

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2011-39820

(P2011-39820A)

(43) 公開日 平成23年2月24日(2011.2.24)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 17/30 (2006.01)	G06F 17/30 110C	5B075
	G06F 17/30 340D	
	G06F 17/30 180D	

審査請求 未請求 請求項の数 14 O L (全 30 頁)

(21) 出願番号 特願2009-187137 (P2009-187137)
 (22) 出願日 平成21年8月12日 (2009.8.12)

(71) 出願人 000005108
 株式会社日立製作所
 東京都千代田区丸の内一丁目6番6号
 (74) 代理人 110000350
 ポレール特許業務法人
 (72) 発明者 勝沼 聡
 東京都国分寺市東恋ヶ窪一丁目280番地
 株式会社日立製作所中央研究所内
 (72) 発明者 今木 常之
 東京都国分寺市東恋ヶ窪一丁目280番地
 株式会社日立製作所中央研究所内
 (72) 発明者 藤原 真二
 東京都国分寺市東恋ヶ窪一丁目280番地
 株式会社日立製作所中央研究所内
 Fターム(参考) 5B075 KK02 QR00 UU40

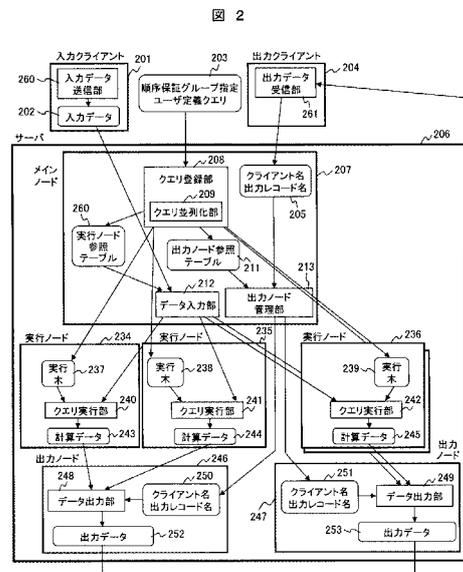
(54) 【発明の名称】 ストリームデータ処理方法及び装置

(57) 【要約】

【課題】オペレータ処理後に、一つのノードで入力時に付加したタイムスタンプに従って並び替え出力して入力された順序を保つと、スループットは低下し、レイテンシは増大する。

【解決手段】ユーザ定義クエリ203によりユーザはグループを指定する。そして、同じグループのデータ間においては入力された順序を保証する。同じグループに属するデータを同一のノード246で出力することで、必要部分に対して入力された順序を保証する。異なるグループに属するデータを別のノード247で出力することで、複数ノードでの処理を実現し性能低下を回避する。

【選択図】図2



【特許請求の範囲】**【請求項 1】**

複数のノードにより、入力されるデータに対する処理を行うストリームデータ処理方法であって、

第一のノードでクエリの登録を行い、第二のノードの集合でタイムスタンプが付加されたデータをオペレータ処理し、第三のノードの集合でオペレータ処理されたデータを前記タイムスタンプに従って並び替えて出力し、

前記第一のノードは、登録した前記クエリから、入力された順序を同じグループのデータで保証するというグループが示された第一の値式の集合を抽出し、オペレータ処理を複数の前記第二のノードでデータごとに分割する方法を示す第二の値式の集合を算出し、

前記第一及び第二の値式の集合から、出力処理を複数の前記第三のノードにデータを分けて行う方法が示された第三の値式の集合を算出し、前記第三の値式の集合から前記第三のノードの集合にデータの振分ける方法を決定し、

前記第三のノードの集合における振分け方法と前記第一の値式の集合から、同じ前記第二のノードで処理されたデータが同じ前記第三のノードで処理されるようにデータを振り分ける方法を決定する、

ことを特徴とするストリームデータ処理方法。

【請求項 2】

請求項 1 記載のストリームデータ処理方法であって、

前記第一及び第二の値式の集合における要素として共に任意の値式 F がある場合に、前記値式 F を第三の値式の集合における要素とする、

ことを特徴とするストリームデータ処理方法。

【請求項 3】

請求項 1 記載のストリームデータ処理方法であって、

第一の値式の集合の要素である任意の値式 F 、第二の値式の集合の要素である任意の値式 G に対して、 F の値が同じストリームデータは必ず G の値が同一になることが示されている場合に、第三の値式の集合における要素を G とする、

ことを特徴とするストリームデータ処理方法。

【請求項 4】

請求項 1 記載のストリームデータ処理方法であって、

前記第三の値式の集合と、前記第二のノードの数、前記第三のノードの数、ストリームデータに関する情報から、第三のノードの集合にデータを振分ける方法を決定し、

前記第三のノードの集合における振分け方法と前記第一の値式の集合、前記第二のノードの数、ストリームデータに関する情報から、第二のノードの集合のデータの振り分け方を、第三のノードの集合の振分け方と同様の方法または、第三のノードの集合より多くのノードにデータを振り分ける方法を解析し決定する、

ことを特徴とするストリームデータ処理方法。

【請求項 5】

請求項 1 に記載のストリームデータ処理方法であって、

前記第一のノードまたはクライアントのノードは、前記第三のノードの集合のデータの振り分け方法に関する情報を参照し、前記第三のノードの集合を検索し、前記第三のノードの集合にデータとクライアントに関する情報を送り、

前記第三のノードの集合は、データと前記クライアントの情報に基づいて出力データを前記クライアントに転送する、

ことを特徴とするストリームデータ処理方法。

【請求項 6】

請求項 1 に記載のストリームデータ処理方法であって、

前記第一のノードは、前記クライアントから、出力するデータ間で許容する入力された順序の時刻ずれ幅が示された場合に、異なる前記第三のノードの集合から出力されるデータに対して、前記第二のノードの集合または前記第三のノードの集合は、前記第三のノード

10

20

30

40

50

の集合から出力されるデータのタイムスタンプ値及び許容する時刻ずれ幅を用いて出力時刻を調整し、入力されたデータの順序のずれを示された範囲に抑制する、ことを特徴とするストリームデータ処理方法。

【請求項 7】

請求項 1 記載のストリームデータ処理方法であって、前記第一ノードは、複数のタイムスタンプ付加ノードを有し、同じ前記第三ノードで処理をするデータは必ず同じ前記タイムスタンプ付加ノードでタイムスタンプ付加処理をするように、複数の前記タイムスタンプ付加ノードにデータを振分ける、ことを特徴とするストリームデータ処理方法。

【請求項 8】

複数のノードにより、入力されるデータに対する処理を行うストリームデータ処理装置であって、

第一のノードでクエリの登録を行い、第二のノードの集合でタイムスタンプが付加されたデータをオペレータ処理し、第三のノードの集合でオペレータ処理されたデータを前記タイムスタンプに従って並び替えて出力し、

前記第一のノードは、登録した前記クエリから、入力された順序を同じグループのデータで保証するというグループが示された第一の値式の集合を抽出し、オペレータ処理を複数の前記第二のノードでデータごとに分割する方法を示す第二の値式の集合を算出し、

前記第一及び第二の値式の集合から、出力処理を複数の前記第三のノードにデータを分けて行う方法が示された第三の値式の集合を算出し、前記第三の値式の集合から前記第三のノードの集合にデータの振分ける方法を示す第一のテーブルを生成し、

前記第三のノードの集合における振分け方法と前記第一の値式の集合から、同じ前記第二のノードで処理されたデータが同じ前記第三のノードで処理されるようにデータを振り分ける方法を示す第二のテーブルを生成する、

ことを特徴とするストリームデータ処理装置。

【請求項 9】

請求項 8 記載のストリームデータ処理装置であって、

前記第一及び第二の値式の集合における要素として共に任意の値式 F がある場合に、前記値式 F を第三の値式の集合における要素とする、

ことを特徴とするストリームデータ処理装置。

【請求項 10】

請求項 8 記載のストリームデータ処理装置であって、

第一の値式の集合の要素である任意の値式 F 、第二の値式の集合の要素である任意の値式 G に対して、 F の値が同じストリームデータは必ず G の値が同一になることが示されている場合に、第三の値式の集合における要素を G とする、

ことを特徴とするストリームデータ処理装置。

【請求項 11】

請求項 8 記載のストリームデータ処理装置であって、

前記第三の値式の集合と、前記第二のノードの数、前記第三のノードの数、ストリームデータに関する情報から、第三のノードの集合にデータを振分ける方法を示す前記第一のテーブルを生成し、

前記第三のノードの集合における振分け方法と前記第一の値式の集合、前記第二のノードの数、ストリームデータに関する情報から、第二のノードの集合のデータの振り分け方を、第三のノードの集合の振分け方と同様の方法または、第三のノードの集合より多くのノードにデータを振り分ける方法とし、

第二のノードの集合のデータの振り分け方法を示す前記第二のテーブルを生成する、

ことを特徴とするストリームデータ処理装置。

【請求項 12】

請求項 8 記載のストリームデータ処理装置であって、

前記第一のノードまたはクライアントのノードは、前記第一のテーブルから前記第三のノ

10

20

30

40

50

ードの集合を検索し、前記第三のノードの集合にデータと前記クライアントに関する情報を送り、

前記第三のノードの集合は、取得した前記クライアントの情報に基づいてデータを前記クライアントに転送する、

ことを特徴とするストリームデータ処理装置。

【請求項 13】

請求項 11 に記載のストリームデータ処理装置であって、

前記第一のノードは、前記クライアントから、出力するデータ間で許容する入力された順序の時刻ずれ幅が示された場合に、異なる前記第三のノードの集合から出力されるデータに対して、前記第二のノードの集合または前記第三のノードの集合は、前記第三のノードの集合から出力されるデータのタイムスタンプ値及び許容する時刻ずれ幅を用いて出力時刻を調整し、入力されたデータの順序のずれを示された範囲に抑制する

ことを特徴とするストリームデータ処理装置。

【請求項 14】

請求項 8 に記載のストリームデータ処理装置であって、

前記第一ノードは、複数のタイムスタンプ付加ノードを有し、同じ前記第三ノードで処理をするデータは必ず同じ前記タイムスタンプ付加ノードでタイムスタンプ付加処理をするように、複数の前記タイムスタンプ付加ノードにデータを振り分ける方法を示すテーブルを生成する、

ことを特徴とするストリームデータ処理装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストリームデータ処理システムにおいて、データごとに複数ノードでオペレータ処理や出力処理をする技術に関する。

【背景技術】

【0002】

時々刻々と到着する大量のデータをリアルタイム処理するデータ処理システムに対する要求が高まっている。例えば、株自動売買、カーブロープ、Webのアクセス監視、製造監視などを挙げることができる。

【0003】

従来、企業情報システムのデータ管理の中心にはデータベース管理システム (Data Base Management System、以下、DBMS とする) が位置づけられていた。DBMS は、処理対象のデータをストレージに格納し、格納したデータに対してトランザクション処理に代表される高信頼な処理を実現している。しかし、DBMS では、新たなデータが到着する度に、全てのデータに対して検索処理を施すため上記のリアルタイム性の要求を満たすことは難しい。例えば、株取引を支援する金融アプリケーションを考えた場合、株価の変動にいかに関与に迅速に反応できるかがシステムの最重要課題の一つである。しかし、データの検索処理が株価変動のスピードに追いつくことができず、ビジネスチャンスを逃してしまうことになりかねない。

