

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6531056号  
(P6531056)

(45) 発行日 令和1年6月12日 (2019.6.12)

(24) 登録日 令和1年5月24日 (2019.5.24)

(51) Int. Cl.

F I

G O 6 F 16/00 (2019.01)

G O 6 F 17/30 3 4 0 Z

G O 6 F 12/00 (2006.01)

G O 6 F 12/00 5 1 3 D

請求項の数 13 (全 17 頁)

(21) 出願番号 特願2016-61823 (P2016-61823)  
 (22) 出願日 平成28年3月25日 (2016.3.25)  
 (65) 公開番号 特開2017-174290 (P2017-174290A)  
 (43) 公開日 平成29年9月28日 (2017.9.28)  
 審査請求日 平成30年3月7日 (2018.3.7)

(73) 特許権者 000208891  
 K D D I 株式会社  
 東京都新宿区西新宿二丁目3番2号  
 (74) 代理人 100092772  
 弁理士 阪本 清孝  
 (74) 代理人 100119688  
 弁理士 田邊 壽二  
 (72) 発明者 斉藤 和広  
 埼玉県ふじみ野市大原二丁目1番15号  
 株式会社 K D D I 研究所内  
 審査官 後藤 昂彦

最終頁に続く

(54) 【発明の名称】 クエリ処理制御装置、制御方法及び制御プログラム

(57) 【特許請求の範囲】

【請求項 1】

クライアントが投稿したクエリに対してストレージから必要なデータを取得し前記クエリの演算をバッファ上で処理して生成した演算結果をユーザに返すクエリ処理装置と、

前記クエリの演算処理に際して利用するメモリ上のバッファ領域を管理して前記バッファのサイズを指定するバッファ管理装置と、を備え、

前記バッファ管理装置は、

前記クエリ処理装置から実行するクエリの実行計画を取得して、事前設定された利用可能なメモリサイズとクエリ処理装置の演算の実装情報から演算毎のバッファの初期設定サイズを確保したバッファ情報を作成するバッファ計画部と、

前記クエリ処理装置から演算指定でバッファ情報の提供依頼を受け、前記バッファ情報を提供するとともに、前記クエリ処理装置から演算指定でバッファ増量依頼を受け、バッファ情報における当該バッファの利用状況に応じて終了した演算のバッファ領域を移行することでバッファのサイズを変更するバッファ制御部と、

前記クエリ処理装置から演算処理によって変化したバッファの利用状況を取得してバッファ情報を更新するバッファ利用量収集部と

を含むことを特徴とするクエリ処理制御装置。

【請求項 2】

前記クエリ処理装置は、クエリ計画部とクエリ実行制御部と演算処理部を備え、

前記クエリ計画部は、クライアントが投稿したクエリを解釈し、クエリの実行計画を生

成して前記バッファ計画部及びクエリ実行制御部に実行計画を渡し、

前記クエリ実行制御部は、前記実行計画を基に前記バッファ制御部からバッファ情報を取得して演算が利用するバッファサイズを指定し、クエリ処理方式に沿って前記演算処理部に演算処理を依頼し、必要に応じて前記バッファ制御部にバッファの増量依頼し、

前記演算処理部は、前記クエリ実行制御部から依頼された演算を処理し、バッファの利用量の変化を前記バッファ利用量収集部に渡す

ことを特徴とする請求項 1 に記載のクエリ処理制御装置。

【請求項 3】

前記バッファ制御部は、

前記クエリ処理装置から演算指定でバッファ増量依頼を受けた際に、増量対象が出力バッファである場合、当該演算の制御バッファを当該演算の出力バッファ領域に移行する請求項 1 又は請求項 2 に記載のクエリ処理制御装置。

10

【請求項 4】

前記バッファ制御部は、

前記クエリ処理装置から演算指定でバッファ増量依頼を受けた際に、増量対象が制御バッファである場合、上限値が設定された当該演算の出力バッファを当該演算の制御バッファ領域に移行する請求項 1 又は請求項 2 に記載のクエリ処理制御装置。

【請求項 5】

前記バッファ制御部は、

前記クエリ処理装置から演算指定でバッファ増量依頼を受けた際に、前記演算以外の他の演算の入力が完了している制御バッファを増量対象の演算のバッファ領域に移行する請求項 1 又は請求項 2 に記載のクエリ処理制御装置。

20

【請求項 6】

前記バッファ制御部は、

前記クエリ処理装置から演算指定でバッファ増量依頼を受けた際に、上限値が設定された空きのバッファを増量対象の演算のバッファ領域に移行する請求項 1 又は請求項 2 に記載のクエリ処理制御装置。

【請求項 7】

クライアントが投稿したクエリに対してストレージから必要なデータを取得し前記クエリの演算をバッファ上で処理して生成した演算結果をユーザに返すクエリ処理において、

30

前記クエリを解釈して生成したツリー構造の最下位の演算から順に処理する実体化型クエリ処理方式によるクエリ演算の実行計画を生成する手順と、

事前設定された利用可能なメモリサイズと演算の実装情報から前記クエリの演算毎のバッファサイズを指定する手順と、

演算結果が格納されたバッファが上限であり且つ下位の出力バッファにデータがある場合にバッファサイズの増量依頼を受け、バッファの利用状況に応じて終了した演算のバッファ領域を空き領域として検索して出力バッファまたは制御バッファのいずれかの対象バッファに付与する手順と、を備えることで、

クエリ実行時に、前記演算毎に利用するバッファサイズを動的に制御することで、ストレージ利用と演算切り換えの回数を削減することを特徴とするクエリ処理制御方法。

40

【請求項 8】

クライアントが投稿したクエリに対してストレージから必要なデータを取得し前記クエリの演算をバッファ上で処理して生成した演算結果をユーザに返すクエリ処理において、

前記クエリを解釈して生成したツリー構造の最上位から順に下位の演算にデータ生成を依頼して処理する要求駆動型クエリ処理方式によるクエリ演算の実行計画を生成する手順と、

事前設定された利用可能なメモリサイズと前記クエリ演算の実装情報から演算毎のバッファサイズを指定する手順と、

制御バッファが上限であり且つ下位の出力バッファにデータがある場合に制御バッファ

50

のバッファサイズの増量依頼を受けるとともに、演算結果が格納されたバッファ（出力バッファもしくは制御バッファ）が上限に達しておらず且つ直下の演算が終了していない場合に出力バッファのバッファサイズの増量依頼を受け、バッファの利用状況に応じて終了した演算のバッファ領域を空き領域として検索して出力バッファまたは制御バッファのいずれかの対象バッファに付与する手順と、を備えることで、

クエリ実行時に、前記演算毎に利用するバッファサイズを動的に制御することで、ストレージ利用と演算切り換えの回数を削減することを特徴とするクエリ処理制御方法。

【請求項 9】

演算に対してバッファサイズの増量依頼を受けた際に、増量対象が出力バッファである場合、当該演算の制御バッファを前記空き領域とする請求項 7 又は請求項 8 に記載のクエリ処理制御方法。

