューに戻す必要はない可能性がある。

【 0 0 6 5 】

ワークユニット間の相対的順序付けが維持されなければならない場合もあり、そのような場合には、パーティショニングルールを変更する際に部分的順序付けを維持するための追加の技法がある。例えば、維持され得る1つの部分的順序付けプロパティは、常に同じMQSに別のワークユニットとして存在していた任意のワークユニットが、そのワークユニットに対して相対的順序で存在することである。ワークユニット間の相対的順序付けが重要な場合のために、「パーティショニングエボック」に基づく例示的な手順を以下でより詳細に説明する。

【 0 0 6 6 】

MQSのキューはまた、複数の部分に分割されてもよく、又は他の理由により別々のキューに分割されてもよい。例えば、各キューは、MQSの準備完了キューからワークユニットを消費するワーカーごとに異なるサブキューに分割され得る。複数のサブキューが存在する可能性があるもう1つの理由は、動的な負荷バランスを容易にするためである。パーティショナは、ワークユニットが送信され得る他のMQSごとに異なるバージョンの再配置キューである別々のサブキューを維持することができる。次いで、それらの再配置サブキューの長さは、各MQSがどのくらいビジーであるかを示す背圧を感知するために使用され得る。また、パーティショナは、所与のワークユニットが、異なるMQSのパーティションに再配置するために、それ自体のパーティションに残るか、又は再配置サブキューのうち1つに移動する必要があるかどうかを決定するために、それ自体の準備完了キューから背圧を感知することができる。

【 0 0 6 7 】

動的なパーティショニングの場合、パーティショニングルールが変更され得、前のパーティショニングルールの下でのパーティショニングの完了前に(すなわち、すべてのMQSの保留キューが空になる前に)ワークユニットが再パーティショニングされる。パーティショニングルールが変更される際には、パーティショニングの性質とMQSの下流の任意の順序要件の両方に応じて、パーティショナが使用できる様々な手順がそれらのMQSを管理する。一般に、パーティショニングルールの変更に応じてパーティショナがワークユニットを再パーティショニングする方法は、上記で画定されたように、ワークユニットへのパーティショニングの割当てが、コンテンツに依存するかコンテンツに依存しないかによって異なる。コンテンツに依存する場合、再パーティショニングするステップは、再配置キューに記憶された任意のワークユニットと特定の保留キューに対応する準備完了キューをその特定の保留キューに戻すステップを含む。戻されるワークユニットは、上記で画定された戻りサブキューに戻される場合がある。コンテンツに依存しない場合、再パーティショニングするステップはそのステップを実施する必要はない。コンテンツに依存する場合に対して特定の順序要件がある場合、再配置キューを排出する追加のステップもある。コンテンツに依存する、又はコンテンツに依存しないいずれのパーティショニングでも、パーティショニングルールを変更すると、パーティションの数が変更され得る。新しいパーティショニングごとに新しいMQSとパーティショナが追加されるが、除去されている古いパーティションは、そのMQSが任意のワークユニットを記憶しない限りMQSとパーティショナを維持し得る。再パーティショニングを実施するための具体的な手順の詳細な例については、以下でより詳細に説明する。

【 0 0 6 8 】

これらの手順を説明する前に、代替のMQS実施形態の例が図3Bに示されている。この例では、MQS 400A、400B、400C、400Dごとに、それぞれ1つの保留キュー402A、402B、402C、402Dが存在する。図3Aの例では、図3Bの

10

20

30

40

50

例はまた、対応する再配置キュー 404A、404B、404C、404D のセット、対応する準備完了キュー 406A、406B、406C、406D のセット、及びそれらの準備完了キューからワークユニットを消費するワーカプール 408A、408B、408C、408D を含む。しかしながら、この例では、所与のシリアルライゼーションキーを有する単一のワークユニットだけがいつでもワーカーによって消費され得るというシリアルライゼーション制約を強制するために、各準備完了キューは、所与のシリアルライゼーションキーを有する1つのワークユニットしか1度に記憶できない。したがって、ビジューキュー 410A、410B、410C、410D (それぞれ保留キュー 402A、402B、402C、402D に対応する) の追加のセットが存在する。シリアルに実施されるように制約されているすべてのワークユニットがシリアルライゼーションキーの同じ値を有するように、各ワークユニットには、シリアルライゼーションキーが割り当てられているものとする (同じ「キーグループ」にあると言われる)。特定のワークユニットに割り当てられたパーティションを決定した後、パーティションは、各保留キューからワークユニットを、それぞれの MQS 内の3つの他のタイプのキューのうちの対応する1つに分離する。再配置される必要があるワークユニットは、依然として対応する再配置キューに移動される。しかし、再配置される必要がないワークユニットは、その準備完了キューが現在、同じキーグループ内に別のワークユニット (すなわち、同じシリアルライゼーションキーを有する) を記憶していない場合のみ、対応する準備完了キューに移動され、その準備完了キューのワーカプールは、現在、同じキーグループ内の別のワークユニットを処理していない。ワーカプールが複数のワークユニットを同時に処理できるとすると、ワークユニットが準備完了キューに入ることが許容されていれば、ワーカプールは準備完了キューから取り出し、両方のワークユニットを同じキーグループから同時に処理することができるため、ワーカプールが現在同じキーグループからのワークユニットを処理している場合、ワークユニットは準備完了キューに入れられない。パーティションは、例えば、準備完了キュー内の異なるシリアルライゼーションキーを有する、又はワーカプールによって処理されているワークユニットの存在を追跡する各シリアルライゼーションキー (例えば、ビットマップ内) ごとに1又は2以上のプールフラグを使用して、これらの条件のいずれかが真であるかどうかを追跡することができる。準備完了キューが同じキーグループ内に別のワークユニットを記憶している場合、又はワーカプールが同じキーグループ内の別のワークユニットを処理している場合、ワークユニットは対応するビジューキューに移動される。パーティション (又は、パーティションとは別のバックグラウンドプロセス) は、同じキーグループからのワークユニットがワーカプールによって準備完了キューから除去され、ワーカプールが同じキーグループからのワークユニットの処理を終了すると、ワークユニットをビジューキューから準備完了キューに移動する。パーティショニングルールに変更がある場合、ビジューキュー内のワークユニットは、それらの相対的順序を変更せずに、準備完了キュー内のワークユニットとともに戻りサブキューに移動される。