【0004】

このようなリアルタイムデータ処理に好適なデータ処理システムとして、ストリームデータ処理システムが提案されている。例えば非特許文献 1 にストリームデータ処理システム "STREAM" が開示されている。ストリームデータ処理システムでは従来の DBMS とは異なり、まずクエリ (問合せ) をシステムに登録しデータの到来と共に該クエリが継続的に実行される。実行されるクエリがあらかじめ把握できるため、新たなデータが到着したら、それまでの処理結果からの差分のみを処理することで高速な処理が可能である。したがって、ストリームデータ処理によって、株取引などにおける高レートで発生するデータをリアルタイムに解析し、ビジネスに有効なイベントの発生を監視して利活用することが可能になる。

10

20

30

40

50

【 0 0 0 5 】

大量のデータに対して高速処理するために、ストリームデータ処理では複数のコンピュータ（ノード）による分散処理が求められている。分散処理では、クエリを構成するオペレータごとに異なるノードで処理する方式（以下、パイプライン並列方式とする）と、同一のオペレータに対してデータごとに複数ノードで処理する方式（以下、データ並列方式とする）が知られている。特に、データ並列方式では、パイプライン並列方式と比してノード数の増加に従って通信オーバーヘッドはさほど増加しないため、スループットを大幅に向上させることが可能である。

【 0 0 0 6 】

データ並列方式においてデータの各ノードへの割り当て方法は、各オペレータの処理方法から算出される。ストリームデータ処理記述用の言語は非特許文献 2 に開示されている CQL (Continuous Query Language) 等、DBMS で広く用いられている SQL (Structured Query Language) に近い言語で書かれていることが多く、RDB に準じた方法でデータの分割方法を計算することが可能である。例えば、CQL には SQL と同様に Join (結合) オペレータや Aggregation (集計) オペレータが存在し、RDB と同様に、それぞれ結合の条件や集計の単位によってデータの分割の仕方が定まる。ストリームデータ処理システムにおけるデータ並列化方式としては、特許文献 1、特許文献 2、非特許文献 3、非特許文献 4、非特許文献 5 が開示されている。

10

【 先行技術文献 】

20

【 特許文献 】

【 0 0 0 7 】

【 特許文献 1 】 米国公開特許 US2007/0288635 号

【 特許文献 2 】 米国公開特許 US2008/0168179 号

【 非特許文献 】

【 0 0 0 8 】

【 非特許文献 1 】 R. Motwani, J. Widom, A. Arasu, B. Babcock, S. Babu, M. Datar, G. Manku, C. Olston, J. Rosenstein and R. Varma 著: "Query Processing, Resource Management, and Approximation in a Data Stream Management System", In Proc. of the 2003 Conf. on Innovative Data Systems Research (CIDR), January 2003

30

【 非特許文献 2 】 A. Arasu, S. Babu and J. Widom 著: "The CQL continuous query language: semantic foundations and query execution", The VLDB Journal, Volume 15, Issue 2, pp. 121 - 142 (June 2006)

【 非特許文献 3 】 T. Johnson, M. S. Muthukrishnan, V. Shkapenyuk, O. Spatscheck 著: "Query-aware partitioning for monitoring massive network data streams", SIGMOD, 2008

40

【 非特許文献 4 】 M. A. Shah, J. M. Hellerstein, S. Chandrasekaran, M. J. Franklin 著: "Flux: an adaptive partitioning operator for continuous query systems", ICDE, 2003

【 非特許文献 5 】 M. Ivanova, T. Risch 著: "Customizable parallel execution of scientific stream queries", VLDB, 2005

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 9 】

50

前述のように、ストリームデータ処理ではまずクエリを登録しデータの到来と共にクエリが実行される。そして、システムに到着したデータは、途中の処理に影響されことなくその到着順を守って結果が出力される。したがって、データ並列化方式では到着順を守るため、システムに到着したデータはまず、タイムスタンプが付加され、その後、複数ノードでオペレータ処理される。そして、データ出力時にタイムスタンプ値に従った並び替えが行われる。

【0010】

出力時の並び替えは、全てのデータに対して一つのノードで処理する。したがって、オペレータ処理のためのノード数を増加させても、出力時の処理によって律速されスループットを一定以上向上させることができない。また、一つのコンピュータの処理遅延によって他のコンピュータで処理した結果が遅延し、レイテンシが増大する。

10

【0011】

実際のアプリケーションでは、一部のデータの間で、入力された順序を保証する必要がある場合が多い。そのようなアプリケーションに対して必要部分で入力された順序を保ちつつ高速に実行するために、複数のクエリを定義し、手動で必要に応じて各クエリの出力データをマージする。しかし、この方式では、ユーザが各クエリの出力を把握していなければ、性能だけでなくデータ順序の整合性を保つこともできず利便性が低い。

【0012】

本願発明の目的は、順序を保証すべきデータに対して順序を保ちつつ高速にストリームデータ処理を行うストリームデータ処理方法及び装置を提供することにある。

20

【課題を解決するための手段】

【0013】

上記の目的を達成するため、本発明においては、複数のノードにより、入力されるデータに対する処理を行うストリームデータ処理方法及び装置であって、第一のノードでクエリの登録を行い、第二のノードの集合でタイムスタンプが付加されたデータをオペレータ処理し、第三のノードの集合でオペレータ処理されたデータを前記タイムスタンプに従って並び替えて出力し、第一のノードは、登録したクエリから、入力された順序を同じグループのデータで保証するというグループが示された第一の値式の集合を抽出し、オペレータ処理を複数の第二のノードでデータごとに分割する方法を示す第二の値式の集合を算出し、且つ複数の第三ノードにデータを振り分ける方法を解析し、同じ第二のノードで処理されたデータが同じ第三のノードで処理されるようにデータを振り分けるストリームデータ処理方法及び装置を提供する。

30

【0014】

すなわち、本発明では、クエリの定義において、ユーザはグループ（以下、順序保証グループと呼ぶ）を指定する。そして、同じ順序保証グループのデータ間においては入力された順序が保証される。同じ順序保証グループに属するデータを同一のノードで出力することで、必要部分に対して入力された順序を保証する。また、異なる順序保証グループに属するデータを別のノードで出力することで、複数ノードでの処理を実現し性能低下を回避する。また、ユーザが指定した順序保証グループから、自動的にクエリの出力処理方法が解析されるので実現が容易である。

40

【0015】

本発明の好適な実施態様では、ストリームデータ処理サーバにクエリの出力処理を並列化する機構（以下、クエリ並列化部と呼ぶ）を追加する。そして、順序保証グループが指定されたクエリがサーバに登録されると、クエリ並列化部が動作する。まず、クエリのオペレータをデータ並列方式で処理する方法を算出した後に、順序保証グループを読み込み、算出したオペレータの処理方法と組み合わせ計算する。そして、クエリの出力処理における複数ノードへのデータの振り分け方法を解析する。また、出力処理の振り分け方から、オペレータ処理における複数ノードへのデータの振り分け方法を解析する。

【0016】

また、出力ノードを管理する機構（以下、出力ノード管理部）をストリームデータ処理

50

エンジンに追加する。そして、クエリ実行時に、オペレータ処理の振分け方法にしたがって複数ノードでオペレータを処理すると共に、出力ノード管理部でクライアントからデータ出力要求を受け取る。そして、出力処理の振分け方法から、データを出力するノードを検索し、クライアントと出力データの情報を出力ノードに伝える。最後に、出力ノードでは対応するデータのオペレータ処理が完了するとクライアントに転送する。

【発明の効果】

【0017】

本発明により、入力された順序を必要に応じて保証しつつ、高スループットかつ低レイテンシなストリームデータ処理を実現できる。

【図面の簡単な説明】

10

【0018】

【図1】各実施例に利用される計算機システムの一例を示すブロック図である。

【図2】第1実施例の計算機システムの動作図である。

【図3】第1実施例における、クエリ定義の一例を示す図である。

【図4】第1実施例におけるクエリ登録時動作のフローチャート図である。

【図5】第1実施例における、オペレータ処理分割キー算出のフローチャートを示す図である。

【図6】第1実施例における、出力処理分割キー算出のフローチャートを示す図である。

【図7】第1実施例における実行木、実行ノード参照テーブル、出力ノード参照テーブル生成のフローチャート図である。

20

【図8】第1実施例における、クエリ登録時処理の動作例を示す図である。

【図9】第1実施例における、実行木、実行ノード参照テーブル、出力ノード参照テーブル配布処理の動作例を示す図である。

【図10】第1実施例における出力処理分割キー算出の動作例を示す図である。

【図11】第1実施例におけるクエリ実行時のフローチャートを示す図である。

【図12】第1実施例における、クエリ実行時の動作例を示す図である。

【図13】第1実施例における、クライアントの出力ノード検索処理のフローチャートを示す図である。

【図14】第1実施例における、クライアントの出力ノード検索処理の動作例を示す図である。

30

【図15】第1実施例における、実行時出力ノード変更処理を示す図である。

【図16】第1実施例における、車渋滞検出クエリ処理のクエリ定義の一例を示す図である。

【図17】第1実施例における、車渋滞検出クエリ処理のクエリ登録時の動作例を示す図である。

【図18】第1実施例における、車渋滞検出クエリ処理のクエリ実行時の動作例を示す図である。

【図19】第2実施例における、出力ノード制御方式の動作図である。

【図20】第2実施例における、実行ノード制御方式の動作図である。

【図21】第2実施例における、クエリ定義の一例を示す図である。

40

【図22】第2実施例における、出力ノード制御方式のフローチャートを示す図である。

【図23】第2実施例における、出力ノード制御方式の動作例を示す図である。

【図24】第2実施例における、実行ノード制御方式のフローチャートを示す図である。

【図25】第2実施例における、実行ノード制御方式の動作例を示す図である。

【図26】第3実施例における、計算機システムの動作図である。

【図27】第3実施例における、フローチャートを示す図である。

【図28】第3実施例における、クエリ登録時の動作例を示す図である。

【図29】第3実施例における、クエリ実行時の動作例を示す図である。

【図30】第3実施例における、クライアントによる振分け処理のフローチャートを示す図である。

50

【図 3 1】第 3 実施例における、クライアントによる振分け処理の動作例を示す図である。

【発明を実施するための形態】

【0019】

以下、図面に基づき本発明の種々の実施形態を説明する。

【0020】

図 1 は各実施例が適用される計算機システムの一構成を示す図である。メモリ 150 上に入力データ送信部 260 を含む入力クライアントノード 201、メモリ 151 上に出力データ受信部 261 を含む出力クライアント 204 は、ネットワーク 127 を介してストリームデータ処理システムが稼動するストリームデータ処理サーバ 206 に接続されている。ネットワーク 127 は、イーサネット（登録商標）、光ファイバなどで接続されるローカルエリアネットワーク（LAN）、もしくは LAN よりも低速なインターネットを含んだワイドエリアネットワーク（WAN）でも差し支えない。また、クライアント計算機 201、204 は、パーソナルコンピュータ（PC）、ブレード型の計算機システムなどの任意のコンピュータシステムでよい。