【請求項 10】

演算に対してバッファサイズの増量依頼を受けた際に、増量対象が制御バッファである場合、上限値が設定された当該演算の出力バッファを前記空き領域とする請求項 7 又は請求項 8 に記載のクエリ処理制御方法。

【請求項 11】

演算に対してバッファサイズの増量依頼を受けた際に、前記演算以外の他の演算の入力が完了している制御バッファを前記空き領域とする請求項 7 又は請求項 8 に記載のクエリ処理制御方法。

【請求項 12】

演算に対してバッファサイズの増量依頼を受けた際に、上限値が設定された空きのバッファを前記空き領域とする請求項 7 又は請求項 8 に記載のクエリ処理制御方法。

【請求項 13】

請求項 7 から請求項 12 のいずれか 1 項に記載の各手順をコンピュータに実行させるクエリ処理制御プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、SQLクエリを処理するデータベースシステムにおいて、利用可能なメモリサイズを超えるSQLクエリを高速に実行可能とするためのクエリ処理制御装置、制御方法及び制御プログラムに関する。

【背景技術】

【0002】

データベースシステムにおけるクエリ処理装置は、ユーザが投稿したクエリを解釈し、必要なデータをストレージから取得してクエリを処理する。多くのデータベースシステムではSQLと呼ばれるクエリ言語を利用しており、クエリ処理装置はSQLクエリを解釈することで、関係代数に沿った演算を抽出し、各演算の処理順序を決めることができる。この処理順序は、関係代数式としてツリー構造をしたクエリ実行計画で表現される。

演算はデータに対する加工処理を表現している。選択演算や射影演算などはデータから不要な部分を削減するが、一方で結合演算や集合演算のように複数のデータを組み合わせる演算や、集計演算を多数行う場合など、元のデータよりサイズを大きくする演算が存在する。

【0003】

クエリ処理装置は、クエリ処理を装置のメモリ上で行う必要がある。以降、クエリ処理において演算が利用するメモリ領域をバッファと呼ぶ。メモリは有限であり、また多数のクエリを同時に実行する場合など、バッファに利用できる物理メモリが小さく制限される場合がある。元のデータが大規模である場合や、前述のようにデータサイズを大きくする演算によって、このようなバッファに利用可能なメモリサイズを超えるデータが発生する場合、OSのスワップ機構の利用によって大幅な遅延が発生する、又は、クエリ実行に失敗

10

20

30

40

50

する事態が発生する。

【 0 0 0 4 】

クエリ処理方式としては、実体化型クエリ処理方式 (Materialized model又はBulk processing model) と、要求駆動型クエリ処理方式 (Volcano model又はPipeline model) が提案されている。

実体化型クエリ処理方式は、非特許文献 1 に示されるように、クエリ実行計画におけるツリー構造の最下位の演算から順に処理する。演算は、入力される全てのデータを処理完了すると、次の上位の演算に処理を移行する。

本方式では、演算が利用するバッファのサイズに制限がないため、演算の処理中に対象のデータサイズが利用可能なメモリサイズを超える場合がある。それによって処理の大幅な遅延やクエリ実行失敗が発生する。

10

【 0 0 0 5 】

要求駆動型クエリ処理方式は、非特許文献 2 に示されるように、クエリ実行計画におけるツリー構造の最上位から順に下位の演算にデータ生成を依頼する。依頼された演算は、データを生成するために演算に必要なデータを同様に下位の演算に依頼する。入力データが得られた演算は、演算処理によって要求されたデータを生成できた段階で上位の演算に処理を移行する。

この要求単位は、1 レコードもしくはストレージのディスクブロックなどの固定サイズであり、利用可能なメモリサイズを超える事象は発生しない。

【 0 0 0 6 】

20

しかし、大規模なデータを対象とする処理において、演算の切り替えが頻繁に発生し、性能が劣化する。また、仮に要求単位を大きく設定した場合、演算毎に所持するデータの規模が膨らみ、同時実行のクエリが増えた場合や、演算の数が多いくエリが実行された場合において、実体化型クエリ処理方式と同様に演算が利用するバッファの制御をしないため、利用可能なメモリサイズを超える現象が発生する。

更には、全データの蓄積を必要とする演算 (整列演算等) や、二項演算 (結合演算等) における片方のデータの蓄積が必要な実装 (ハッシュ結合等) においては、実体化型クエリ処理方式と同様に、その蓄積データが利用可能なメモリサイズを超える現象が発生する。

【 0 0 0 7 】

30

本発明者は、実体化型クエリ処理方式をベースに、大規模データにおいても確実にクエリ処理を完了する分割実行型クエリ処理方式の提案を行った (特許文献 1)。

この方式は、利用可能なメモリサイズを超えた場合に、OSのスワップ機構とは別の仕組みでストレージを活用することで、確実にクエリ処理を完了する。更に、ストレージの読み込み回数を最小化する演算処理手法を導入することで、性能劣化を抑止している。

【 先行技術文献 】

【 特許文献 】

【 0 0 0 8 】

【 特許文献 1 】 特願 2 0 1 5 - 0 0 2 2 9 1

【 非特許文献 】

40

【 0 0 0 9 】

【 非特許文献 1 】 F. Groffen, N. Nes, S. Nes, and S. Kersten, "MonetDB: Two Decades of Research in Column-oriented Database Architectures", Bulletin of the IEEE Computer Society Technical Committee on Data Engineering, pp. 1-6, 2012.

【 非特許文献 2 】 G. Graefe, "Volcano-an extensible and parallel query evaluation system", IEEE Transactions on Knowledge and Data Engineering, vol. 6, no. 1, pp. 120-135, 1994

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 1 0 】

50

しかしながら、特許文献 1 に記載の実体化型クエリ処理方式を利用した処理によれば、データサイズを削減するような演算においても、入力データが利用可能メモリサイズを超える場合、要求駆動型クエリ処理方式のように上位の演算に処理後のデータを渡すことをせず、処理後の少量データであっても一時的にストレージに格納するため、無駄なストレージ利用が発生し性能が悪化する可能性があった。

【 0 0 1 1 】

本発明は上記実情に鑑みて提案されたもので、要求駆動型クエリ処理方式及び実体化型クエリ処理方式のクエリ処理制御装置において、ストレージの利用回数と演算の切り替え回数を最小化して高速なクエリ処理を実現することで、大規模データを対象とするクエリ処理を確実に実行可能とするクエリ処理制御装置、制御方法及び制御プログラムを提供する

10

【課題を解決するための手段】

【 0 0 1 2 】

上記目的を達成するため本発明は、演算単位のバッファの利用量が上限に達した段階で次の演算に処理を推移させてストレージの利用回数を削減し、更にバッファサイズをクエリ実行時に動的に制御して演算の切り替えを最小限にすることで、確実にクエリ処理を完了しつつ性能の向上を図っている。