【0069】

場合によっては、パーティションによって行われた決定が単純化され得る。例えば、ラウンドロビン又はランダムパーティショニングルールを使用するコンテンツに依存しないケースでは、パーティションがそのパーティショニングを決定 (又は確認) せずに、保留キューの再循環サブキューから移動されているワークユニットが、準備完了キューに直接移動され得る。これは、最初のラウンドロビン又はランダムな決定が、ワークユニットが保留キューの保留中のサブキューから移動されたときに1度だけ行われる必要があるためである。次いで、ワークユニットが再循環サブキューに入っているとき、パーティショニングルールが変更されていない限り、その決定は依然として正しいと見なされ得る。

【0070】

いくつかの実施形態では、データ処理システム100のノード間でより長命の参照ユニットをパーティショニングするための技法は、より短命のワークユニットをパーティショニングするための技法とは異なる。MQSを管理するために、システム100はワークユニットパーティションを生成し、参照ユニットを管理するために、システム100は、パ

10

20

30

40

50

パーティショニングルールに基づいて参照ユニットをパーティショニングするために適切なノードにおいて参照ユニットパーティショナを生成する（例えば、変数が、そのキーに対するタスクが生成されるコンピューティングノード102上に配置されることを確実にするためにコンパイラによって提供されるキー付き共有変数のキーに基づいて）。個々の参照ユニットは、一般に個々のワークユニットよりも大きい可能性があるため、ノード間で参照ユニットを転送することはより時間がかかる場合がある。しかしながら、ワーカーは、ワークユニットを消費するためにタスクを実施している間に参照ユニットにアクセスする必要があり得る。したがって、参照ユニットの動的な再パーティショニングのための技法は、参照ユニットを非同期的に移行するパーティショニング手順を含み、任意の特定の参照ユニットが古いノード、新しいノード、又は古いノードと新しいノードとの間のトランジットに配置され得る可能性を可能にする、移行中に参照ユニットにアクセスする方法を提供する。

10

【0071】

例えば、パーティショニング手順中にコンピューティングノード102上でパーティショニングされたキー付き共有変数にアクセスするために、データ処理システム100は、どのコンピューティングノード102が特定のキーを有する変数を記憶しているかを決定し、次いで、そのコンピューティングノード102上のその変数にアクセスするためのタスクを生成する。どのコンピューティングノード102が変数を記憶しているかを決定するために、パーティショニング手順は、古いパーティションと新しいパーティションの両方を計算するためにキーに適用される関数（例えば、ハッシュ関数に基づく）を提供する。共有変数にアクセスする前に、古いパーティションを決定し、その変数の古いパーティションの古いストレージ位置をチェックするために、そのキーが使用される。古いストレージ位置に変数が見つからない場合、新しいパーティションを決定し、変数の新しいパーティションの新しいストレージ位置をチェックするために、そのキーが使用される。新しいストレージ位置に変数が見つからない場合、（チェックするための手順がストレージ位置間の通過時間を考慮するために十分に長く待機する限り）存在しないと決定される。パーティショニング手順中に変数がアクセスされると、その変数はアクセスが完了するまで一時的に移行が遅れる場合がある。パーティショニング手順が完了した後、古いパーティションを計算するための関数はもはや必要なくなる。いくつかの実施形態は、複数のパーティション手順が同時に進行されるようにし、その場合、変数の存在をチェックするためのストレージ位置の数は、それが見つかるまで最も古いものから最も新しいものまで増加し、見つからない場合は存在しないと決定される。

20

30

【0072】

参照ユニットのパーティショニング手順の変形は、参照ユニットがすべて同時に移行されない増分パーティショニング手順である。例えば、キー付き共有変数の場合、キーを所定の数のグループ（パーティショニングから独立している）に分割するために使用されるハッシュ関数によって決定されるように、参照ユニットの複数のグループが存在する。次に、グループのうちの1つにあるキーを有する変数だけが1度に移行され、そのグループにない他の任意の変数は、特定のストレージ位置（そのグループがすでに移行されたかどうかに応じて、古いストレージ位置又は新しいストレージ位置のいずれか）にあることが保証される。この方法では、移行されるグループ内のキーを有する変数だけが、変数を見つける前に両方のストレージ位置をチェックする必要がある可能性がある。したがって、増分パーティショニング手順は、移行中に参照ユニットの多くがアクセスされる場合に必要となる全体的な通信を削減する。

40

【0073】

図4Aは、明示的に示されたMQS内の異なるタイプのキュー間の接続を有する2つのMQSのセットの例を示す。MQS420A、420Bごとに、それぞれ1つの保留キュー422A、422Bが存在する。対応する再配置キュー424A、424Bのセット、対応する準備完了キュー426A、426Bのセットも存在する。この例では、再配置キュー424Aと保留キュー422Bとの間の接続、及び再配置キュー424Bと保留キュー

50

ー 4 2 2 A との間の接続が明示的に示されている。準備完了キュー 4 2 6 A、4 2 6 B からワークユニットをそれぞれ消費するワーカプール 4 0 8 A、4 0 8 B も存在する。この例では、各 MQS の保留キューにワークエレメントを提供する入来パーティシヨナ 4 3 2 とともに、ワークエレメントの上流ソース 4 3 0 も示されている。或いは、他の例では、各保留キューは異なるソースからワークエレメントを受信することができ、パーティシヨニングは、入来パーティシヨナ 4 3 2 を必要とせずに、各 MQS に関連付けられる個々のパーティシヨナによって完全に達成され得る。図 4 B は、MQS 4 2 0 A、4 2 0 B のセットの様々なキューを横断するときワークユニット 4 4 0 が取るパスの例を示す。ワークユニット 4 4 0 は、MQS のセットを通して進むにつれて、相対時間 T 1、T 2、T 3、T 4、T 5、T 6 によってラベル付けされる。この例では、ワークユニット 4 4 0 は、(時間 T 2 において) MQS 4 2 0 B のパーティシヨナ B に最初に割り当てられるが、最終的にワーカプール 4 2 8 A によって処理(又は「消費」)される前に(時間 T 6 において)、MQS 4 2 0 A のパーティシヨナ A に動的に再パーティシヨニングされる。