10

【0021】

ストリームデータ処理サーバ 206 は、メインノード 207、複数の実行ノード 234、複数の出力ノード 246 で構成され、各ノードは外部ネットワーク 127、もしくは、内部ネットワークを介して接続されている。ノード 207、234、246 は、インタフェース部を構成する I/O インタフェース 136 ~ 138、それぞれ処理部を構成する中央処理部（CPU）130 ~ 132、記憶部であるメモリ 139 ~ 141 が、バスで結合された計算機であり、ブレード型計算機システム、PC サーバなどの任意のコンピュータシステム、データ処理装置でよい。

20

【0022】

ノード 207、234、246 はそれぞれ I/O インタフェース 136 ~ 138 を介してクライアント計算機 201、204 にアクセスする。ノード 207、234、246 で、ストリームデータ処理結果、処理の中間結果、システム動作に必要な設定データを不揮発性のストレージに格納する場合には、それぞれ記憶部であるストレージ装置 133 ~ 135 を用いることができる。ストレージ装置 133 ~ 135 は、それぞれ I/O インタフェース 136 ~ 138 を介して直接接続される、もしくはそれぞれ I/O インタフェース 136 ~ 138 よりネットワークを介して接続される。ノード 207、234、246 の各々のメモリ 139 ~ 141 に、ストリームデータ処理をする各ストリームデータ処理エンジン（ソフトウェア）103、115、120 がマッピングされている。

30

【0023】

なお、本明細書において、メインノード 207、複数の実行ノード 234、複数の出力ノード 246 をそれぞれ第一のノード、第二のノード、第三のノード、あるいは第一のノードの集合、第二のノードの集合、第三のノードの集合と呼ぶ場合がある。また、以上の説明では、サーバ 206 に関して複数のノードで成り立つマルチコンピュータの構成を示したが、サーバの構成については一つのノードに複数の CPU で構成されるマルチプロセッサコンピュータに複数のストリームデータ処理エンジン 103、115、120 が動作しても構わない。また、マルチコンピュータとマルチプロセッサコンピュータを組み合わせた構成でも構わない。

40

【0024】

以下でストリームデータ処理エンジン 103、115、120 の各モジュールを説明する。ストリームデータ処理エンジン 103 には、クエリ並列化部 209 を含むクエリ登録部 208、出力ノード参照テーブル格納領域 109、実行ノード参照テーブル格納領域 110、データ入力部 212、出力ノード管理部 213 がある。通常、ストリーム処理サーバ 206 の第一のノード 207 のストリームデータ処理エンジン 103 において、入力したデータに対してタイムスタンプ付加の処理を行うが、ここでは説明の煩雑さのため説明を省略する。一方、ストリームデータ処理エンジン 115 には、クエリ実行部 240、実

50

行木格納領域 117 がある。ストリームデータ処理エンジン 120 には、データ出力部 248 がある。なお、各モジュールの動作については後述する。

【実施例 1】

【0025】

以下、第 1 実施例について図面に基づき説明する。

【0026】

図 2 は第 1 実施例の動作を説明するための図である。アプリケーションによっては、クエリが処理するデータに対し必ずしも入力された順序を保証することが求められない。第 1 実施例では、そのような場合に、あらかじめ指定させたグループが同じであるデータのみ、入力された順序を保証することでスループットの向上及びレイテンシの削減を実現する。

10

【0027】

まず、ユーザによって定義されたクエリ 203 はメインノード 207 が受信する。ユーザ定義クエリ 203 はデータの処理方法の記述に加え、処理するデータが属する順序保証グループが指定される。同じ順序保証グループに属するデータは入力された順序が保証される。メインノード 207 ではクエリ登録部 208 がクエリの各オペレータを解釈すると共に、クエリ登録部 208 のクエリ並列化部 209 は、順序保証グループ及び、クエリのデータ並列実行方法を用いて、オペレータ処理及び出力処理を複数のノードで行うためのデータの振分け方法を解析する。そして、得られたデータ振分け方法から、実行木 237 ~ 239、実行ノード参照テーブル 260、及び出力ノード参照テーブル 211 を生成する。実行木 237 ~ 239 はそれぞれ実行ノード 234 ~ 236 において処理するオペレータを記述したもので、出力ノード参照テーブル 211 はデータのレコード名をインデックスとして、各データを出力するノード名を参照するテーブルである。実行ノード参照テーブル 260 はデータのレコード名をインデックスとして、各データをオペレータ処理するノード名を参照するテーブルである。クエリ登録部 208 は実行木 237 ~ 239 をそれぞれ実行ノード 234 ~ 236 の実行木格納領域に格納し、出力ノード参照テーブル 211 をメインノード 207 の出力ノード参照テーブル格納領域に格納する。また、実行ノード参照テーブル 260 をメインノード 207 の実行ノード参照テーブル格納領域に格納する。

20

【0028】

30

そして、入力クライアント 201 の入力データ送信部 260 からメインノード 207 のデータ入力部 212 へ入力データ 202 を転送することで登録したクエリの実行を開始する。データ入力部 212 は入力データ 202 にタイムスタンプを付加し、データを実行ノード参照テーブル 260 に従って実行ノード 234 ~ 236 に振り分ける。実行ノード 234 ~ 236 はそれぞれクエリ実行部 240 ~ 242 がデータを受信し、実行木 237 ~ 239 で指定されたオペレータで処理する。そして、オペレータ処理により実行ノード 234、235 でそれぞれ生成された計算データ 243、244 を出力ノード 246 に転送し、実行ノード 236 で生成された計算データ 245 を出力ノード 247 に転送する。

【0029】

40

なお、実行ノード 234 ~ 236 で生成された計算データ 243 ~ 245 は、実行木によっては出力ノードに直接転送されない。実行木 237 によっては、計算データ 243 を実行ノード 235、236 の実行ノードに転送し、クエリ実行部 241、242 でオペレータ処理する。そして、オペレータ処理により生成された計算データを出力ノード 246、247 に転送する。計算データ 244、245 についても同様に実行ノード 234 ~ 236 に転送される。

そして、最終的に出力ノード 246、247 が受信した計算データ 243 ~ 244、245 をそれぞれデータ出力部 248、249 で、計算データに付加されたタイムスタンプに従って並び替え、それぞれ出力データ 252、253 を生成する。

【0030】

一方、出力データ 252、253 を取得する出力クライアント 204 は、取得する出力

50

データのレコード名（例えば、クエリ301ならば銘柄名としてA電、B銀、H電等）と出力クライアントの名前205をメインノード207に転送する。そして、メインノード207の出力ノード管理部213は、クエリ登録時に生成した出力ノード参照テーブル211によりレコード名に対応するデータを出力するノードの名前を検索する。そして、検出された出力ノード246、247にそれぞれ出力先のクライアントの名前と出力レコード名250、251を転送する。出力ノード250、251は、それぞれ出力データ252、253のレコード名と出力レコード名250、251とのマッチングをし、一致していたら出力クライアント204にそれぞれ出力データ252、253を転送する。

【0031】

第1実施例では、既存のデータ並列方式で出力ノード処理がボトルネックになっている場合にスループットを向上させる。例えば、図18に示すように、実行ノード数6、出力ノード数2で処理する時、一つの出力ノードにおける処理可能な入力レートをT（ダブル/秒）とすると、既存のデータ並列方式では、全体で処理可能な入力レートもT（ダブル/秒）になるのに対し、第1実施例では2T（ダブル/秒）となる（一つの処理ノードの処理可能な入力レートがT/3（ダブル/秒）以上の場合）。

10

【0032】

続いて、第1実施例の動作の詳細を述べる。

【0033】

ユーザによるクエリの定義方法について説明する。図3はユーザが定義したクエリの一例である。本実施例では、順序保証グループが示された値式の集合をクエリに対して指定させる。指定の方法は、クエリに記述しても、設定ファイルを介して指示されていてもその他の方法でも構わない。図3では、ストリームの定義において、ストリームデータ処理記述言語（CQL）に新たに追加した構文 Order Partition Byを用いて値式 銘柄が指定されている。値式の指定は、ストリームの定義以外にもクエリ定義やその他の箇所で行われても構わない。

20

【0034】

クエリ301では、銘柄、取引所、注文数、口座、業種をカラムとするストリーム 株注文を定義する。ストリーム 株注文は株取引に関する情報を配信しており、クエリq1はストリーム 株注文を入力し、銘柄及び取引所ごとの前一分間の合計注文数を常時出力する。

30

【0035】

クエリ301を定義したユーザは、同じ銘柄に対しては取引の順序にしたがって集計結果を得る必要があるが、異なる銘柄に対して取引の順序を要求していないとする。この時、ユーザは、入力ストリーム 株注文の定義において、構文Order Partition Byを用いてストリームのカラム 銘柄を指定することで、入力された順序を異銘柄に対して保証する必要がないことをシステムに伝える。

【0036】

ユーザによって定義されたクエリの登録時動作について述べる。図4はメインノードにおけるクエリ登録時のフローチャートである。最初に、ユーザ定義クエリを読み込み、全てのオペレータを複数のノードにデータを分けて並列処理する方法を示す値式の集合（以下、オペレータ処理分割キーと呼ぶ）を算出する。値式はストリームのカラムを項とする演算式であり、後述する値式も同様である（400）。

40

【0037】

そして、算出したオペレータ処理分割キーと、Order Partition By構文で指定された値式の集合（以下、OrderPartitionByキーと呼ぶ）から、オペレータ処理及び出力処理におけるデータの振分け方法を解析する。振分け方法は、1．出力処理においてデータを複数のノードに振り分ける。2．同じノードでオペレータ処理されたデータを必ず同じノードで出力処理する、を満たすことで、高スループットな出力処理を実現する。なお、アプリケーションが要求するスループットを満たすならば、同じノードでオペレータ処理されたデータを異なるノードで出力処理する場合があっても構わない。

50

【 0 0 3 8 】

動作としては、まず、オペレータ処理分割キーとOrderPartitionByキーから出力時の処理を複数のノードにデータを分けて並列に行う方法が示された値式の集合（以下、出力処理分割キーと呼ぶ）を算出する（401）。そして、算出したオペレータ処理及び、出力処理分割キーから、各々の処理における各ノードへのデータの振分け方を解析する。なお、データの振分け方法の解析は、数値解析による方法など、出力処理分割キーの算出をしない方法でも構わない。そして、その振分け方から実行木、実行ノード参照テーブル、出力ノード参照テーブルを生成し各ノードに配信する（402）。

【 0 0 3 9 】

図5は図4の400の動作詳細を示すフローチャートである。まず、ユーザ定義クエリを読み込み（501）、クエリを構成する各オペレータが複数のノードにデータを分けて並列処理する方法を示す値式の集合（以下、処理分割キーと呼ぶ）を抽出する（502）。処理分割キーの抽出方法は既存のRDB(Relational Data Base)に従い、オペレータの処理方法から抽出する。例えば、オペレータJoinなら結合の等価条件で指定された値式を処理分割キーの要素とし、オペレータaggregationならGroup By構文で指定された値式を処理分割キーの要素とする。