【 0 0 1 3 】

すなわち、請求項 1 のクエリ処理制御装置は、

クライアントが投稿したクエリに対してストレージから必要なデータを取得し前記クエリの演算をバッファ上で処理して生成した演算結果をユーザに返すクエリ処理装置と、

20

前記クエリの演算処理に際して利用するメモリ上のバッファ領域を管理して前記バッファのサイズを指定するバッファ管理装置と、を備え、

前記バッファ管理装置は、

前記クエリ処理装置から実行するクエリの実行計画を取得して、事前設定された利用可能なメモリサイズとクエリ処理装置の演算の実装情報から演算毎のバッファの初期設定サイズを確保したバッファ情報を作成するバッファ計画部と、

前記クエリ処理装置から演算指定でバッファ情報の提供依頼を受け、前記バッファ情報を提供するとともに、前記クエリ処理装置から演算指定でバッファ増量依頼を受け、バッファ情報における当該バッファの利用状況に応じて終了した演算のバッファ領域を移行することでバッファのサイズを変更するバッファ制御部と、

30

前記クエリ処理装置から演算処理によって変化したバッファの利用状況を取得してバッファ情報を更新するバッファ利用量収集部とを含むことを特徴としている。

【 0 0 1 4 】

請求項 2 は、請求項 1 のクエリ処理制御装置において、

前記クエリ処理装置は、クエリ計画部とクエリ実行制御部と演算処理部を備え、

前記クエリ計画部は、クライアントが投稿したクエリを解釈し、クエリの実行計画を生成して前記バッファ計画部及びクエリ実行制御部に実行計画を渡し、

前記クエリ実行制御部は、前記実行計画を基に前記バッファ制御部からバッファ情報を取得して演算が利用するバッファサイズを指定し、クエリ処理方式に沿って前記演算処理部に演算処理を依頼し、必要に応じて前記バッファ制御部にバッファの増量依頼し、

40

前記演算処理部は、前記クエリ実行制御部から依頼された演算を処理し、バッファの利用量の変化を前記バッファ利用量収集部に渡すことを特徴としている。

【 0 0 1 5 】

請求項 3 は、請求項 1 又は請求項 2 のクエリ処理制御装置において、

前記バッファ制御部は、

前記クエリ処理装置から演算指定でバッファ増量依頼を受けた際に、増量対象が出力バッファである場合、当該演算の制御バッファを当該演算の出力バッファ領域に移行することを特徴としている。

【 0 0 1 6 】

50

請求項 4 は、請求項 1 又は請求項 2 のクエリ処理制御装置において、  
前記バッファ制御部は、

前記クエリ処理装置から演算指定でバッファ増量依頼を受けた際に、増量対象が制御バッファである場合、上限値が設定された当該演算の出力バッファを当該演算の制御バッファ領域に移行することを特徴としている。

【0017】

請求項 5 は、請求項 1 又は請求項 2 のクエリ処理制御装置において、  
前記バッファ制御部は、

前記クエリ処理装置から演算指定でバッファ増量依頼を受けた際に、前記演算以外の他の演算の入力が完了している制御バッファを増量対象の演算のバッファ領域に移行することを特徴としている。

10

【0018】

請求項 6 は、請求項 1 又は請求項 2 のクエリ処理制御装置において、  
前記バッファ制御部は、

前記クエリ処理装置から演算指定でバッファ増量依頼を受けた際に、上限値が設定された空きのバッファを増量対象の演算のバッファ領域に移行することを特徴としている。

【0019】

請求項 7 のクエリ処理制御方法は、

クライアントが投稿したクエリに対してストレージから必要なデータを取得し前記クエリの演算をバッファ上で処理して生成した演算結果をユーザに返すクエリ処理において、

20

前記クエリを解釈して生成したツリー構造の最下位の演算から順に処理する実体化型クエリ処理方式によるクエリ演算の実行計画を生成する手順と、

事前設定された利用可能なメモリサイズと演算の実装情報から前記クエリの演算毎のバッファサイズを指定する手順と、

演算結果が格納されたバッファが上限であり且つ下位の出力バッファにデータがある場合にバッファサイズの増量依頼を受け、バッファの利用状況に応じて終了した演算のバッファ領域を空き領域として検索して出力バッファまたは制御バッファのいずれかの対象バッファに付与する手順と、を備えることで、

クエリ実行時に、前記演算毎に利用するバッファサイズを動的に制御することで、ストレージ利用と演算切り換えの回数を削減することを特徴としている。

30

【0020】

請求項 8 のクエリ処理制御方法は、

クライアントが投稿したクエリに対してストレージから必要なデータを取得し前記クエリの演算をバッファ上で処理して生成した演算結果をユーザに返すクエリ処理において、

前記クエリを解釈して生成したツリー構造の最上位から順に下位の演算にデータ生成を依頼して処理する要求駆動型クエリ処理方式によるクエリ演算の実行計画を生成する手順と、

事前設定された利用可能なメモリサイズと前記クエリ演算の実装情報から演算毎のバッファサイズを指定する手順と、

制御バッファが上限であり且つ下位の出力バッファにデータがある場合に制御バッファのバッファサイズの増量依頼を受けるとともに、演算結果が格納されたバッファ（出力バッファもしくは制御バッファ）が上限に達しておらず且つ直下の演算が終了していない場合に出力バッファのバッファサイズの増量依頼を受け、バッファの利用状況に応じて終了した演算のバッファ領域を空き領域として検索して出力バッファまたは制御バッファのいずれかの対象バッファに付与する手順と、を備えることで、

40

クエリ実行時に、前記演算毎に利用するバッファサイズを動的に制御することで、ストレージ利用と演算切り換えの回数を削減することを特徴としている。

【0021】

請求項 9 は、請求項 7 又は請求項 8 のクエリ処理制御方法において、

演算に対してバッファサイズの増量依頼を受けた際に、増量対象が出力バッファである

50

場合、当該演算の制御バッファを前記空き領域とすることを特徴としている。

【0022】

請求項10は、請求項7又は請求項8のクエリ処理制御方法において、

演算に対してバッファサイズの増量依頼を受けた際に、増量対象が制御バッファである場合、上限値が設定された当該演算の出力バッファを前記空き領域とすることを特徴としている。

【0023】

請求項11は、請求項7又は請求項8のクエリ処理制御方法において、

演算に対してバッファサイズの増量依頼を受けた際に、前記演算以外の他の演算の入力が完了している制御バッファを前記空き領域とすることを特徴としている。

10

【0024】

請求項12は、請求項7又は請求項8のクエリ処理制御方法において、

演算に対してバッファサイズの増量依頼を受けた際に、上限値が設定された空きのバッファを前記空き領域とすることを特徴としている。

【0025】

請求項13のクエリ処理制御プログラムは、請求項7から請求項12のいずれか1項に記載の各手順をコンピュータに実行させることを特徴としている。

【発明の効果】

【0026】

請求項1のクエリ処理制御装置によれば、クエリ処理装置から演算指定でバッファ増量依頼を受け、バッファ情報における当該バッファの利用状況に応じて終了した演算のバッファ領域を移行することでバッファのサイズを変更するので、クエリ処理実行時に演算毎に利用するメモリサイズを制御することができる。その結果、演算が利用するバッファを動的に制御して、利用可能なメモリサイズを超えることを抑制し、ストレージの利用と演算の切り替えを最小限にすることで、大規模データの確実なクエリ処理実行及び処理の高速化を実現することができる。