【 0 0 7 4 】

データ処理システム 1 0 0 はまた、データ処理プログラム内で必要とされるバッファリングの種類に基づいて MQS を管理する際にパーティシヨナによって使用されるべきパーティシヨニングルールの種類を選択するように構成される。図 5 A は、内側ループ 5 0 2 及び外側ループ 5 0 4 を含む、ネストされた実行セットを有するデータ処理グラフ 5 0 0 の一例を示し、その各々は、ワークユニットの入来ストリームにわたって反復するループにおいて計算を実施する。コンパイラは、システム 1 0 0 の複数のコンピューティングノード 1 0 2 にわたってループをパーティシヨニングするための MQS のセットにそれぞれ対応するバッファコンポーネントを挿入することができる。図 5 B は、コンパイラが、ネストされた実行セットのためのパーティシヨニングを扱うためにバッファコンポーネントを追加した後に、対応するデータ処理グラフ 5 5 0 が(少なくとも概念的に)どのように現れるかを示す。具体的には、バッファコンポーネント 5 5 2 A、5 5 2 B、5 5 2 C、及び 5 5 2 D の各々は、実行時に割り振られ、コンピューティングノード 1 0 2 を介して分散される MQS のそれぞれのセットを表す。いくつかの実施形態では、コンパイラによって追加されたこれらのバッファコンポーネントはユーザに表示されない。

【 0 0 7 5 】

この例では、MQS が提供する動的なパーティシヨニング機能が、コンパイラと実行時システムに、内側ループ及び外側ループを並列化するための効率的な選択を行うための柔軟性を与える。システム 1 0 0 は、内側ループ及び外側ループのパーティシヨニングが、互換性があり効率的になるように、MQS を処理するパーティシヨナを構成することができる。システム 1 0 0 は、パーティシヨニングルールがコンテンツに依存する必要があるか、又はコンテンツに依存しない必要があるかを考慮に入れることもできる。例えば、コンテンツに依存しないパーティシヨニングルールの場合、内側ループのオプションは、(1) 外側ループパーティシヨンのすべてのワークエレメントを、関連付けられる「ローカル」ノードに維持する(すなわち、再パーティシヨニングしない)こと、又は(2) すべてのノードにわたる外側ループパーティシヨンのローカルノードからのワークエレメントをパーティシヨニングすることを含むことができる。オプション 1 は潜在的によりスケラブルであるが、オプション 2 は、処理されているワークエレメントの数に応じて、潜在的により高速である。外側ループが大きい(すなわち、反復回数が多い)場合、及び内側ループが小さい(すなわち、反復回数が少ない)場合、オプション 1 はオプション 2 より高速である可能性がある。これについての 1 つの説明は、オプション 2 の場合、(小さい)内側ループのワークエレメントをパーティシヨニングするオーバーヘッドコストを回収するのに十分な時間がないということである。一方、外側ループが小さく、内側ループが大きい場合、オプション 1 はオプション 2 よりも高速ではない可能性がある。これについての 1 つの説明は、(小さい)外側ループがコンピューティングノード 1 0 2 のうちのいくつかのみパーティシヨニングすることができ、オプション 2 のパーティシヨニングは、(大きい)内側ループをより高速に実施することを助けるために、残りのコンピューティ

10

20

30

40

50

ングノード102が使用されることを可能にする。

【0076】

動的なパーティショニング機能は、システム100が実行時まで内側ループ及び外側ループの相対的なサイズに関する十分な情報を有していない可能性があるため、適切なオプションを選択するために有用であり得る。コンパイラは、特定の前提条件に基づいて、パーティショナにデフォルトのパーティショニングルールを使用し得る。例えば、内側ループの方が、それらの外側ループに比べて小さくなるという前提条件があり得る。さらに、多数のワークエレメントを生成する可能性のある特定のソース（例えば、過去の実行からの統計に基づいて）、及び/又は生成される可能性が高いワークエレメントの数を推定するために分析され得るソース（例えば、ファイルの場合、それが含むワークユニットの数を推定するためにファイルサイズが使用され得る）が存在する場合がある。次いで、実行時に、システム100は、いつデフォルトのパーティショニングルールが変更されるべきかを動的に決定することができる。例えば、内側ループによって処理されるワークエレメントの数がしきい値（あらかじめ定められたしきい値、又は外側ループに対するしきい値）を超える場合、パーティショニングルールは、本明細書に記載の技法を使用して変更され得る。

10

【0077】

上述のように、パーティショニングルールの変更への対応方法を決定する際には、複数のオプションもある。最良のオプションは、パーティショニングの性質と、MQSのセットで実装されているバッファの下流で発生する処理の順序要件の両方に依存する。例えば、

20

【0078】

オプション1：パーティショニングルールがコンテンツに依存しない場合 - つまり、すべてのワークユニットにあらかじめ定められた正しいパーティションがない場合 - 準備完了キュー又は再配置キューにすでに作成されたワークユニットを再パーティショニングする必要さえない。この例としては、すでにパーティショニングされているワークユニットを再パーティショニングする必要がないラウンドロビン又は負荷バランシングのパーティショニングルールがある。

【0079】

オプション2：上記で説明したこのオプションでは、MQS内のワークユニットが回収され、保留キューに戻される（保留キューの戻りサブキュー部分を使用してもよい）。このオプションは、パーティショニングルールがコンテンツに依存し、下流処理が順序付けられていない場合、最も効率的なオプションである。これは、例えば、ワークユニットがキーによってパーティショニングされているが、処理は同じキーを共有するワークユニットに対して特定の順序付けを必要としなかった場合に発生する。

30

【0080】

オプション3：このオプションは、パーティショニングルールがコンテンツに依存し、下流処理が順序付けされている場合に適用される。これは、例えば、ワークユニットが、キーによってパーティショニングされており、同じキーを有し、同じソースからのワークユニットが順番に処理されなければならない場合に発生する。このオプションでは、ワークユニットを再パーティショニングすることは、ワークユニットが正しいパーティショニングで処理されることを保証するために十分であるが、順序が維持されることを保証するためには不十分である。特定の部分的順序付けを保証するために、「エポック移行手順」と呼ばれる複数ステップの手順が使用され得る。この手順は、特定のMQSのグループの一部であるMQSを管理しているすべてのパーティショナについてパーティショニングルールが同じである期間に対応する「パーティショニングエポック」間の移行のために使用される。隣接するパーティショニングエポックの間には、エポック移行手順が発生する「移行期間」がある。場合によっては、1又は2以上のパーティショナが新しいパーティショニングルールを使用し始めるか、依然として古いパーティショニングルールを使用している移行期間中に、いくつかのパーティショニングが依然として発生している場合がある

40

50

。1つのパーティショニングエポックから次のパーティショニングエポックへの進展を調整することによって、システム100は、以下でより詳細に説明するように、パーティショニングルールの変更によって部分的順序付けを保存することができる。