10

【 0 0 4 0 】

そして、各オペレータの処理分割キーの積集合Opをオペレータ処理分割キーとする（503）。Opが空集合であれば、クエリのオペレータの処理を各ノードで並列に実行することが不可能であることを通知し登録動作を終了する。または、オペレータ処理において計算データをノード間で再割当てすることを許可する場合には、一部のオペレータの積集合をオペレータ処理分割キーとして抽出してもよい。なお、オペレータ処理分割キーは上記の方法で算出する以外にも、クエリを定義したユーザなどによって指定されたキーを抽出するなど、その他の方法で導出されても構わない。

20

【 0 0 4 1 】

図6は図4の401の動作詳細を示すフローチャートである。まず、Order Partition By構文で指定した値式の集合（OrderPartitionByキー）Orを抽出する（601）。そして、オペレータ処理分割キーOp、OrderPartitionByキーOrから出力処理分割キーOuを算出する。まず、Op、Orに対し、 $F(c) \quad Op \quad F(c) \quad Or$ （ $F(c)$ は、ストリームの任意のカラムcを項とする演算式）を満たす $F(c)$ が存在する場合は（602）、 $F(c)$ をOuの要素にする。条件を満たす $F(c)$ が複数存在する場合には全ての $F(c)$ を要素とする（603）。また、 $c/n \quad Op \quad c/m \quad Or$ （ n, m は整数）を満たす c が存在する場合には（604）、 $c / s(n, m)$ （ $s(n, m)$ は n と m の最小公倍数）をOuの要素にする。条件を満たす c が複数存在する場合には全ての $c / s(n, m)$ を要素とする（605）。

30

【 0 0 4 2 】

さらに、ユーザによって $G(c2) \text{ including } F(c1)$ （ $F(c1)$ 、 $G(c2)$ はストリームの任意のカラム $c1$ 、 $c2$ を項とする値式）という関係が指定されている場合には、 $F(c1) \quad Op \quad G(c2) \quad Or$ を満たすならば、 $G(c2)$ をOuの要素にする。 $G(c2) \text{ including } F(c1)$ とは、値式 $F(c1)$ が同じストリームデータは必ず値式 $G(c2)$ が同一になることを示す。条件が複数存在するならば、全ての $G(c2)$ を要素とする（606、607）。そして、Ouが空集合でなければ（608）、Ouを出力処理の分割キーとして算出する（610）。Ouが空集合ならば出力処理のデータ並列化が不可能であることを通知し登録動作を終了する（609）。401の動作は、統計的解析など前述の方法以外でも、値式 $F(c1)$ が同じストリームデータは必ず値式 $G(c2)$ が同一になることを示されれば、 $G(c2)$ をOuの要素にして構わない（ $F(c1)$ 、 $G(c2)$ はストリームの任意のカラム $c1$ 、 $c2$ を項とする値式で、 $F(c1) \quad Op \quad G(c2) \quad Or$ を満足）。

40

【 0 0 4 3 】

図7は図4の402の動作詳細を示すフローチャートである。まず、ユーザから指定された実行ノード数 n 、出力ノード数 m 、値式 $F(c)$ （ $F(c)$ は、ストリームの任意のカラム c を項とする演算式）とそのとりうる値（または、値の数）を読み込む（701）。なお

50

、実行ノード数及び、出力ノード数はOS等のシステム情報から取得するなど、その他の方法で取ってきて構わない。また、 $F(c)$ とそのとりうる値もデータから統計的に分析する等、その他の方法で取得しても構わない。そして、算出したオペレータ処理及び出力処理の分割キーを併せて用いることで、各ノードへの振り分け方法を決定し、実行木及び実行ノード参照テーブル、出力ノード参照テーブルを生成する。

まず、出力処理分割キーの要素である値式のとりうる値の数が、実行ノード数 n よりも大きい場合には、データを n 個の出力ノードに振り分ける方法を決定し、出力ノード参照テーブルを生成する。振り分け方法は、値式の大小や辞書順で一定の領域に区切って振り分ける方法や、 mod (剰余演算)等のハッシュ関数で値式からハッシュ値を求め、ハッシュ値に従って振り分ける方法や、その他、データの偏りを考慮した振り分け方法など任意の方法で構わない。そして、オペレータ処理におけるデータの振り分け方も出力処理の振り分け方と同様の方法にし、実行ノード参照テーブルを生成する。また、同じ実行木を n 個、生成する(704)。こうして、データをオペレータ処理した後、同じノードで続けて出力処理をするようにデータを振り分けることで、オペレータ処理及び出力処理の並列化によって高スループット化しつつ、オペレータ処理と出力処理の間におけるノード間通信を発生させず低レイテンシ化する。

【0044】

また、出力処理分割キーの要素である値式のとりうる値の数が、実行ノード数 n よりも小さい場合には、データを最大で m 個の出力ノードに振り分ける方法を決定し、出力ノード参照テーブルを生成する(702)。振り分け方法は、値を一定の領域に区切って振り分ける方法や、ハッシュ関数でハッシュ値を求め、ハッシュ値に従って振り分ける方法や、その他、任意の方法で構わない。そして、出力処理の振り分け方法では同じ出力ノードに振り分けられるデータに対して、オペレータ処理では複数の実行ノードに振り分けることで、オペレータ処理で最大 n 個の実行ノードに振り分けるよう方法を解析する。具体的には、オペレータ処理分割キーの要素である値式を用いて複数の領域に分ける、または、値式からハッシュ値を求め、ハッシュ値で分けるなどして、複数の実行ノードに振り分ける方法を解析する。そして解析結果から実行ノード参照テーブルを生成する。また、実行ノード数と同数の実行木を生成する(705)。一般的にオペレータ処理は出力処理よりも計算処理コストが大きいので、オペレータ処理を出力処理よりも処理するノードを増やすことでスループットが向上する。しかし、そのためにオペレータ処理と出力処理の間に通信が発生しレイテンシが増大する。スループットの向上よりもレイテンシを低くすることを優先するアプリケーションなどでは、オペレータ処理を出力処理と同様にデータを振り分けるなどその他の方法でも構わない。

最後に、生成された実行木と出力ノード参照テーブル、実行ノード参照テーブルを各ノードに転送する(706)。

【0045】

図8、図9は、クエリ301を例にしたクエリ登録時の動作である。クエリ301はオペレータ range 、 group by 、 istream から構成される(801)。オペレータ range 、 istream については全てのデータの処理を並列に実行可能であり、したがって、処理分割キーの要素として任意の値式を取る。また、オペレータ group by は、集計の単位より、ストリーム株取引のカラム 銘柄または取引所が異なっているデータに対しては並列に実行可能である。したがって、オペレータ group by の処理分割キーは{銘柄、取引所}(802)となる(502の動作)。したがって、各オペレータの処理分割キーの積集合を取ることで、オペレータ処理分割キーは{銘柄、取引所}(803)と算出される(503、503の動作)。

【0046】

また、構文 $\text{Order Partition By}$ より、 OrderPartitionBy キーは{銘柄}(804)と抽出される(601の動作)。したがって、オペレータ処理分割キーと OrderPartitionBy キーに共通の要素としてカラム 銘柄を含む値式 銘柄が存在するので、出力処理分割キーは{銘柄}(805)と算出される(602、603の動作)。

10

20

30

40

50

【 0 0 4 7 】

そして、算出した出力処理分割キー、ユーザから指定された実行ノード数 2 (8 0 6)、出力ノード数 2 (8 1 0)、出力処理分割キーの要素である「銘柄」のとりうる値の数 1 0 4 3 から出力処理の振分け方法を決定する (7 0 1 の動作)。まず、実行ノード数よりも「銘柄」のとりうる値の数のほうが大きいので、2つの実行ノードにデータを振分ける。振分け方法としては、出力分割キーの要素である「銘柄」に対しハッシュ関数 (文字列をビット列で表現した値に対し、2で除算した剰余を返す関数) で求めたハッシュ値を用いる。ハッシュ値が0であるデータをノード1に、ハッシュ値が1であるデータをノード2に振分ける。そして、その振り分け方法を出力ノード参照テーブル 8 0 9 として実現する (7 0 2 の動作)。

10

さらに、出力処理を2個 (実行ノード数) のノードに振分けることが可能であったので (7 0 3 の動作)、オペレータ処理のデータの振分け方法を、出力処理のデータの振り分け方法と同様にし、実行ノード参照テーブルを生成する。また、2個 (実行ノード数) の実行木 (8 0 8) を生成する (7 0 5 の動作)。

【 0 0 4 8 】

生成された出力ノード参照テーブル 9 0 6 は、ストリームデータ処理サーバ 9 0 1 のメインノード #0 (9 0 2) の出力ノード参照テーブル格納領域 9 0 5 に格納され、実行ノード参照テーブル 9 2 3 は、実行ノード参照テーブル格納領域 9 2 2 に格納される。また、実行木 9 2 0、9 2 1 はそれぞれ実行・出力ノード # 1 (9 0 7)、# 2 (9 1 0) の実行木格納領域 9 1 8、9 1 9 に格納される (7 0 6 の動作)。

20

【 0 0 4 9 】

クエリ 3 0 1 は1つのクエリで構成されていたが、実際のストリームデータ処理のクエリでは複数のクエリが連なっていることが一般的である。本実施例では、複数のクエリが連なっている場合にも同様の方法でオペレータ処理分割キーの算出が可能である。また、クエリ 3 0 1 では、オペレータ処理分割キーと OrderPartitionBy キーの要素にストリームの共通のカラムを含む値式が存在した。しかし、たとえ共通のカラムを含む値式が存在しなかったとしても出力処理分割キーの算出が可能である。以下は、図 1 0 のクエリ 1 0 0 1 を例にした出力処理分割キー算出の動作である。

【 0 0 5 0 】

図 1 0 に示すように、クエリ 1 0 0 1 は、クエリ q2、q3 から成る。クエリ q2 はクエリ q1 と同じくストリーム 株注文を入力とし、取引所、銘柄ごとの分単位の合計注文数を出力する。そして、クエリ q3 はクエリ q2 で出力した銘柄毎の取引所ごとの合計注文数から銘柄毎の最大注文数を出力する。クエリ q2 ではオペレータ group by から処理分割キーが { 取引所、銘柄 } (1 0 0 2) と抽出され、クエリ q3 では同じくオペレータ group by により処理分割キーが { 銘柄 } (1 0 0 3) と抽出される。したがって、クエリ q2、q3 の処理分割キーの積集合を取ることによりオペレータ処理分割キーは { 銘柄 } (1 0 0 4) と算出される (5 0 2、5 0 3 の動作)。