20

【0027】

請求項2のクエリ処理制御装置によれば、クエリ実行制御部が実行計画を基にバッファ制御部からバッファ情報を取得して演算が利用するバッファサイズを指定し、クエリ処理方式に沿って演算処理部に演算処理を依頼し、必要に応じてバッファ制御部にバッファの増量依頼することで、バッファのサイズを変更することができる。

30

【0028】

請求項7のクエリ処理制御方法によれば、クエリを解釈して生成したツリー構造の最下位の演算から順に処理する実体化型クエリ処理方式により、クエリ演算の実行計画を生成することができる。

【0029】

請求項8のクエリ処理制御方法によれば、クエリを解釈して生成したツリー構造の最上位から順に下位の演算にデータ生成を依頼して処理する要求駆動型クエリ処理方式により、クエリ演算の実行計画を生成することができる。

【0030】

請求項3のクエリ処理制御装置及び請求項9のクエリ処理制御方法によれば、演算の増量対象が出力バッファである場合に、当該演算の制御バッファを当該演算のバッファ領域に移行することでバッファのサイズを変更することができる。

40

【0031】

請求項4のクエリ処理制御装置及び請求項10のクエリ処理制御方法によれば、演算の増量対象が制御バッファである場合に、上限値が設定された当該演算の出力バッファを当該演算のバッファ領域に移行することでバッファのサイズを変更することができる。

【0032】

請求項5のクエリ処理制御装置及び請求項11のクエリ処理制御方法によれば、他の演算の入力が完了している制御バッファを増量対象の演算のバッファ領域に移行することで

50

バッファのサイズを変更することができる。

【 0 0 3 3 】

請求項 6 のクエリ処理制御装置及び請求項 1 2 のクエリ処理制御方法によれば、上限値が設定された空きのバッファを増量対象の演算のバッファ領域に移行することでバッファのサイズを変更することができる。

【 0 0 3 4 】

請求項 1 3 のクエリ処理制御プログラムによれば、コンピュータ上にクエリ処理制御装置を構築させることができる。

【図面の簡単な説明】

【 0 0 3 5 】

【図 1】本発明の一実施形態に係るクエリ処理制御装置の構成を示すブロック図である。

【図 2】バッファ情報の構成を表す説明図である。

【図 3】バッファ計画部における処理フローを示すフローチャート図である。

【図 4】バッファ制御部におけるバッファ増量処理の処理フローを示すフローチャート図である。

【図 5】バッファ利用量収集部における処理フローを示すフローチャート図である。

【図 6】本発明の他の実施形態に係るクエリ処理制御装置の構成を示すブロック図である。

。

【図 7】実体化型クエリ処理方式におけるクエリ実行制御部の処理フローを示すフローチャート図である。

【図 8】要求駆動型クエリ処理方式におけるクエリ実行制御部の処理フローを示すフローチャート図である。

【図 9】実体化型クエリ処理方式及び要求駆動型クエリ処理方式における演算処理部の処理フローを示すフローチャート図である。

【発明を実施するための形態】

【 0 0 3 6 】

本発明の一実施形態のクエリ処理制御装置について、図 1 の構成ブロック図を参照しながら説明する。

クライアント 1、ストレージ 2、バッファ 3、クエリ処理装置 10 及びバッファ管理装置 20 は、基本プログラムや各種の基本デバイスが記憶された ROM と、各種のプログラムやデータが記憶されるハードディスクドライブ装置 (HDD) と、CR-ROM や DVD 等の記憶媒体からプログラムやデータを読み出すメディアドライブ装置と、プログラムを実行する CPU と、この CPU にワークエリアを提供する RAM と、外部装置と通信するパラレル/シリアル I/F とを主要部分とする一般的な構成を備えたコンピュータ上に構築されている。

例えば、上述した構成を有するコンピュータにおいて、クエリ処理を実行するためのクエリ処理制御プログラムがメディアドライブ装置を介して HDD にインストールされることでクエリ処理制御装置が構築される。

【 0 0 3 7 】

クライアント 1 は、クエリ処理装置 10 に対してクエリを投稿し、処理結果を取得する。クエリ処理装置 10 は投稿クエリの処理結果を出力し、バッファ管理装置 20 はクエリの演算処理を行うバッファ (バッファ領域) 3 を管理する。

クエリ処理装置 10 は、クライアント 1 が投稿したクエリに対して、ストレージ 2 から必要なデータを取得して、バッファ管理装置 20 が指定するサイズのバッファ (バッファ領域) 3 上でクエリの演算を処理することで、処理結果を生成し、ユーザに処理結果を返す。

【 0 0 3 8 】

ストレージ 2 は、クエリ処理装置 10 上及びネットワーク経由先の磁気ディスク、SSD (Solid State Drive)、フラッシュメモリ、物理メモリ、外部記憶媒体を含む装置であり、クエリ処理対象のデータを所持している。

10

20

30

40

50



## 【 0 0 3 9 】

バッファ管理装置 2 0 は、バッファ 3 のサイズを決定するバッファ計画部 2 1 と、バッファ 3 の利用状況からサイズを変更するバッファ制御部 2 2 と、バッファ 3 の利用状況を取得するバッファ利用量収集部 2 3 から構成されている。

バッファ計画部 2 1 は、クエリ処理装置 1 0 から実行するクエリの実行計画を取得して、事前設定された利用可能なメモリサイズ（利用可能メモリ 2 6 ）と、クエリ処理装置 1 0 の演算の実装情報（演算の実装情報 2 5 ）から演算毎のバッファ 3 の初期設定サイズを決定し、バッファ情報 2 7 を作成する。

## 【 0 0 4 0 】

バッファ制御部 2 2 は、クエリ処理装置 1 0 から演算指定でバッファ情報提供依頼を受け、バッファ情報 2 7 を提供するとともに、クエリ処理装置 1 0 から演算指定でバッファ増量依頼を受け、バッファ情報 2 7 における当該バッファの利用状況からバッファ 3 のサイズを変更する。バッファサイズを変更することで、演算が利用するバッファを動的に制御し、演算の切り替えを最小限にすることができる。また、バッファを増量することで、演算が利用可能なメモリサイズを超えることを抑制することができる。

## 【 0 0 4 1 】

バッファ利用量収集部 2 3 は、クエリ処理装置 1 0 から演算処理によって変化したバッファ 3 の利用状況を取得してバッファ情報 2 7 の更新を行うことで、演算処理で利用するメモリ上のバッファ 領域を管理することができる。