【0081】

パーティショナがMQSを管理する方法は、移行期間ではなくパーティショニングエポックの間に異なる場合がある。パーティショニングエポックの間、再循環サブキューに入るワークユニットのパーティションは、以前のコンピューティングノードで計算されたパーティションが正しいため、新しいコンピューティングノードにおいて再計算される必要はないが、新しいコンピューティングノードにおいて計算されたパーティションと必ずしも同じではない。一方、移行期間中に、再循環サブキューに入るワークユニットのパーティションが再計算される。例えば、ランダムパーティショニングの場合、同じパーティショニングエポックの間の第2のパーティション決定は、最初にランダムに選択されるパーティションに再配置されたばかりのワークユニットには不適切である。その代わりに、再循環サブキューに入る任意のワークユニットは、単に準備完了キューに直接移動され得る（1つずつ、又は複数のワークユニットのグループにおいて）。

10

【0082】

MQSは、ネットワーク104を介してコンピューティングノード102間にワークユニットを分散する。データ処理プログラムは、ワークユニットをMQSに注入する1又は2以上のソース（例えば、データ処理グラフの処理コンポーネント又はソースコンポーネントに対応する）と、MQSからワークユニットを消費する1又は2以上のワーカプール（例えば、データ処理グラフの処理コンポーネント又はシンクコンポーネントに対応する）を定義する。パーティショナが同じパーティショニングルールを使用するMQSの特定のセットは、コンピューティングノード102の特定のサブネットワーク上でのパーティショニングを管理すると考えられる。このMQSのセットは、例えば、上述したように、データ処理グラフ内のバッファコンポーネントとして表され得る。MQSのセットによって使用されるサブネットワークを介したルーティングは、それらの共通のパーティショニングルール、及び順序制約などのデータ処理プログラムの他の制約によって決定される。MQSの特定のセットは、それらのルーティング制約に基づいて、3つのタイプのうちの1つを有するものとして分類され得る。

20

・タイプ1：最小限の制約付きルーティングにより、ワークユニットが任意のワーカプールによって処理されることが可能になる。データ処理プログラムは、ワークユニットを特定のワーカプール（又は、特定のコンピューティングノード102）にルーティングするためにMQSのセットを必要とせず、同じソースからのワークユニットが互いに順番に処理される必要はない。したがって、MQSを管理するパーティショナは、ラウンドロビンパーティショニング又は負荷バランシングパーティショニングなどのパーティショニングルールを使用することができる。

30

・タイプ2：いくつかのデータ処理プログラムは、ワークユニットから導出されたキー値に基づいて、ワークユニットが特定のワーカプール（又は、特定の計算ノード102）によって処理されることを必要とするが、同じソースからのワークユニットがワーカプールにおいて順番に処理されることを必要としない。MQSを管理するパーティショナは、キーベースのパーティショニングルールを使用する。

40

・タイプ3：最も制約の厳しいルーティングは、ワークユニットが、ワークユニットから導出されたキー値に基づいて、特定のワーカプール（又は、特定のコンピューティングノード102）によって処理されることを必要とし、また、同じキーを有する同じソースからのワークユニットが、それらがソースにおいて生成されたのと同じ順序でワーカプールにおいて処理されることを必要とする。MQSを管理するパーティショナは、キーベースのパーティショニングルールを使用し、順序が保存されることを保証するために追加のステップを実行する。

【0083】

場合によっては、サブネットワークが動作している間にパーティショニングルールを変

50

更する必要がある。これは、ワーカープールをあるコンピューティングノード102から別のコンピューティングノード102に移動させるために、又はワークユニットのキーとそれを処理するワーカープールとの間のマッピングを変更するために、ワーカープールを追加又は除去するために行われ得る。パーティショニングルールが変更されると、特定のパーティションに関連付けられるMQSに記憶されたワークユニットが、新しいパーティショニングルールに従って誤ったパーティションに存在する可能性がある。それらのワークユニットは、それらの新しく正しいパーティションに移動される必要があり、順序を保存するために特定のステップが実行される。

【0084】

パーティショニングエポック内では、どのワークユニットもそのキーを1度だけ調べることによってその宛先にルーティングされ得る。パーティショナは、ワークユニットを（保留キューの）保留中のサブキューからそのMQSの準備完了キュー又は再配置キューのいずれかに移動させるときに、ワークユニットのキーを調べるだけでよい。ワークユニットが、現在のパーティショニングエポックの全体にわたって有効であった同じパーティショニングルールに基づいてこのMQSにルーティングされたことがわかっているため、（保留キューの）再循環サブキューに到着するワークユニットは、それらのキーを再検査せずに、準備完了キューに直接移動され得る。

10

【0085】

1つのワークユニットがソースAからワーカープールPまでのサブネットワークを通るパスをたどる場合、ソースAにおいて発生し、同じパーティショニングエポックの間にワーカープールPによって消費されるすべてのワークユニットは、AからPまでの同じパスに従う。AからPに送信されるすべてのワークユニットは同じパスに従い、パーティショナは移行中のワークユニットの相対的順序を保存するように構成されているため、ソースにおいてワークユニットが生産された順序は、ソースとワーカープールのすべてのペアについて、ワーカープールによって消費された順序で保存される。異なるソースにおいて発生するか、又は異なるワーカープールによって消費される2つのワークユニットは、必然的にサブネットワークを通る異なるパスを取るため、それらのワークユニット間には順序は画定されていない。

20

【0086】

セット内のMQSを管理しているパーティショナのうちの少なくとも1つに対してパーティショニングルールが変更されると、現在のパーティショニングエポックが終了し、新しいパーティショニングエポックへの移行期間が開始される。移行期間中、順序が保存されていることを保証する不変量が侵害されることが可能になる。データ処理プログラムの順序要件及びパーティショニング要件に応じて、新しいパーティショニングエポックに安全に移行するために、異なる手法が必要である。

30

【0087】

単純で低コストの移行戦略は、再循環サブキュー、再配置キュー、又は準備完了キュー内にすでにある任意のワークユニットを再処理するために時間を費やすことなしに、新しいパーティショニングルールを使用するようにパーティショナを単純に更新することである。この戦略は、タイプ1のMQSセットに対してのみ許容される。パーティショニングエポックNの間にキーが検査された、いくつかのワークユニットは、エポックN+1の間に検査された場合にルーティングされたワーカープールとは異なるワーカープールによって、パーティショニングエポックN+1の間に消費されるが、これは、特定のプールへのワークユニットのルーティングは、そのタイプのMQSセットの要件ではないため、タイプ1のMQSセットについて許容される。