30

【 0 0 5 1 】

また、クエリ 1 0 0 1 では、構文 Order Partition by によりストリームのカラム 業種を指定している。したがって、OrderPartitionBy キーは { 業種 } (1 0 0 5) となり (6 0 1 の動作)、オペレータ処理分割キー { 銘柄 } とは同じカラムを含む値式は存在しない。しかし、クエリ 1 0 0 1 ではさらに構文 including で { 銘柄 } が指定されることで、{ 銘柄 } が同じストリームデータは必ず { 業種 } が同じになることを示している。したがって、出力処理分割キーとして { 業種 } (1 0 0 6) が算出される (6 0 6 の動作)。

40

【 0 0 5 2 】

クエリ実行時の動作詳細を示す。図 1 1 は実行時動作のフローチャートである。本実施例では、出力ノードが複数存在するため、出力クライアントが取得するデータが出力されるノードをあらかじめ把握できない。したがって、以下の動作をストリームデータ処理システムに追加する。

【 0 0 5 3 】

50

まず、出力クライアント204は、取得するデータのレコード名及びクライアント名をメインノード207に送信する(1101)。そして、メインノード207では、クエリ登録時に格納された出力ノード参照テーブル211によりレコード名から、対応するデータを出力するノードを検索する。(1102)。さらに、メインノード207では検出した出力ノードにレコード名と出力クライアント名を送信する(1103)。出力ノードでは、オペレータ処理された出力データのレコード名とメインノードから受信したレコード名を比較する(1104)。そして、レコード名が一致していたならば、メインノード207から転送された名前のクライアントにデータを出力する(1104)。

【0054】

図12は図3に示したクエリ301を例にしたクエリ実行時の動作である。クエリ301に対して入力クライアント201からデータが入力されると実行が開始される。まず、サーバ901のメインノード#0(902)で入力データ1206-1209にタイムスタンプを付加し(1210)、付加されたデータを、実行ノード参照テーブル(出力ノード参照テーブル)906に従って処理・出力ノード#1(907)、#2(910)に送信する。今、銘柄A電のハッシュ値を0、銘柄H電、B銀のハッシュ値を1とする。その時、銘柄A電のデータ1223は実行・出力ノード#1(907)に送信し、銘柄H電、B銀のデータ1213は実行・出力ノード#2(910)に送信する。そして、実行・出力ノード#1、#2では受信されたデータを各々実行木909で指定されたオペレータ(range、group by、istream)で処理する。

【0055】

一方、出力クライアント1.2.3.4(IPアドレス)(1205)では取得するデータのレコード名A電、H電、B銀及びクライアント名1.2.3.4(1219)をメインノード#0に送信する(1101の動作)。そして、メインノード#0では出力ノード参照テーブル906を参照することで、受信したレコード名のデータを出力するノードを検索する(1211)。

【0056】

すなわち、出力ノード参照テーブル906から銘柄名A電のデータは処理・出力ノード#1、銘柄名H電、B銀のデータは実行・出力ノード#2から出力されると検出される(1102の動作)。そして、各出力ノードに、出力クライアント名1.2.3.4と各々のノードに対応する銘柄名を送信する。実行・出力ノード#1は「銘柄A電」(1220)が送信され、実行・出力ノード#2には「銘柄B銀、H電」(1222)が送信される(1103の動作)。

【0057】

実行・出力ノード#1ではオペレータを処理(1217)後に、メインノード#0からの情報1220に基づいて、銘柄名A電のデータ(1221)が出力クライアント1.2.3.4に出力される。また、#2では、オペレータを実行(1218)後、#0からの情報1222に基づいて、銘柄名B銀、H電のデータ(1223~1225)が出力クライアント1.2.3.4に出力される(1104の動作)。

【0058】

続いて、第1実施例における拡張方を述べる。まず、クライアント側の出力ノード管理について説明する。基本動作では出力ノード参照テーブル211をメインノード207に格納し、その出力ノード参照テーブル211を用いて出力ノードを検索した。しかし、出力ノード導出処理はサーバのノードではなく出力クライアントで行っても構わない。ストリームデータ処理においては出力先のクライアントが多数存在することも多く、そのような場合にはサーバで出力ノードを検索するのと比べて、新たにサーバ・ノードを追加することなしに出力ノード管理を並列実行させることができる。図13はその拡張方式のフローチャートである。

【0059】

まず、図13の上段において、クエリの実行開始時に出力クライアント204はメインノード207に出力ストリームの名前を送信する(1301)。そして、メインノード2

10

20

30

40

50

07では、受信した出力ストリーム名に対応する出力ノード参照テーブルを出力クライアントに送信する(1302)。最後に、出力クライアントでは受信した出力ノード参照テーブルを出力データ受信部に格納する(1303)。

【0060】

次に、図13下段に示すクエリ実行時に、出力クライアントは出力ノード参照テーブルを参照し、取得をするデータのレコード名から出力ノードを検索する(1304)。そして、対応する出力ノードにデータのレコード名とクライアント名を送信する(1305)。出力ノードでは、オペレータ実行されたデータのレコード名と、出力クライアントから受信したレコード名を比較し一致すると、そのデータを出力クライアントに出力する(1306)。

10

【0061】

図14は、クエリ301の例を用いたクライアントによる出力ノード管理動作である。

【0062】

まず、図14上段に示すクエリの実行開始時に、出力クライアント1.2.3.4(1205)は出力ストリーム名q1(1402)をメインノード#0(902)に転送する(1301の動作)。#0ではq1の登録時に、出力ストリーム名q1(1402)に対応する出力ノード参照テーブルを生成し出力ノード参照テーブル格納領域に格納する。その格納されているテーブルを1.2.3.4に転送する(1302の動作)。

【0063】

そして、図14下段に示すクエリ実行時に、1.2.3.4は出力ノード参照テーブル906を参照し、取得する銘柄A電のデータの出力ノードを#1、同じく取得する銘柄H電、B銀のデータの出力ノードを#2と検出する(1304の動作)。そして、1.2.3.4はクライアント名「1.2.3.4」の情報(1411、1412)と共に、#1に「銘柄A電」の情報(1411)、そして、#2に「銘柄H電、B銀」の情報(1412)を送信する(1305の動作)。#1、#2では、入力データに対してオペレータを処理後、#1は「銘柄A電」の情報(1411)より銘柄名A電のデータ(1409)を1.2.3.4に出力し、#2は「銘柄H電、B銀」の情報(1412)より銘柄名H電、B銀のデータ(1410)を1.2.3.4に出力する(1306の動作)。

20

【0064】

図15は、実行時に出力ノードを変更する処理に関する拡張方式である。ストリームデータ処理システムでは、実行時にデータが入力される。そして、データの入力レートが大きくなる、または、入力されるデータの種類の偏りが生じる等でノードが過負荷になるため、実行時に出力ノードにおけるデータの振分けの仕方を変更する必要性が生じる。

30

【0065】

クライアントに出力ノード参照テーブルを配置する場合の実行時出力ノード変更方法を述べる。図15上段のフローにおいて、データの割振り方法を変更する時に、まず、メインノードにある出力ノード参照テーブルが更新され、その更新された出力ノード参照テーブルとその出力ストリームの名前を出力クライアントに送信する(1501)。そして、出力クライアントでは、受信した出力ストリーム名に対応する出力ノード参照テーブルを受信したテーブルに置換える(1502)。

40

【0066】

図15下段により、クエリ301を例にした実行時出力ノード変更処理の動作を述べる。変更前、#1に銘柄%2が0であるデータ、#2に銘柄%2が1であるデータが振分けられている。そして、変更によって新たに出力ノード#3が追加され、#2に銘柄%4が1であるデータ、#3に銘柄%4が3であるデータが割当てられたとする。その時、メインノード#0(902)で更新された出力ノード参照テーブル1506を出力ストリーム名q1(1402)と共に出力クライアント1.2.3.4(1205)に転送する(1501の動作)。そして、1.2.3.4では出力ストリーム名q1の出力ノード参照テーブル906を更新された出力ノード参照テーブル1506に置き換える(1502の動作)。その後は、1.2.3.4では置換された出力ノード参照テーブルを参照することで

50

、取得データの出力ノードを把握する。

【0067】

続いて、車渋滞検出クエリの例を用いて第1実施例の動作とその効果を述べる。特に、最小公倍数を用いて出力処理分割キーを算出する動作(605、606の動作)を説明する。

【0068】

図16は車渋滞検出クエリである。クエリでは、カラムとして車両ID、x、yを持つストリーム 車両情報1601を入力とする。車両IDは車両ごとに割当てられており、x、yはその各車両の位置情報を表しており、ストリーム 車両情報では定期的に各車両の位置情報が配信される。クエリでは各々の車両の現在と一つ前の位置情報(x、y)から車両速度を計算する(1604)。また、複数の車両の位置情報により各々の車両の周りの車両密度を計算する(1605)。そして、車両速度と車両密度から各々の車両周辺の渋滞量を計算する(1606)。

10

【0069】

上記のクエリに対し、ユーザは、近くを走行中の車両同士の渋滞情報は時刻順に更新されることを求めているとする。すなわち、入力ストリーム 車両情報(1601)に対して、構文 Order Partition Byで $(int)x/300$ を指定し、 $(int)x/300$ ごと、すなわち、xがそれぞれ0~299、300~599の範囲でそれぞれ渋滞情報を順序通りに処理することをシステムに伝える。

【0070】

20

図17はクエリ1602を用いた登録時の動作例である。

【0071】

まず、オペレータ処理分割キーを算出する。クエリ1602は、車両速度計算(1604)、車両密度計算(1605)、渋滞量計算(1606)に関するクエリで構成され、それぞれのクエリの処理分割キーは、クエリの各オペレータの処理分割キーの積集合を取ること、{車両ID, $(int)x/100$ } (1701)、{ $(int)x/100$ } (1702)、{(任意)} (1703)となる(502、503の動作)。したがって、オペレータ処理分割キーは各クエリの処理分割キーの積集合{ $(int)x/100$ } (1704)となる(503、504の動作)。

次に、出力処理分割キーを算出する。まず、OrderPartitionByキーが、Order Partition By構文より $\{(int)x/300\}$ (1705)と求まる(601の動作)。そして、オペレータ処理分割キー $\{(int)x/100\}$ とOrderPartitionByキー $\{(int)x/300\}$ は、ともに入力ストリームのカラムxを含む値式 $(int)x/100$ 、 $(int)x/300$ を要素に持ち、かつ、共に x/n (nは整数)の形式を取ることで、100と300の最小公倍数をとることにより出力処理の分割キーは $\{(int)x/300\}$ (1706)と算出される(604、605の動作)。

30

【0072】

そして、各処理キーからデータの割り当て方法を解析する。まず、ユーザによって指定された実行ノード数6 (1711)、出力処理分割キーの要素の値式である $(int)x/300$ の取りうる値数2 (1708)という情報より(701の動作)、実行ノード数の方が値式の取りうる値の数よりも大きいので、最大で出力ノード数2個(1707)のノードに振り分ける方法を決定し、出力ノード参照テーブル1710を生成する(702の動作)。そして、オペレータ処理において最大6個(実行ノード数)に振り分ける方法を解析する。今、オペレータ処理分割キー $\{(int)x/100\}$ では最大6ノードにデータを割り当てるのが可能のため、オペレータ処理を6ノードにデータを振り分ける方法を決定し、実行ノード参照テーブル1712を生成する。また、実行木1709を6個(実行ノード数分)生成する(705の動作)。最後に、生成した実行木1709を実行ノード#1~#6に配布し、出力ノード参照テーブル1710及び実行ノード参照テーブル1712をメインノード#0に配布する(706の動作)。