## 【 0 0 4 2 】

次に、バッファ管理装置 2 0 においてバッファ情報 2 7 の作成及び変更を行う場合の処理について、図 2 ～ 5 を参照して説明する。

まず、バッファ情報 2 7 の構成について、図 2 を参照して説明する。op[i]はi番目の演算処理（演算を実際に処理する動作）を示し、B[i]は演算処理 op[i]が利用するバッファを示す。

バッファ B は、投稿クエリの演算単位に生成される。バッファは、処理結果を格納する出力バッファoutbufと、演算処理に必要な中間結果を格納する制御バッファctrlbufから構成される。

演算（処理の内容）には、「結合」「直積」「集計」「整列」「集合」「選択」「射影」等の関係代数で定義されているものが存在し、演算専用のバッファが用意される。演算で処理される入力データは、その演算の直下の演算の出力バッファと、自身専用の制御バッファから取得する。

## 【 0 0 4 3 】

通常の実行計画では、最下位の演算はストレージ 2 からデータを取得するScan演算であり、最上位の演算はクライアント 1 に処理結果を返すOutput演算である。バッファ利用の例外として、Scan演算は入力先がストレージ 2 であり、Output演算は専用のバッファが存在しない。従って、演算毎のバッファの数は、演算の数 - 1 となる。また、演算の実装アルゴリズムによって、制御バッファの有無が異なる。

例えば、結合演算のアルゴリズムがHash Joinの場合、片方の入力データから生成したハッシュテーブルを制御バッファに保管する必要がある。一方で、射影演算のように、一時的にデータを保存する必要のない演算の場合、制御バッファは不要となる。このために、クエリ処理装置 1 0 が処理可能な演算それぞれに対する制御バッファの有無を、演算の実装情報 2 5 に定義する。

## 【 0 0 4 4 】

バッファ管理装置 2 0 のバッファ計画部 2 1 における処理手順について、図 3 のフローチャート図を参照して説明する。

クエリ処理装置 1 0 側からクエリの実行計画を取得し（ステップ 3 1 ）、演算の構成から、バッファ情報 2 7 の元となるツリー構造を生成する。

次に、事前に設定された利用可能なメモリサイズと、演算の実装情報を元に、バッファ情報の各演算にバッファのサイズを設定していく（ステップ 3 2 ）。演算の実装情報には

10

20

30

40

50

、実行計画で出現する演算の種類毎に、制御バッファの有無や、サイズの設定方法が設定されている。

【 0 0 4 5 】

最も単純なバッファサイズの決定方法は、利用可能なメモリサイズをバッファの数で割った値をバッファサイズに設定する。制御バッファがあるバッファは、そのバッファサイズを半分にしてお出力バッファと制御バッファのサイズに設定する。

各演算のバッファをメモリ上に確保し、バッファ情報 2 7 を作成する ( ステップ 3 3 )

。

【 0 0 4 6 】

バッファ制御部 2 2 におけるバッファ増量処理の処理手順について、図 4 のフローチャート図を参照して説明する。

バッファ制御部 2 2 がクエリ処理装置 1 0 から演算のバッファの増量依頼を受けると ( ステップ 4 1 )、バッファ情報 2 7 から空いているバッファを検索し ( ステップ 4 2 )、バッファの空き領域の有無を判定する ( ステップ 4 3 )。

ステップ 4 3 でバッファの空き領域を確保できた場合、対象のバッファに付与し ( ステップ 4 4 )、かつ当該空き領域を開放する。これらの変更をバッファ情報 2 7 に適用する ( ステップ 4 5 )。

【 0 0 4 7 】

ステップ 4 3 で空き領域を確保できなかった場合、増量対象が制御バッファか出力バッファのどちらなのかを判断する ( ステップ 4 6 )。

増量対象が制御バッファの場合は、ストレージ 2 の領域を付与して ( ステップ 4 7 )、バッファ情報を更新する ( ステップ 4 5 )。

増量対象が出力バッファの場合は、増量不可としてクエリ処理装置 1 0 に返す ( ステップ 4 8 )。なお、出力バッファは拡張しないことで処理が止まることはないため、ストレージ 2 の領域を付与しない。そのため、ストレージ 2 のバッファ領域としての利用を最小限に抑えることができる。

【 0 0 4 8 】

ステップ 4 3 におけるバッファの空き領域の有無の判定基準は複数存在し、以下に列挙する。これらは、それぞれ単体として利用可能であり、また組み合わせて利用することも可能である。

・演算処理が終了しているバッファ ( 出力バッファ及び制御バッファ ) の未使用領域を「空き」と判定する。

・増量依頼されたバッファが出力バッファの場合、当該演算の制御バッファの未使用領域を「空き」と判定する。ただし、制御バッファへの入力が完了していることを前提とする。

。

・増量依頼されたバッファの演算以外の入力が完了している制御バッファの未使用領域を「空き」と判定する。

・増量依頼されたバッファが制御バッファの場合、当該演算の出力バッファの未使用領域を「空き」と判定する。ただし、まだ利用される可能性があるため、全領域を取得されないために、空き領域として取得可能な上限サイズを出力バッファサイズの割合 ( 例えば、5 0 % まで等 ) で設定する必要がある。

・増量依頼されたバッファの演算以外のバッファの未使用領域を「空き」と判定する。ただし、まだ利用される可能性があるため、全領域を取得されないために、空き領域として取得可能な上限サイズを出力バッファサイズの割合 ( 例えば、5 0 % まで等 ) で設定する必要がある。

・増量依頼されたバッファが制御バッファの場合、外部記憶装置としてのストレージの空き領域を「空き」と判定してもよい。

【 0 0 4 9 】

次に、バッファ利用量収集部 2 3 における処理手順について、図 5 のフローチャート図を参照して説明する。

10

20

30

40

50

まず、クエリ処理装置 10 の演算処理におけるバッファの利用状況を取得する（ステップ 51）。取得する利用状況の情報は、入力元のバッファ（直下の演算の出力バッファもしくは当該演算の制御バッファ）、入力データサイズ、出力先のバッファ（当該演算の出力バッファもしくは制御バッファ）、出力データサイズ、該当の演算の処理完了フラグ、制御バッファがある場合の制御バッファへの入力完了フラグを含む。

次に、取得した利用状況の情報から、バッファ情報における入力元のバッファのサイズと出力先のバッファの利用サイズ、空き領域、完了状況を更新する（ステップ 52 及びステップ 53）。

#### 【0050】

上述したクエリ処理制御装置によれば、クエリ処理装置 10 でのクエリ処理実行時ににおいて、バッファ管理装置 20 のバッファ情報 27 で利用状況を管理することで、クエリの演算毎に利用するバッファサイズを動的に制御し、利用可能なメモリサイズを超えることを抑制し、ストレージ 2 の利用と演算の切り替えを最小限にすることで、大規模データの確実なクエリ処理実行及び処理の高速化を実現することができる。

#### 【0051】

続いて、クエリ処理装置 10 の実施形態として、実体化型クエリ処理方式と、要求駆動型クエリ処理方式に適用した場合の構成について、図 6 を参照して説明する。

図 6 のクエリ処理制御装置のクエリ処理装置 10 は、クエリの実行計画（ツリー構造）を生成するクエリ計画部 11 と、バッファの増量依頼をするクエリ実行制御部 12 と、バッファの利用量の変化を把握する演算処理部 13 とを備えて構成されている。