40

【0088】

よりコストの高い移行戦略は、移行期間中にワークユニットを再処理することを含む。タイプ2とタイプ3のMQSセットは、現在のパーティショニングルールに従ってワークユニットのキーから決定される特定のワーカープールによってワークユニットが消費されることを必要とする。パーティショニングルールが変更されると（移行期間中に）、ワー

50

クユニットの正しい宛先ワーカールも変更されることがある。したがって、パーティシヨナはワークユニットを再処理し、これは、新しいパーティシヨニングルールに従って、ワークユニットキーを再検査することと、必要に応じてワークユニットをそれらの正しい宛先に潜在的に再ルーティングすることを含む。順序が保存されるか(タイプ3)、否か(タイプ2)によって、移行期間中の再処理を処理するための2つの方法がある。

【0089】

タイプ2のMQSセットの場合、ワークユニットの再検査と潜在的な再ルーティングが必要であるが、順序を保存することは必要ない。この場合、セット内のすべてのMQSのすべての準備完了キュー、再循環サブキュー、及び再配置キュー内のワークユニットを回収し、それらを適切な保留中のサブキューに挿入することによって、ワークユニットが再処理され得る。この再処理の後、次のパーティシヨニングエポックにおいて、新しいパーティシヨニングルールを使用してワークユニットがパーティシヨニングされ、正しい宛先ワーカールにおいて消費される。

10

【0090】

以下は、それぞれのMQS上で各パーティシヨナによって実施される、タイプ2のMQSセットのタイプ2エポック移行手順(ルーティングエポックNからN+1まで)の例である。

- ・あらゆるキュー(又はサブキュー)からワークユニットの処理を中断する。
- ・再循環サブキュー、再配置キュー、及び準備完了キューからワークユニットを保留中のサブキューに回収する。
- ・エポックN+1をパーティシヨニングするためのパーティシヨニングルールに従ってパーティシヨニングを切り替える。
- ・すべてのキュー(及びサブキュー)からワークユニットの処理を再開する。

20

【0091】

パーティシヨニングエポックNの間に送信されたが、パーティシヨニングエポックN+1までに配信されなかったワークユニットは、パーティシヨニングエポック移行の前にそれらがすでに再配置されたか(再配置キューから再循環サブキューに)どうかに応じて、サブネットワークを通る2つの異なるパスに従い得るので、順序付けは、このエポック移行手順において保存されることが保証されていない。ワークユニットが両方のパスに沿って宛先ワーカールに同時に収束すると、それらがインターリーブされ得、この場合、相対的順序が失われることになる。このため、この手順はタイプ2のMQSセットでは受け入れられるが、タイプ3のMQSセットでは受け入れられない。

30

【0092】

2つのMQS - MQS1とMQS2 - のMQSセットの場合、ワークユニットの状態は、すべて同じ特定のキーKを有し、パーティシヨニングエポックNの間は、MQS1を供給するソースSから同じパスをたどるように、パーティシヨン2に割り当てられ、次いでMQS2に再配置され、次いでMQS2においてワーカールによって消費されると見なす。この例では、パーティシヨニングエポックNの終了時のキューのコンテンツは以下の通りであり、ソースSから提供された順序に従ってレコード番号R#によってラベル付けされたワークユニットを有し、ワークユニットは左側の先頭から除去されて、右側の末尾に追加され、空白のリストは空のキュー又はサブキューを示す。

40

```

MQS 1      保留中のサブキュー： R 7、 R 8
           再循環サブキュー：
           再配置キュー： R 6
           準備完了キュー
MQS 2      保留中のサブキュー：
           再循環サブキュー： R 4、 R 5
           再配置キュー：
           準備完了キュー： R 3

```

【0093】

50

次いで、エポック移行手順の後、各MQS内のすべてのワークユニットが回収され、保留中のサブキューに移動され、以下ようになる。

MQS 1 保留中のサブキュー：R 6、R 7、R 8
 再循環サブキュー：
 再配置キュー：
 準備完了キュー

MQS 2 保留中のサブキュー：R 3、R 4、R 5
 再循環サブキュー：
 再配置キュー：
 準備完了キュー

10

【0094】

エポックN+1をパーティショニングするための新しいパーティショニングルールに基づいて、キーKを有するワークユニットは、MQS2におけるワーカプールではなく、MQS1においてワーカプールによって消費されるべきパーティション1に割り当てられ、したがって、パーティショニングエポックN+1の間、ワークユニットはMQS2の再配置キューに移動され、次いで、MQS1の再循環サブキューに送られ、最終的にそのワーカプールによって消費されるためにMQS1の準備完了キューに移動される。それらのワークユニットはすべてMQS1準備完了キューにおいて終了するが、ワークユニットの相対的順序は保証されない。以下は、MQS1及びMQS2のパーティショナがいくつかのワークユニットを処理した後の有効な状態である。

20

MQS 1 保留中のサブキュー：R 8
 再循環サブキュー：R 4
 再配置キュー：
 準備完了キュー：R 6、R 3、R 7（ワークユニットが順番通りでない）

MQS 2 保留中のサブキュー：
 再循環サブキュー：
 再配置キュー：R 5
 準備完了キュー：

【0095】

30

タイプ3のMQSセットは、パーティショニングエポック移行にわたって順序付けが保存されることを必要とする。これを保証するために、エポックNの間に処理されたが消費されなかったすべてのワークユニットが、任意のワークユニットがエポックN+1の間に処理され得る前に、新しい最終宛先において再循環キュー又は準備完了キューに順番に配信される必要がある。

【0096】

以下は、それぞれのMQS上の各パーティショナによって実施されるタイプ3のMQSセットのタイプ3のエポック移行手順（ルーティングエポックNからN+1まで）の例であり、必要な順序保証を実現している。

・準備完了キューからワークユニットを消費することを中断する（依然として新しいワークユニットを受信することができる）。

40

・保留中のサブキューを2つのキューに分割する。エポックN+1の保留中のサブキューは、保留中のサブキューにあったすべてのワークユニットを含み、新しく到着した作業ユニットが追加され得るが、そこからワークユニットは除去されない。エポックNの保留中のサブキューは空になり始めるが、1度回収されたワークユニットがそこに置かれると、再検査され、潜在的に再ルーティングされ得る。

・すべての再配置キューが完全に排出されるまで、エポックNをパーティショニングするためのパーティショニングルールに従って処理を続行する。

【0097】

この時点で、同じキーを有するすべてのワークユニットは、同じMQSの再循環サブキ

50

ュー又は準備完了キューのいずれかにまとめられる。それらは、異なるキーを有するワークユニット、又は異なるソースからの同じキーを有する他のワークユニットと混在されるが、同じソースからの同じキーを有するワークユニットの相対的順序は保存される。