40

【0073】

図18はクエリ1602を用いた実行時の動作例である。前述したようにクエリ160

50

2では、#1～6(1802-1807)をオペレータ処理するノードとし、#7(1808)、#8(1809)を、オペレータ処理結果を出力するノードとする。メインノード#0で、実行ノード参照テーブル1709を用いて振り分けられたデータは、#1～6で車両速度計算、車両密度計算(1810)、渋滞量計算(1811)に関するクエリによって実行される。例えば、#3ではxが200～300の範囲の車両情報データが入力され、1812のように渋滞情報データが出力される。また、#4ではxが300～400の範囲のデータが処理され1815のように出力される。そして、#1～3で出力された渋滞情報データは#7に送信され、#4～6の出力データは#8に送信される。#7、#8ではその送信データを入力順に並び替え各々クライアントに出力する。例えば、#7では#1～#3の最新データのタイムスタンプ値を比較し(1816)、最もタイムスタンプが古いデータ(1817)をクライアントに出力する(1818)。また、#8では#4～#6の最新データを比較することで(1819)、同じく最もタイムスタンプ値が古いデータ(1820)をクライアントに出力する(1821)。

10

20

30

40

50

【実施例2】

【0074】

次に第2の実施例について述べる。第1実施例では、入力された順序を保証するのを指定した順序保証グループ内に限定することで高スループット、低レイテンシを実現した。したがって、ノード間の負荷の違い等によって異なるデータグループでは時刻が大幅にずれることがある。しかし、アプリケーションによっては入力された順序を厳密に守る必要はないものの、大きくずれることは好まれないことがある。例えば、株取引の集計に関するクエリでは、複数のユーザが複数の銘柄情報を参考に投資行動を決めることを想定すると、銘柄別に出力時刻が大幅にずれることは、正確性、公平性の観点から回避が必要である。A電株がB銀株よりも情報が遅く入った場合には、A電株を見ている投資家はB銀株を見ている投資家より行動が遅れる。また、A電株とB銀株を比べて有利な行動を選択しようとする投資家は判断を誤る。また、車渋滞検出のクエリにおいては、指定した順序保証グループの境界(クエリ1602なら $x=300$)で順序が保たれない。したがって、例えば、データグループの境界付近を走行中の車両は渋滞情報が時系列順に表れないことがあるため、カーナビで表示する場合に見にくくなるだけでなく、走行経路の選択を誤ることもつながる。

【0075】

そこで、第2実施例では、順序保証グループが異なるデータ間に対しても時刻のずれ幅の一定以内に抑制する。図19、図20はそれぞれ第2実施例の構成を示す。図19は出力ノードでずれ幅を制御する方式を表しており、図20は処理ノードでずれ幅を制御する方式を表している。出力ノード制御方式ではデータ出力時、処理ノード246、247におけるストリームデータ処理エンジンのデータ処理部248、249に出力時刻調整部1901、1902を追加し、お互いの出力データの時刻を交換することで出力時刻を調整する。一方、実行ノード制御方式では、実行ノード235のストリームデータ処理エンジンに追加した出力時刻調整部2001で調整する。つまり、クエリ実行部241で処理された計算データ244を出力ノード246と出力247の両方に、送信時刻を調整しつつ転送する。

【0076】

第2実施例では、時刻のずれ幅を一定以内に抑制しつつ、既存のデータ並列方式よりもスループットを向上させる。例えば、実行ノード数6、出力ノード数2で処理する時(図23、図25)、一つの出力ノードにおける処理可能な入力レートを T (タプル/秒)とすると、既存のデータ並列方式では、全体で処理可能な入力レートが T (タプル/秒)になるのに対し、第2実施例では $1.5T$ (タプル/秒)となる(一つの処理ノードの処理可能な入力レートが $T/3$ (タプル/秒)以上の場合)。

【0077】

以下、第2実施例の動作の詳細を述べる。第2実施例では、Order Partition By構文で指定したデータグループが異なるデータに対して許容する時刻ずれ幅を指定させるために

、ストリームデータ処理記述言語に新たに構文 limit を追加する。例えば、図 2 1 におけるクエリ 2 1 0 1 では、構文 Order Partition By で { 銘柄 } をキーに指定し limit 構文で 1 秒と指定することで、異なる銘柄間においても出力時刻のずれが 1 秒以内となるようにする。許容する時刻ずれ幅の指定は、順序保証グループの指定と同様、ストリームの定義以外にもクエリ定義やその他の箇所で行われても構わない。Limit 構文は、クエリ登録時に OrderPartitionBy 構文と共に解釈され、実行木や出力ノード参照テーブルと同様に、許容時刻ずれ幅に関する情報を該当するノードに転送する。

【 0 0 7 8 】

図 2 2 は出力ノード制御方式のフローチャートである。まず、実行ノードで第 1 実施例と同様にオペレータを処理する。そして、処理完了後 (2 2 0 1)、計算データを出力ノードに送信する (2 2 0 2)。一方、出力ノードではお互いに出力データのタイムスタンプ値を交換する。そして、他の出力ノードからタイムスタンプ値を受信する (2 2 0 3)、または、実行ノードから計算データを受信したら (2 2 0 3)、各実行ノードから送信された最新データの時刻を比較し最古のデータを抽出する (2 2 0 7)。そして、抽出データのタイムスタンプ値と、他ノードから転送された最新のタイムスタンプ値に許容時間ずれ幅を加算した時刻を比較する。そして、抽出データのタイムスタンプ値のほうが古ければ (2 2 0 4)、クライアントにデータを出力する (2 2 0 6)。またその際に、他の出力ノードに出力したデータのタイムスタンプ値を通知する (2 2 0 5)。

10

【 0 0 7 9 】

図 2 3 はクエリ 2 1 0 1 を用いた出力ノード制御方式の動作例である。# 1 ~ 6 (2 3 0 2 ~ 2 3 0 7) はオペレータを処理するノードであり、# 7, 8 (2 3 0 8、2 3 0 9) は処理結果を出力するノードである。# 1 ~ 3 でオペレータ処理された計算データは # 7 に転送され、# 4 ~ 6 で処理されたデータは # 8 に転送される。

20

【 0 0 8 0 】

今、# 3 では銘柄 A 電のデータ (2 3 1 0) が処理され (2 2 0 1 の動作)、その処理結果が # 7 に出力される (2 2 0 2 の動作)。また、# 4 では銘柄 H 電、B 銀のデータ (2 3 1 1) が処理され (2 2 0 1 の動作)、その処理結果が # 8 に出力される (2 2 0 2 の動作)。また、# 7 ではデータ 2 3 2 1 を出力し、データ 2 3 2 1 のタイムスタンプ値は 0 ' 2 であることから、# 7 は # 8 にタイムスタンプ値 0 ' 2 (2 3 2 7) を転送した (2 2 0 3 の動作)。一方、# 8 ではデータ 2 3 2 2 を出力し、データ 2 3 2 2 のタイムスタンプ値は 0 4 であることから # 8 は # 7 にタイムスタンプ値 0 ' 4 (2 3 1 5) を転送した (2 2 0 3 の動作)。

30

【 0 0 8 1 】

そして、# 7 は # 1 ~ # 3 から送信された計算データを保持するキュー 2 3 4 0 ~ 2 3 4 2 の先頭データのタイムスタンプ値を比較し、最も古いデータ (2 3 1 4) を出力データとする (2 2 0 7 の動作)。一方、# 8 は # 4 ~ # 6 に対応するキュー 2 3 4 3 ~ 2 3 4 5 の先頭データのタイムスタンプ値を比較し、最も古いデータ (2 3 1 8) を出力データとする (2 2 0 7 の動作)。

【 0 0 8 2 】

そして、# 7 では、出力データ 2 3 1 4 のタイムスタンプ値 0 3 と、# 8 から受け取ったタイムスタンプ値 0 ' 4 (2 3 1 5) に許容時刻 1 0 を足した値を比較する (2 3 1 2)。そして、出力データ 2 3 1 4 のタイムスタンプ値のほうが古かったので出力データ 2 3 2 0 をクライアントに転送する (2 2 0 4、2 2 0 6 の動作)。そして、出力データ 2 3 2 0 のタイムスタンプ値 0 3 を # 8 に転送する (2 2 0 5 の動作)。

40

【 0 0 8 3 】

一方、# 8 では、出力データ 2 3 1 8 のタイムスタンプ値 1 ' 5 と、# 7 から受け取ったタイムスタンプ値 (2 3 2 7) に許容時刻 1 0 を足した値を比較する (2 3 1 9)。そして、出力データ 2 3 1 8 のタイムスタンプ値のほうが新しかったので出力データ 2 3 1 8 はクライアントに出力されない (2 2 0 4 の動作)。

【 0 0 8 4 】

50

次に、本実施例における、実行ノード出力制御方式の動作を述べる。図24は実行ノード出力制御方式のフローチャートである。実行ノードでは第1実施例と同様にオペレータを処理する。そして、オペレータ処理完了後(2401)、出力ノードに計算データを送信する(2405)。その時に、実行ノードに共有キューが存在するノードでは(2402)、第1実施例とは異なる動作を取る。共有キューが存在する実行ノードでは複数の出力ノードへ転送する。その時に、共有キューでは複数の出力ノードへ転送するデータを保持する。共有キューを持つノードは、データ出力時に計算データのタイムスタンプ値が、他のキューの先頭データのタイムスタンプ値に許容時間を加算した時刻よりも古ければ出力ノードにデータを送信し(2403)、新しければデータを送信せずキューにデータを保留する(2404)。そして、出力ノードに保留中と通告する(2404)。そして、データを受信した(2406)出力ノードでは、他の実行ノードが保留中でなければ(2407)、実行ノードから送信された最新データの時刻を比較し最古のデータを出力する(2408)。

【0085】

図25はクエリ2101を用いた実行ノード出力制御方式の動作例である。図23と同様に、#1~6(2302~2307)は実行ノードであり、#7,8(2308、2309)は出力ノードである。#1~3、#4~6で処理されたデータはそれぞれ#7、#8に転送される。例えば、#3では銘柄A電のデータ(2510)が処理され、#7に転送され(2517)クライアントに出力される(2321)。また、#4では共有キュー2512を持ち、銘柄H電、B銀のデータ(2513、2514)を処理する。処理された計算データは、#8に転送されクライアントに出力する(2322)。また、同時に#7に計算データのタイムスタンプ値を転送する。