#### 【0052】

クエリ計画部 11 は、クライアント 1 が投稿したクエリを解釈し、クエリの実行計画を生成してバッファ計画部 21 及びクエリ実行制御部 12 に実行計画を渡す。

クエリ実行制御部 12 は、実行計画を基に、バッファ制御部 22 からバッファ情報を取得して演算が利用するバッファサイズを指定し、クエリ処理方式に沿って演算処理部 13 に演算処理を依頼し、必要に応じてバッファ制御部 22 にバッファの増量依頼をする。

演算処理部 13 は、クエリ実行制御部 12 から依頼された演算を処理し、バッファの利用量の変化をバッファ利用量収集部 23 に渡す。

#### 【0053】

実体化型クエリ処理方式を採用した場合のクエリ処理制御装置におけるクエリ実行制御部 12 の処理手順について、図 7 のフローチャート図を参照して説明する。

クエリ計画部 11 から実行計画を取得し（ステップ 701）、ツリー構造の最下位の演算から処理を開始する（ステップ 702）。通常、最下位の演算は、ストレージ 2 からデータを取得する Scan 演算であり、これを選択する。

ここで、選択された演算とその直下の演算のバッファ情報をバッファ制御部 22 から取得する（ステップ 703）。なお、Scan 演算の場合は直下の演算は存在しないため、直下の演算のバッファ情報は取得しない。同様に Output 演算も自身のバッファが存在しないため、取得しない。

#### 【0054】

バッファ情報を基に対象演算の実行判定をする（ステップ 704）。

実行可能である場合、演算処理部 13 に当該演算の実行を依頼する（ステップ 705）。

実行完了後、その処理によるバッファ情報の変化を知るために、バッファ制御部 22 からバッファ情報を再度取得する（ステップ 706）。

ここで、当該演算の増量判定を行い（ステップ 707）、増量が必要な場合は、バッファ制御部 22 に増量を依頼する（ステップ 708）。

増量ができた場合（ステップ 709）、再度演算処理に戻る（ステップ 705）。

増量ができなかった場合、対象の演算を直上の演算に変更し（ステップ 710）、再度バッファ情報の取得（ステップ 703）と実行判定（ステップ 704）に戻る。

#### 【0055】

ステップ707のバッファの増量判定において、増量が不要であった場合、当該演算の終了判定を行う(ステップ711)。

終了したと判定された場合、当該演算を終了設定する(ステップ712)。ここでもし実行計画における最上位の演算であった場合、クエリ処理完了となる。

#### 【0056】

ステップ711で終了していないと判定された場合、次に処理する演算を判定(推移判定)し、対象の演算を直上の演算(ステップ710)又は直下の演算(ステップ714)に変更し、再度バッファ情報の取得(ステップ703)と実行判定(ステップ704)に戻る。

また、ステップ704の演算の実行判定において、実行しないと判定された場合は、推移判定(ステップ713)によって対象の演算を変更し(ステップ710又はステップ714)、再度バッファ情報の取得(ステップ703)と実行判定(ステップ704)に戻る。

#### 【0057】

実体化型クエリ処理方式におけるクエリ実行制御部12の実行判定(ステップ704)は、以下の条件(1)及び(2)を全て満たす場合に「実行する」と判定される。全て満たさない場合は、「実行しない」と判定される。

(1) 出力バッファが空である。

(2) 下位の出力バッファにデータがある、または、入力データがストレージで読み込みデータの残りがあ、または、制御バッファにデータがある。

#### 【0058】

実体化型クエリ処理方式におけるクエリ実行制御部12の増量判定(ステップ707)は、以下の条件(1)及び(2)を全て満たす場合に増量すると判定される。満たさない場合は、増量しないと判定される。

(1) 演算結果が格納されたバッファが上限である。

(2) 下位の出力バッファにデータがある。

#### 【0059】

実体化型クエリ処理方式におけるクエリ実行制御部12の終了判定(ステップ711)は、以下の条件(1)～(3)を全て満たす場合に「終了する」と判定される。全て満たさない場合は、「終了しない」と判定される。

(1) 下位の演算が終了している、または、入力がストレージである。

(2) 下位の出力バッファが空である、または、ストレージの読み込みデータの残りがない。

(3) 制御バッファがない、または、制御バッファが空である。

#### 【0060】

実体化型クエリ処理方式におけるクエリ実行制御部12の演算推移判定(ステップ713)は、以下の条件(1)及び(2)を全て満たす場合に「直下の演算」に推移すると判定される。全て満たさない場合は、「直上の演算」に推移すると判定される。

(1) 出力バッファが空である。

(2) 直下の演算が未終了である。

#### 【0061】

要求駆動型クエリ処理方式を採用した場合のクエリ処理制御装置におけるクエリ実行制御部12の処理手順について、図8のフローチャート図を参照して説明する。

クエリ計画部11から実行計画を取得し(ステップ801)、ツリー構造の最上位の演算から処理を開始する(ステップ802)。通常、最上位の演算は、クライアントに処理結果を渡すOutput演算であり、これを選択する。

ここで、選択された演算とその直下の演算のバッファ情報をバッファ制御部22から取得する(ステップ803)。なお、Scan演算の場合は直下の演算は存在しないため、直下の演算のバッファ情報は取得しない。同様にOutput演算も自身のバッファが存在しないため、取得しない。

10

20

30

40

50

## 【 0 0 6 2 】

バッファ情報を基に対象演算の実行判定をする（ステップ 8 0 4 ）。

実行不可である場合、直下の演算を対象の演算とし（ステップ 8 1 4 ）、出力バッファのサイズ分のデータ作成を依頼して、再度バッファ情報の取得（ステップ 8 0 3 ）と実行判定（ステップ 8 0 4 ）に戻る。

実行可能である場合、演算処理部 1 3 に当該演算の実行を依頼する（ステップ 8 0 5 ）。

実行完了後、その処理によるバッファ情報の変化を知るために、バッファ制御部 2 2 からバッファ情報を再度取得する（ステップ 8 0 6 ）。

## 【 0 0 6 3 】

ここで、当該演算の制御バッファサイズと、直下の演算の出力バッファサイズ（作成要求するデータサイズ）の増量判定を行う（ステップ 8 0 7 ）。

当該演算の制御バッファサイズに増量が必要な場合は、バッファ制御部 2 2 にこの増量を依頼する（ステップ 8 0 8 ）。

ステップ 8 0 9 で制御バッファサイズの増量ができただけの場合は、再度演算処理に戻る（ステップ 8 0 5 ）。制御バッファサイズの増量ができなかった場合は、対象の演算を直上の演算に変更し（ステップ 8 1 0 ）、再度バッファ情報の取得（ステップ 8 0 3 ）と実行判定（ステップ 8 0 4 ）に戻る。