・エポックN + 1をパーティショニングするためのパーティショニングルールに従ってパーティショニングに切り替える。

・ワークユニットを再循環サブキュー及び準備完了キューからエポックNの保留中のサブキューに回収し、それらの元の順序付けを保存する(同じソースにおいて発生したワークユニットに対して)。

・準備完了キューからワークユニットを消費することを再開する。

・すべてのエポックNの保留中のサブキューが完全に排出され、再配置キューが空になるまで待機する。

10

【0098】

この時点では、すべてのワークユニットが新しいパーティショニング機能のためにそれらの正しい宛先にあり、順序はパーティショナによって保存される。パーティショナは、順序を保存するために、回収プロセス中に特定のルールに従うことができる。例えば、各パーティショナが、先頭における準備完了キューからのワークユニット、次に再循環サブキューからのワークユニット、及び末尾における保留中のサブキューにおける元のワークユニット(回収される際に並べ替えられていないワークユニット)、を有するワークユニットをエポックNの保留中のサブキューに戻すと、順序制約が保存され、したがって、回収プロセスが完了した後、同じソースからの、及び同じキーグループ内のワークユニットが依然として順番に並ぶようになる。相互に並べ替えられた可能性のある唯一のワークユニットは、異なる保留中のサブキュー又は再循環サブキューから来たワークユニットであり、それらが元々異なるソースから来たものであるか、異なるキーを有するため異なるMQSに送られた場合にのみ起こりうる。いずれにしても、ワークユニットは、画定された順序制約を有するために同じソースと同じキーの両方を有していなければならないので、その並べ替えが許容される。

20

・空のエポックNの保留中のサブキューを破棄し、パーティショナがパーティショニングエポックN + 1に完全に移行したので、現在は通常の保留中のサブキューであるエポックN + 1の保留中のサブキューからの読取りを再開する。

【0099】

30

以下は、タイプ3のエポック移行手順が、タイプ2のエポック移行手順の例において上記で使用されたものと同じ開始ワークユニット設定を扱う方法であり、ここで、パーティショニングエポックNの終了時のキューのコンテンツは以下の通りである。

MQS 1 保留中のサブキュー：R 7、R 8
再循環サブキュー：
再配置キュー：R 6
準備完了キュー：

MQS 2 保留中のサブキュー：
再循環サブキュー：R 4、R 5
再配置キュー：
準備完了キュー：R 3

40

【0100】

保留中のサブキューを分割した後、コンテンツは次の通りである。

MQS 1 保留中のサブキューN：
保留中のサブキューN + 1：R 7、R 8
再循環サブキュー：
再配置キュー：R 6
準備完了キュー：

MQS 2 保留中のサブキューN：
保留中のサブキューN + 1：

50

再循環サブキュー：R 4、R 5
 再配置キュー：
 準備完了キュー：R 3

【0101】

次いで、パーティショナは、再配置キューが排出されるまで（ワークユニットR 4が再循環サブキューの先頭から準備完了キューの末尾に移動することを許容された状態で）、エポックNをパーティショニングするためのパーティショニングルールで処理を続行する。

MQS 1	保留中のサブキューN： 保留中のサブキューN + 1：R 7、R 8 再循環サブキュー： 再配置キュー： 準備完了キュー：	10
MQS 2	保留中のサブキューN： 保留中のサブキューN + 1： 再循環サブキュー：R 5、R 6 再配置キュー： 準備完了キュー：R 3、R 4	

【0102】

再配置キューが空になった後、再循環サブキュー及び準備完了キューは、エポックNの保留中のサブキューに回収され得る。

MQS 1	保留中のサブキューN： 保留中のサブキューN + 1：R 7、R 8 再循環サブキュー： 再配置キュー： 準備完了キュー：	
MQS 2	保留中のサブキューN：R 3、R 4、R 5、R 6 保留中のサブキューN + 1： 再循環サブキュー： 再配置キュー： 準備完了キュー：	30

【0103】

次いで、パーティショナは、エポックN + 1パーティショニングルールをパーティショニングするためのパーティショニングルールを使用して正常に処理を再開し、完全に排出され、再配置キューが空になるまで、エポックNの保留中のサブキューからワークユニットを除去し、これは最終的に以下の状態になる（ワークユニットR 3 ~ R 5は、再循環サブキューから準備完了キューに移動することができる）。

MQS 1	保留中のサブキューN： 保留中のサブキューN + 1：R 7、R 8 再循環サブキュー：R 6 再配置キュー： 準備完了キュー：R 3、R 4、R 5	40
MQS 2	保留中のサブキューN： 保留中のサブキューN + 1： 再循環サブキュー： 再配置キュー： 準備完了キュー：	

【0104】

すべてのワークユニットが正しい宛先MQSに再配置され、エポックNの保留中のサブキューと再配置キューとが空になった後、パーティショナは空のエポックNの保留中のサ

ブキューを破棄して、エポックN + 1の保留中のサブキューを唯一の保留中のサブキューとして使用して正常に処理を続行することができる。

MQS 1 保留中のサブキュー：R 7、R 8
再循環サブキュー：R 6
再配置キュー：
準備完了キュー：R 3、R 4、R 5

MQS 2 保留中のサブキュー：
再循環サブキュー：
再配置キュー：
準備完了キュー：

10

【0105】

この時点で、パーティションは新しいパーティショニングエポックに正常に移行し、ワークユニットを正しい宛先と正しい順序で有している。

【0106】

上述の動的なパーティショニング手法は、例えば、適切なソフトウェア命令を実行するプログラム可能なコンピューティングシステムを使用して実装されてもよく、フィールドプログラマブルゲートアレイ(FPGA, field-programmable gate array)などの適切なハードウェア又はある種のハイブリッド形式で実装されてもよい。例えば、プログラムされた手法では、ソフトウェアは、1又は2以上のプログラムされた又はプログラム可能なコンピューティングシステム(分散型、クライアント/サーバ型、又はグリッドなどの様々なアーキテクチャのものであり得る)で実行される1又は2以上のコンピュータプログラムにおける手順を含み得、それぞれが、少なくとも1つのプロセッサと、少なくとも1つのデータストレージシステム(揮発性及び/又は不揮発性メモリ並びに/或いはストレージエレメントを含む)と、少なくとも1つのユーザインターフェース(少なくとも1つの入力デバイス又はポートを使用して入力を受信するための、及び少なくとも1つの出力デバイス又はポートを使用して出力を提供するための)とを含む。ソフトウェアは、例えば、データフローグラフの設計、構成、及び実行に関連するサービスを提供する、より大きなプログラムの1又は2以上のモジュールを含み得る。プログラムのモジュール(例えば、データフローグラフのエレメント)は、データ構造又はデータリポジトリに記憶されたデータモデルに適合する他の組織化データとして実装され得る。