【0086】

今、#4で処理された計算データ2514を共有キュー2512に投入した。その時、#4では、#7へ転送されるデータを保持するキューの先頭データ2513のタイムスタンプ値と、#8へ転送されるデータを保持するキューの先頭データ2514のタイムスタンプ値に許容時刻ずれ幅10を足した値を比較する(2511)。そして、2513のタイムスタンプ値のほうが古かったので#7にタイムスタンプ値2518を転送する。#7では#1~#3からの計算データを保持するキュー(2340~2342)の先頭データ(2517)とタイムスタンプ値2518を比較する。そして、データ2517が最古データなので、クライアントに出力する(2320)。

【0087】

また、#4では、計算データ2514を共有キュー2512に投入した時に、#8へ転送されるデータを保持するキューの先頭データ2514のタイムスタンプ値と、#8へ転送されるデータを保持するキューの先頭データ2513のタイムスタンプ値に許容時刻ずれ幅10を足した値を比較する(2515)。そして、2514のタイムスタンプ値のほうが新しかったので、#8に保留中(2343)と通告し、#8からは、#5、#6からの計算データを保持するキュー(2344、2345)のデータは出力されない。

【実施例3】

【0088】

最後に、タイムスタンプ付加を複数のタイムスタンプ付加ノードで行う、第3の実施例を説明する。図26は第3の実施例の動作を説明するための図である。第1及び第2の実施例においては、説明は省略したが1つのノードでタイムスタンプ付加の処理をしていたため、ノード数が増加した場合にはタイムスタンプ付加処理が過負荷になりスループットが一定以上向上しない可能性があった。そこで、本実施例では、タイムスタンプ付加を第一のノードに追加された複数のタイムスタンプ付加ノードで処理することでスループットをさらに向上させる。

【0089】

第3実施例では、図26に示すようにストリームデータ処理サーバ206において、ストリームデータ処理エンジンにタイムスタンプ付加部2605、2606を含むタイムス

10

20

30

40

50

タイムスタンプ付加ノード 2603、2604 が第一のノードの集合として追加される。また、メインノード 207 のストリームデータ処理エンジンにタイムスタンプ付加ノード管理部 2602 がある。

【0090】

動作としては、まず、メインノード 207 で読み込まれたユーザ定義クエリ 203 に対し、タイムスタンプ付加処理を分散するよう指定されていたならば、クエリ並列化部 209 で、クエリ 203 のタイムスタンプ付加処理を複数のノードで行うためのデータの振分け方を解析する。そしてデータの振分け方法からタイムスタンプ付加ノード参照テーブル 2601 を生成する。タイムスタンプ付加ノード参照テーブル 2601 は、データのレコード名をインデックスとしてそのデータのタイムスタンプを付加するノード名を参照するテーブルである。

10

【0091】

そして、クエリの実行時に、入力データ 202 がメインノード 207 に取り込まれた時に、そのタイムスタンプ付加ノード管理部 2602 でタイムスタンプ付加ノード参照テーブル 2601 を参照し、入力データ 202 のレコード名からデータのタイムスタンプを付加するノード名を検索し、タイムスタンプ付加ノード 2603、2604 に入力データを割当てる。そして、タイムスタンプ付加ノード 2603、2604 のタイムスタンプ付加部 2605、2606 でそれぞれデータにタイムスタンプを付加し、実行ノード 234、235 に転送する。転送後、第 1 及び第 2 実施例と同様にオペレータを処理する。

【0092】

20

第 3 実施例では、第 1 実施例でタイムスタンプ付加処理がボトルネックになっている場合にスループットを向上させる。例えば、第 1 実施例で、実行ノード数 4、出力ノード数 2 で処理する時 (図 12)、メインノードにおける処理可能な入力レートを T (ダブル/秒) とすると、第 1 実施例では、全体で処理可能な入力レートが T (ダブル/秒) になる (一つの処理・出力ノードの処理可能な入力レートが $T/2$ (ダブル/秒) 以上の場合)。一方、第 3 実施例で、タイムスタンプ付加ノード数 2 で処理する時 (図 29)、全体で処理可能な入力レートは $2T$ (ダブル/秒) となる (一つの処理・出力ノードの処理可能な入力レートが $T/2$ (ダブル/秒) 以上の場合)。

【0093】

第 3 実施例の動作詳細を述べる。図 27 はクエリの登録時及び実行時のフローチャートである。

30

【0094】

最初に、図 27 上段に基づき、クエリ登録時の動作を述べる。クエリにはあらかじめタイムスタンプ付加を複数ノードで行うか否かについて印付けされている。印付けは、新たな構文を用いても、設定ファイルに記述しても、その他の方法を用いてもよい。そして、タイムスタンプ付加の分散が指定されていたならば、まず、第 1 実施例と同様に出力処理におけるデータの振分け方法を解析 (2701) 後、同じノードで出力処理をするデータは必ず同じノードでタイムスタンプの付加をように、タイムスタンプ付加処理のデータの振分け方法を解析する (2702)。そして、タイムスタンプ付加処理のデータの振分け方からタイムスタンプ付加ノード参照テーブルを生成しメインノードに配信する (2703)。

40

【0095】

次に、第 27 図下段によりクエリ実行時の動作を述べる。まず、クライアントから入力されたデータをメインノードが受信し (2705)、タイムスタンプ付加ノード参照テーブルにより、入力データのレコード名からタイムスタンプ付加するノード名を検索し、対応するタイムスタンプ付加ノードに入力データを送信する (2706)。そして、各タイムスタンプ付加ノードでは受信したデータにタイムスタンプを付加し (2707)、実行ノードにタイムスタンプを付加したデータを送信する (2708)。その後、実行ノードでは、受信したデータに対し第 1 実施例と同様にオペレータの処理をする。

【0096】

50

図 28 はクエリ 301 を用いたクエリ登録時の動作例である。まず、第 1 実施例と同様に出力処理におけるデータの振分け方としてカラム 銘柄によって振分けの方法 2801 を決める (2701 の動作)。そして、出力処理を同じノードとするデータは必ずタイムスタンプの付加を同じノードできるように、2802 のように出力処理と同様にカラム 銘柄によるデータ振分け方法に決定する (2702 の動作)。そして、タイムスタンプ付加処理のデータの振分け方からタイムスタンプ付加ノード参照テーブル 3509 を生成する (2703 の動作)。最後に、第 1 実施例と同様に生成した実行木や出力ノード参照テーブルと共にタイムスタンプ付加ノード参照テーブル 3509 を該当するノードに転送する (2703 の動作)。

【0097】

10

図 29 は、クエリ 301 を用いたクエリ実行時の動作例である。まず、メインノード 902 でタイムスタンプ付加ノード参照テーブル 2901 を用いて、#1 (907)、#2 (910) にそれぞれ銘柄 % 2 が 0 であるデータ (銘柄 A 電 (2903))、銘柄 % 2 が 1 であるデータ (銘柄 B 銀、H 電のデータ (2902)) を振り分ける (2706 の動作)。そして、#1、#2 でそれぞれ振り分けられたデータにタイムスタンプを付加し (2904、2906) (2707 の動作)、付加されたデータに対して、第 1 実施例と同様にそれぞれ実行木 2905、2907 を処理し、処理されたデータを出力クライアントに転送する。

【0098】

20

第 3 実施例ではサーバ上のメインノードでデータのタイムスタンプ付加ノードへの振分けをする。しかし、サーバでなくクライアントによって入力データをあらかじめ振り分けることも可能である。ストリームデータ処理においては入力先のクライアントが多数存在することもあり、そのような場合にはサーバによる振分けと比較して新たに振分けのためのノードを追加することなしに、スループットを向上させることができ有用である。入力クライアントによる振分け動作について図 30 のフローチャートを用いて説明する。

【0099】

まず、同図上段に示すクエリの実行開始時の動作を述べる。メインノードではクエリ登録時にタイムスタンプ付加ノード参照テーブルが格納される。まず、入力クライアントはメインノードに入力ストリーム名を送信する (3001)。メインノードは、受信した入力ストリームに対応するタイムスタンプ付加ノード参照テーブルを入力クライアントに送信する (3002)。そして、入力クライアントは受信したタイムスタンプ付加ノード参照テーブルを入力データ送信部に格納する (3003)。

30

【0100】

次に、同図下段に示すクエリ実行時の動作を述べる。まず、入力クライアントは、データ入力時に、実行開始時に取得したタイムスタンプ付加ノード参照テーブルから入力データの振り分け先を検索し (3004)、各タイムスタンプ付加ノードに入力データを送信する (3005)。そして、各タイムスタンプ付加ノードでは受信データに対してタイムスタンプを付加する。以降の動作は、第 3 実施例に従う。

【0101】

入力クライアントによる割当て方式の動作例についてクエリ 301 を用いて説明する。図 31 はクエリ実行開始及び実行時の動作例を示す図である。

40

【0102】

まず、同図上段のクエリ実行開始時の動作例を説明する。まず、入力クライアント 1 . 2 . 3 . 5 (IP アドレス) (3101) は入力ストリーム名 株注文 (3103) をメインノード #0 (902) に転送する (3001 の動作)。そして、#0 では株注文に対応するタイムスタンプ付加ノード参照テーブル 906 を 1 . 2 . 3 . 5 に転送する (3002 の動作)。

【0103】

そして、同図下段に示すクエリ実行時に、入力クライアントではタイムスタンプ付加ノード参照テーブル 906 を参照し、入力データを振り分ける。すなわち、銘柄 % 2 が 0 で

50

あるデータ(3111)はタイムスタンプを付加するノード#1(3107)に、銘柄%2が1であるデータ(3110)は、同じくタイムスタンプを付加するノード#2(3108)に割当てられる(3004、3005の動作)。そして、#1、#2では、それぞれ銘柄%2が0であるデータ(3112)、銘柄%2が1であるデータ(3113)にタイムスタンプを付加しオペレータを処理する(3006の動作)。

【産業上の利用可能性】

【0104】

本発明は、ストリームデータ処理システムにおいて、データごとに複数ノードでオペレータ処理や出力処理をする技術として有用である。

【符号の説明】

10

【0105】

- 103 ... ストリームデータ処理エンジン(メイン)
- 109 ... 出力ノード参照テーブル格納領域
- 110 ... 実行ノード参照テーブル格納領域
- 115 ... ストリームデータ処理エンジン(実行部)
- 117 ... 実行木格納領域
- 120 ... ストリームデータ処理エンジン(出力部)
- 127 ... ネットワーク
- 130、131、132 ... CPU
- 133、134、135 ... ストレージ
- 136、137、138 ... I/Oインタフェース
- 139、140、141、150、151 ... メモリ
- 201 ... 入力クライアント
- 202 ... 入力データ
- 203 ... 順序保証グループ指定ユーザ定義クエリ
- 204 ... 出力クライアント
- 205 ... クライアント名、出力レコード名
- 207 ... メインノード
- 208 ... クエリ登録部
- 209 ... クエリ並列化部
- 211 ... 出力ノード参照テーブル
- 260 ... 実行ノード参照テーブル
- 212 ... データ入力部
- 213 ... 出力ノード管理部
- 216 ... ストリームデータ処理サーバ
- 234、235、236 ... 実行ノード
- 246、247 ... 出力ノード
- 237、238、239 ... 実行木
- 240、241、242 ... クエリ実行部
- 243、244、245 ... 計算データ
- 248、249 ... データ出力部
- 250、251 ... クライアント名、出力レコード名
- 252、253 ... 出力データ
- 260 ... カデータ送信部
- 261 ... 出力データ受信部
- 1901、1902、2001 ... 出力時刻調整部
- 2601 ... タイムスタンプ付加ノード参照テーブル
- 2602 ... タイムスタンプ付加ノード管理部
- 2603、2604 ... タイムスタンプ付加ノード
- 2605、2606 ... タイムスタンプ付加部。