## 【 0 0 6 4 】

ステップ 8 0 7 において直下の演算の出力バッファサイズの増量が必要な場合も、同様にバッファ制御部 2 2 にこの増量を依頼する（ステップ 8 0 8 ）。

ステップ 8 0 9 で出力バッファサイズの増量できなかった場合は、同様に直上の演算に移行する（ステップ 8 1 0 ）が、増量できた場合は、次に処理する演算を判定（推移判定）する（ステップ 8 1 2 ）。

ステップ 8 1 2 で直上の演算と判断され場合は、対象の演算を変更して（ステップ 8 1 0 ）、再度バッファ情報の取得（ステップ 8 0 3 ）と実行判定（ステップ 8 0 4 ）に戻る。

ステップ 8 1 2 で「直下の演算」と判断された場合は、「直下の演算」にデータ生成を依頼する（ステップ 8 1 4 ）。

バッファの増量判定（ステップ 8 0 7 ）において、増量が不要であった場合、当該演算の終了判定を行う（ステップ 8 1 1 ）。

## 【 0 0 6 5 】

終了したと判定された場合、当該演算を終了設定する（ステップ 8 1 3 ）。ここでもし実行計画における最上位の演算であった場合、クエリ処理完了となる。

ステップ 8 1 1 で終了していないと判定された場合、次に処理する演算を判定（推移判定）し（ステップ 8 1 2 ）、対象の演算を変更して（ステップ 8 1 0 、 8 1 4 ）、再度バッファ情報の取得（ステップ 8 0 3 ）と実行判定（ステップ 8 0 4 ）に戻る。

## 【 0 0 6 6 】

要求駆動型クエリ処理方式におけるクエリ実行制御部 1 2 の実行判定（ステップ 8 0 4 ）は、以下の条件（ 1 ）及び（ 2 ）を全て満たす場合に「実行する」と判定される。全て満たさない場合は、「実行しない」と判定される。

（ 1 ）出力バッファが上限でない。

（ 2 ）下位の出力バッファにデータがある、または、入力データがストレージで読み込みデータの残りが残っている。

## 【 0 0 6 7 】

要求駆動型クエリ処理方式におけるクエリ実行制御部 1 2 の増量判定（ステップ 8 0 7 ）は、以下の条件を全て満たす場合に、対象の演算の制御バッファを「増量する」と判定される。全てを満たさない場合は、「増量しない」と判定される。

（ 1 ）制御バッファが上限である。

（ 2 ）下位の出力バッファにデータがある。

10

20

30

40

50

一方で、以下の条件を全て満たす場合に、直下の演算の出力バッファを「増量する」と判定される。

(1) 演算結果が格納されたバッファ(出力バッファもしくは制御バッファ)が上限に達してない。

(2) 直下の演算が終了していない。

【0068】

要求駆動型クエリ処理方式におけるクエリ実行制御部12の終了判定(ステップ811)は、以下の条件(1)~(3)を全て満たす場合に「終了する」と判定される。全て満たさない場合は、「終了しない」と判定される。

(1) 下位の演算が終了している、または、入力がストレージである。

10

(2) 下位の出力バッファが空である、または、ストレージの読み込みデータの残りが無い。

(3) 制御バッファがない、または、制御バッファが空である。

【0069】

要求駆動型クエリ処理方式におけるクエリ実行制御部12の推移判定(ステップ812)は、以下の条件(1)及び(2)を全て満たす場合に直上に推移すると判定される。全て満たさない場合は、直下に推移すると判定される。

(1) 出力バッファが上限に達する(要求量のデータ作成完了する)、または、下位の演算が終了する。

(2) 出力バッファが上限に達する(要求量のデータ作成完了する)、または、出力バッファが空である。

20

【0070】

実体化型クエリ処理方式及び要求駆動型クエリ処理方式における演算処理部13の処理手順について、図9のフローチャート図を参照して説明する。

クエリ実行制御部12から演算の実行依頼を受け(ステップ91)、演算の種類に応じて処理を開始する(ステップ92)。

処理開始後、出力先のバッファ(演算の出力バッファ、もしくは制御バッファ)が上限に達するか、入力先のバッファ(演算の制御バッファ、もしくは直下の演算の出力バッファ)が空になった場合、演算処理を終了する(ステップ93)。

その後、バッファ利用量収集部23に、バッファ利用量収集部23が収集している利用状況の情報を渡す(ステップ94)。

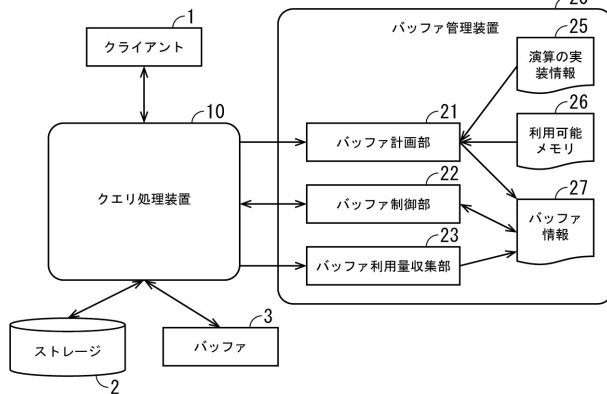
30

【符号の説明】

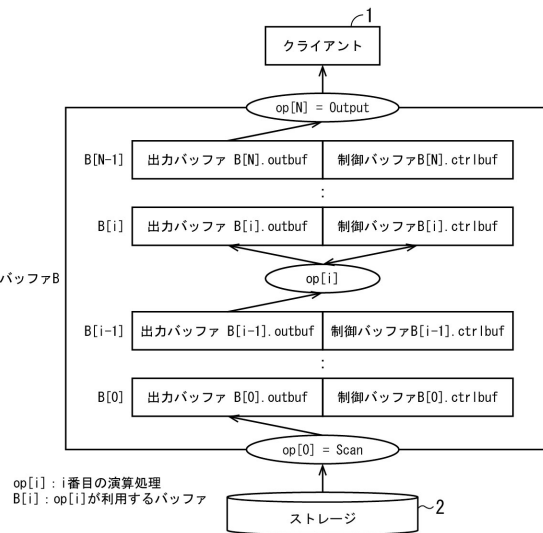
【0071】

1...クライアント、 2...ストレージ、 3...バッファ、 10...クエリ処理装置、  
11...クエリ計画部、 12...クエリ実行制御部、 13...演算処理部、 20...バッファ管理装置、  
21...バッファ計画部、 22...バッファ制御部、 23...バッファ利用量収集部、  
25...演算の実装情報、 26...利用可能メモリ、 27...バッファ情報。