20

30

【0107】

ソフトウェアは、媒体の物理的特性(例えば、表面ピット及びランド、磁区、又は電荷)を使用して、一定期間(例えば、ダイナミックRAMなどのダイナミックメモリデバイスのリフレッシュ期間の間の時間)、揮発性又は不揮発性ストレージ媒体、或いは他の任意の非一時的媒体に具体化されるような非一時的形態で記憶され得る。命令のロードの準備において、ソフトウェアは、CD-ROM又は他のコンピュータ可読媒体(例えば、汎用又は専用コンピューティングシステム或いはデバイスによって読取り可能)などの有形の非一時的媒体上に提供されてもよく、ネットワークの通信媒体を介して、それが実行されるコンピューティングシステムの有形の非一時的媒体に配信(例えば、伝搬された信号で符号化)されてもよい。処理の一部又は全部は、専用コンピュータ上で、或いはコプロセッサ又はフィールドプログラマブルゲートアレイ(FPGA)或いは専用のデータ処理グラフ専用集積回路(ASICs)などの専用ハードウェアを使用して実施され得る。処理は、ソフトウェアによって指定された計算の異なる部分が異なるコンピューティングエレメントによって実施される分散方式で実装され得る。そのようなコンピュータプログラムの各々は、本明細書に記載される処理を実施するためにコンピュータによってストレージデバイス媒体が読み取られたときに、コンピュータを構成し、動作させるために、汎用又は専用のプログラム可能コンピュータによってアクセス可能なストレージデバイスのコンピュータ可読ストレージ媒体(例えば、固体メモリ又は媒体、或いは磁気又は光学媒体)に記憶される、又はダウンロードされることが好ましい。本発明のシステムは、コンピュータプログラムで構成された有形の非一時的媒体として実装されると考えられてもよく

40

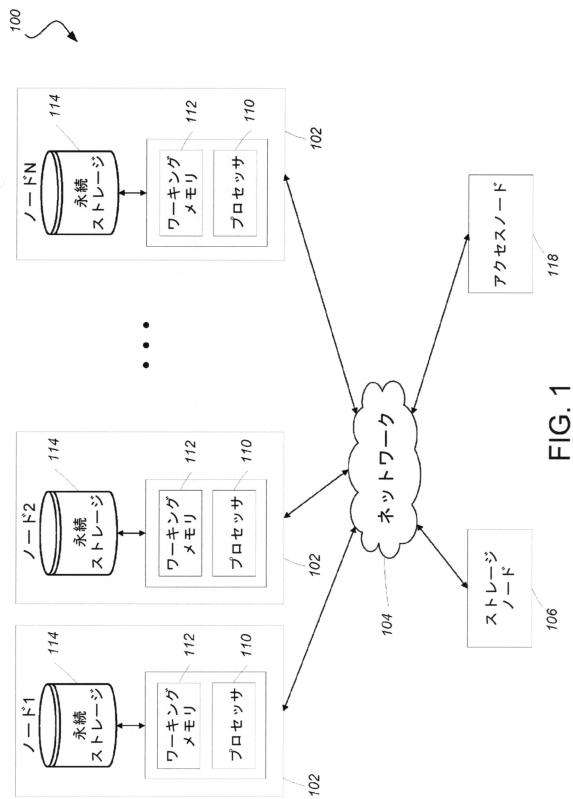
50

、そのように構成された媒体は、コンピュータに本明細書に記載される処理ステップのうち
の1又は2以上を実施するために特定のあらかじめ画定された方法で動作させる。

【0108】

本発明の多くの実施形態について説明した。それにもかかわらず、前述の説明は、本発
明を説明することが意図されるものであり、本発明の範囲を限定するものではないことが
理解されるべきであり、本発明の範囲は、特許請求の範囲によって画定される。したがっ
て、他の実施形態も以下の特許請求の範囲内にある。例えば、本発明の範囲から逸脱する
ことなしに、様々な変更が行われてよい。さらに、上述のステップのうちいくつかは、
順序に依存しないものであってよく、したがって、記載された順序とは異なる順序で実施
され得る。