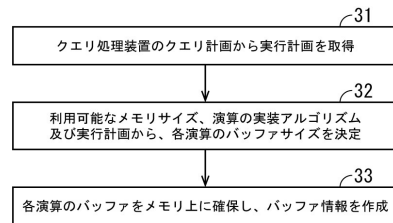
【図 1】



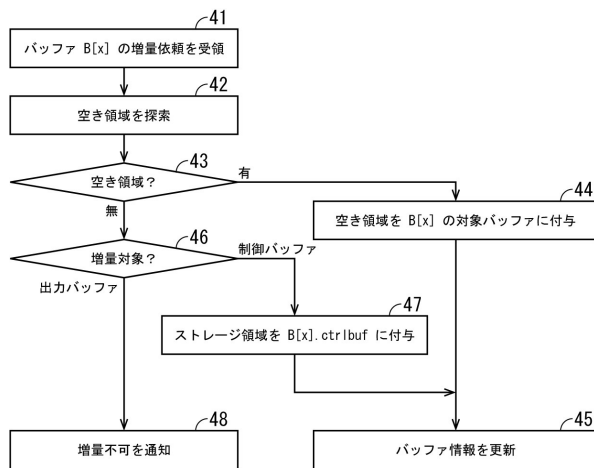
【図 2】



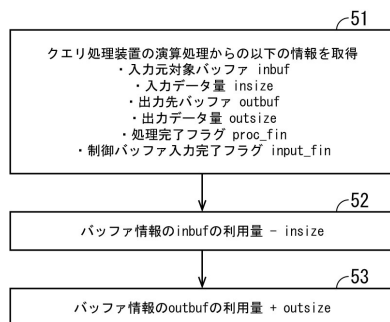
【図 3】



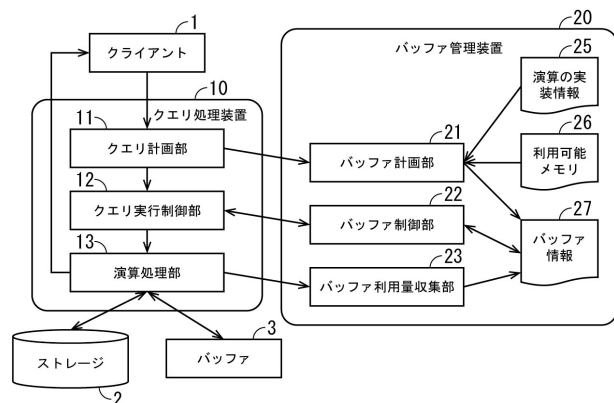
【図 4】



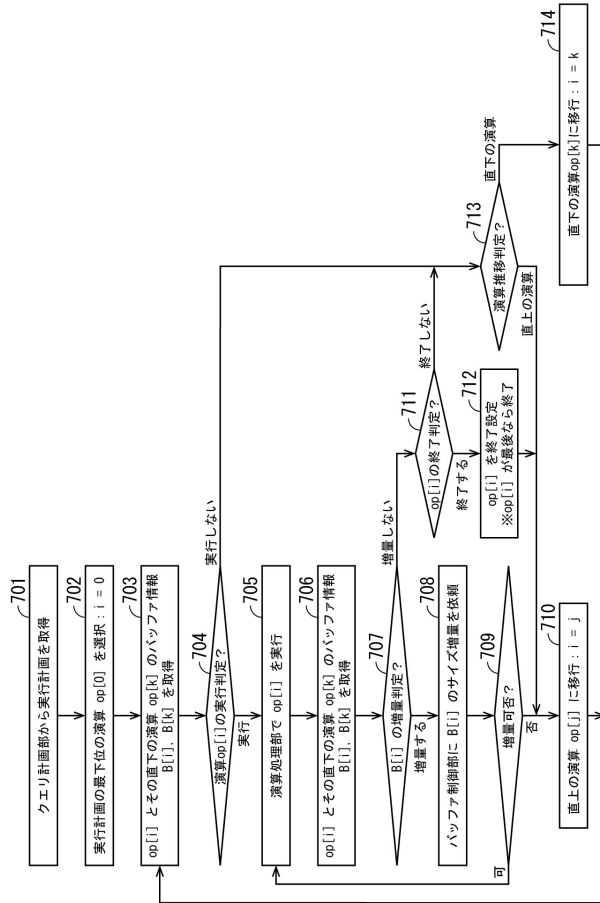
【図 5】



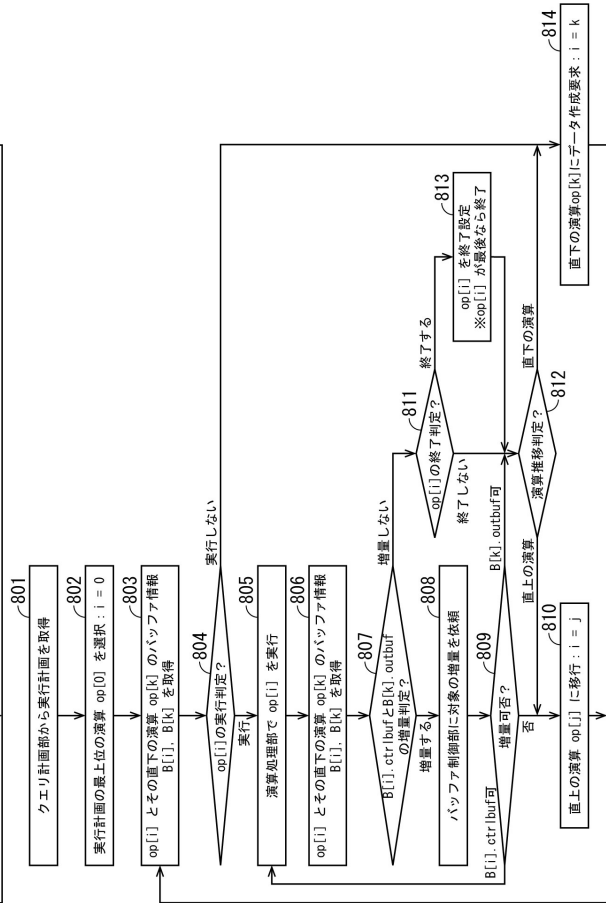
【図 6】



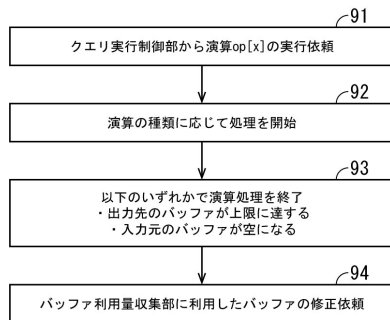
【図 7】



【図 8】



【図 9】





---

フロントページの続き

- (56)参考文献 特開平 09 - 305614 (JP, A)  
米国特許第 9104663 (US, B1)  
米国特許出願公開第 2008 / 0147599 (US, A1)  
国際公開第 2013 / 161081 (WO, A1)  
大西 元、他, “データベースの多様な応用分野を対象とする並列処理システム SMASH - メモリ資源割り当ての計算方式の実現 - ”, 第 38 回 (平成元年前期) 全国大会講演論文集 (II), 社団法人情報処理学会, 1989 年 3 月 17 日, pp. 966 - 967

- (58)調査した分野(Int.Cl., DB 名)  
G06F 16/00 - 16/958  
G06F 12/00