【図1】



【図2】

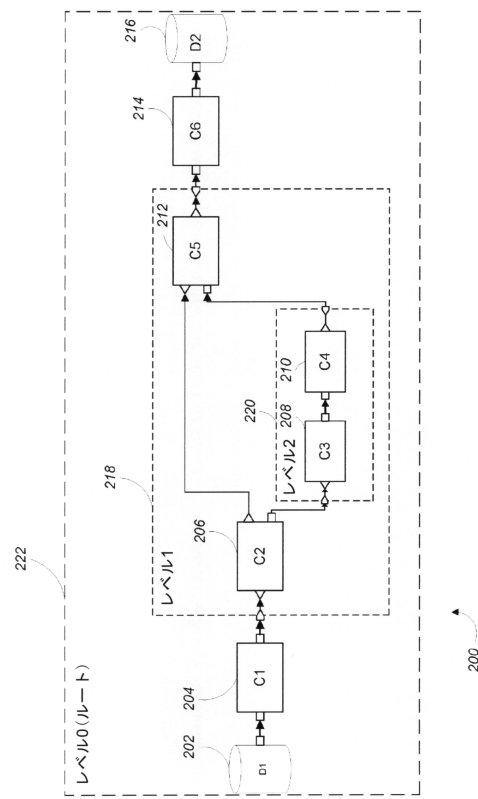


FIG. 2

【図 3 A】

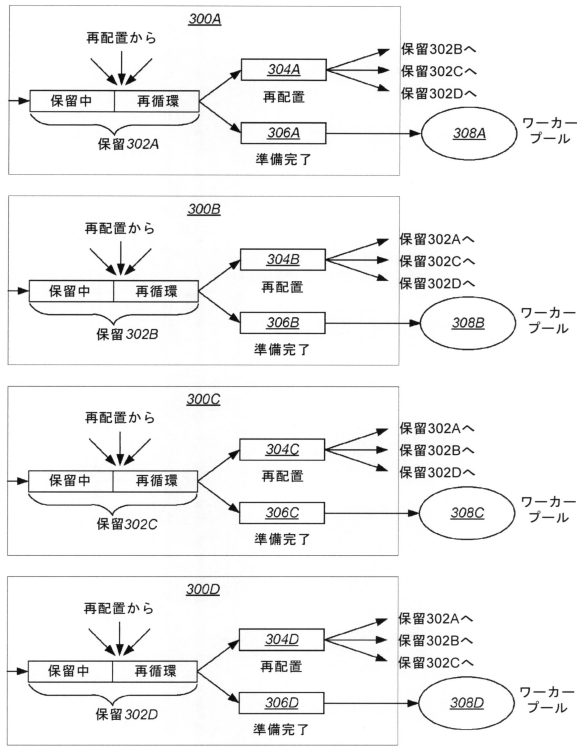


FIG. 3A

【図 3 B】

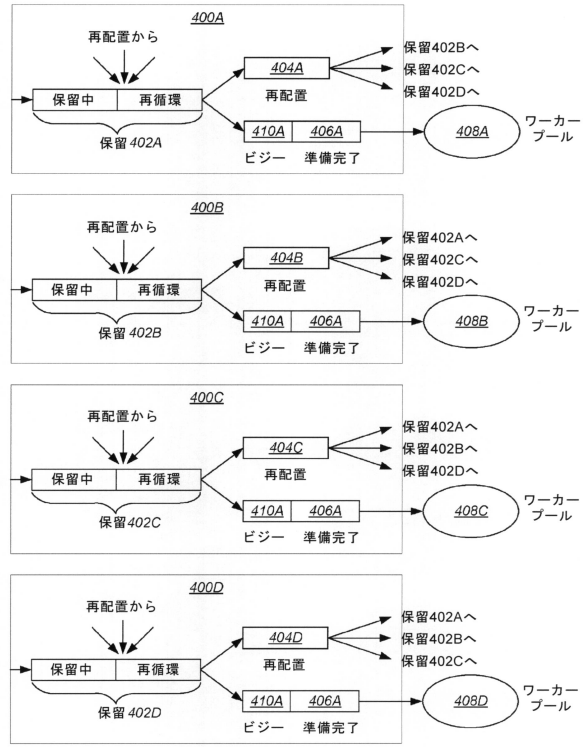


FIG. 3B

【図 4 A】

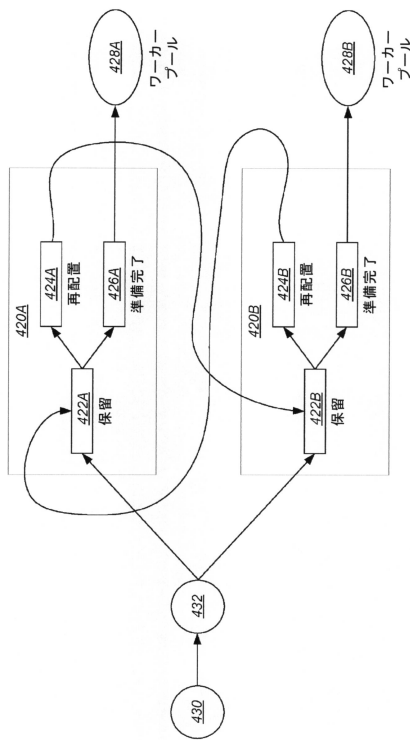


FIG. 4A

【図 4 B】

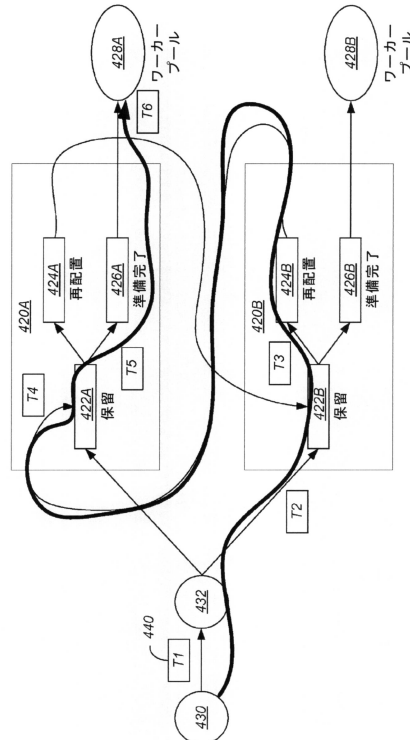


FIG. 4B

【図 5 A】

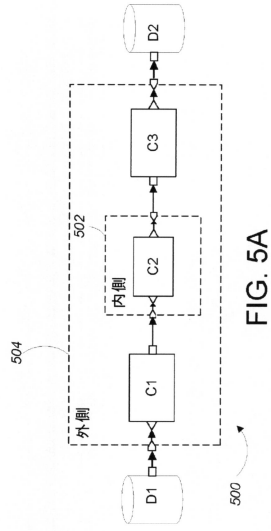


FIG. 5A

【図 5 B】

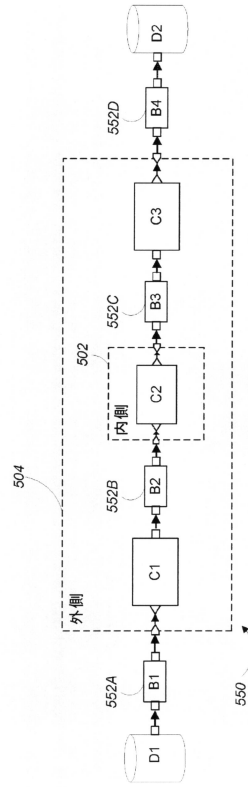


FIG. 5B

フロントページの続き

- (74)代理人 100131093
弁理士 堀内 真
- (74)代理人 100150902
弁理士 山内 正子
- (74)代理人 100141391
弁理士 園元 修一
- (74)代理人 100198074
弁理士 山村 昭裕
- (74)代理人 100145920
弁理士 森川 聡
- (74)代理人 100096013
弁理士 富田 博行
- (72)発明者 ニューバーン ジェフェリー
アメリカ国 01776 マサチューセッツ サドベリー ケンドールロード20
- (72)発明者 スタンフィル クレイグ ダブリュー.
アメリカ国 01773 マサチューセッツ リンカン ハックルベリーロード43

審査官 漆原 孝治

- (56)参考文献 特開2006-259812(JP,A)
特開2008-159081(JP,A)
特開2013-239199(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G06F 9/54
G06F 9/50