20

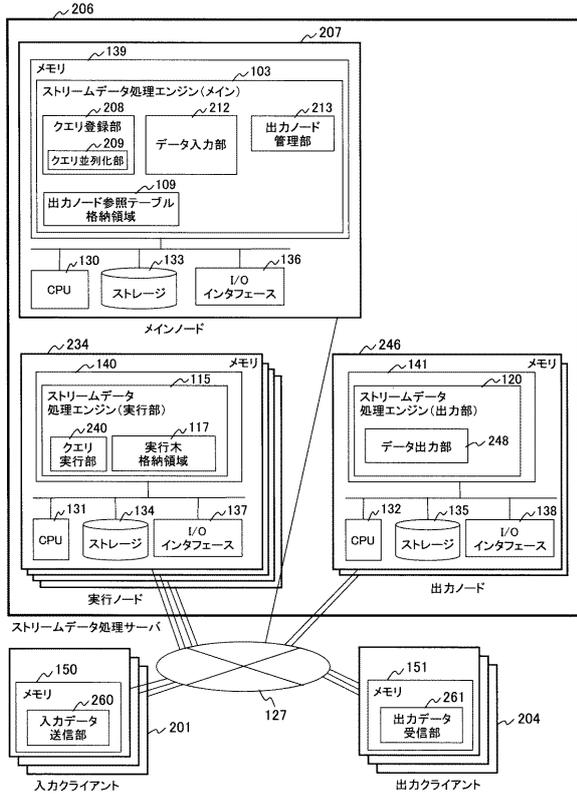
30

40

50

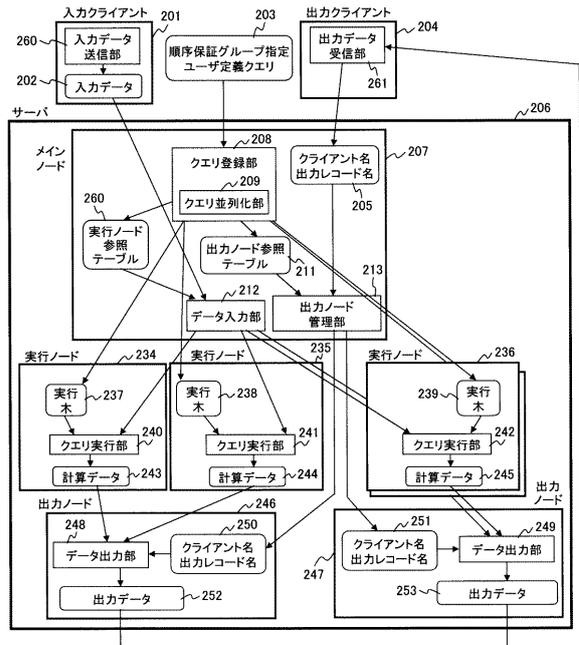
【 図 1 】

図 1



【 図 2 】

図 2



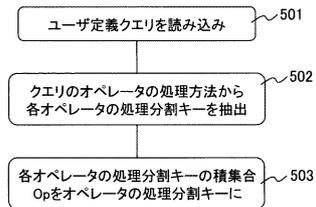
【 図 3 】

図 3



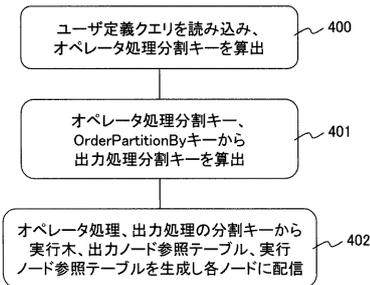
【 図 5 】

図 5

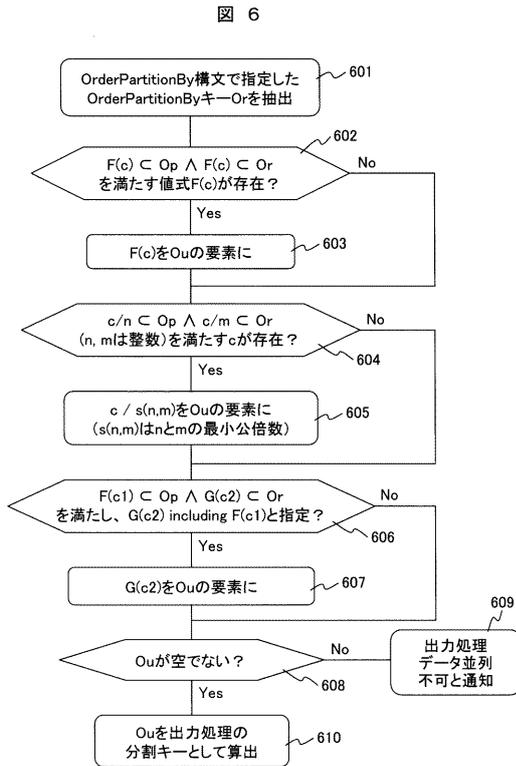


【 図 4 】

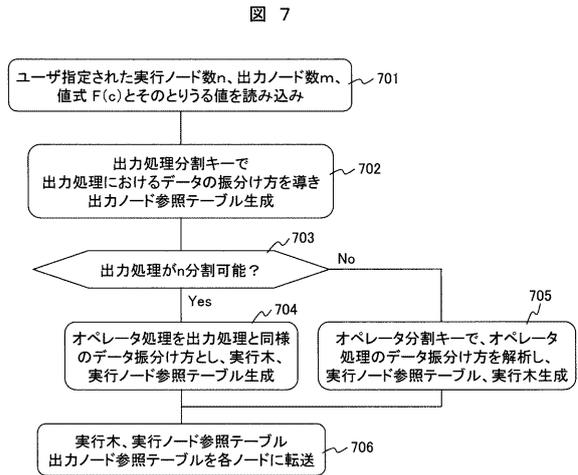
図 4



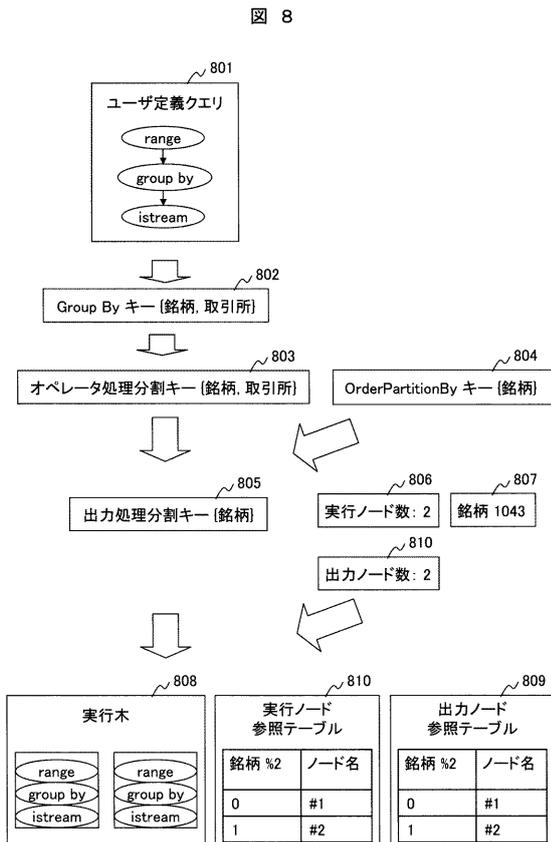
【 図 6 】



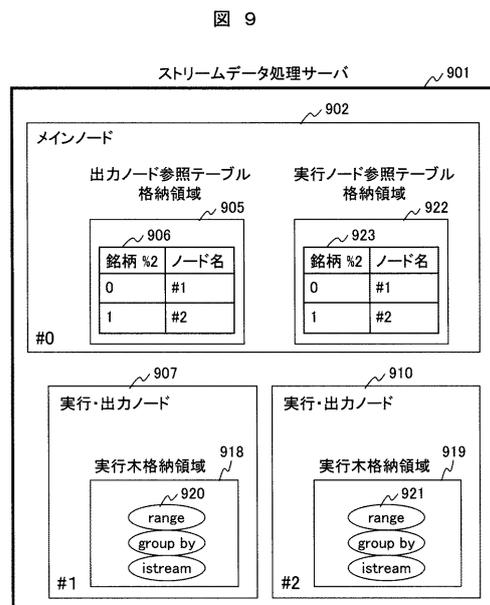
【 図 7 】



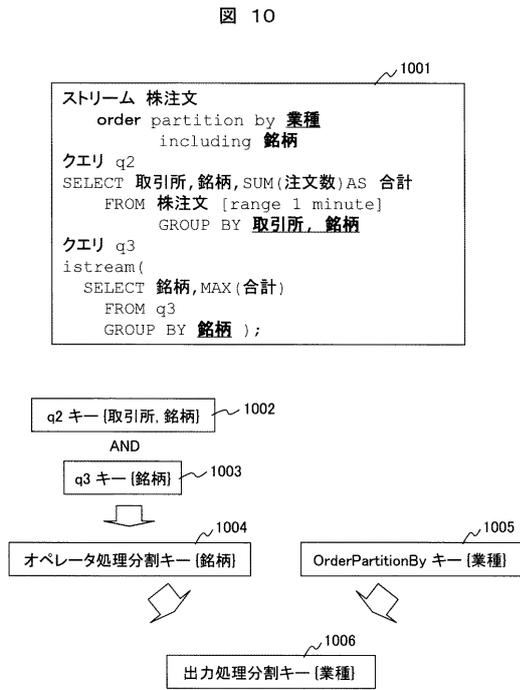
【 図 8 】



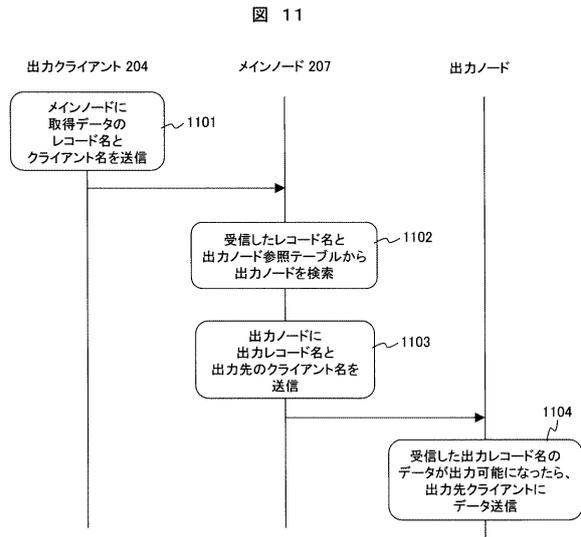
【 図 9 】



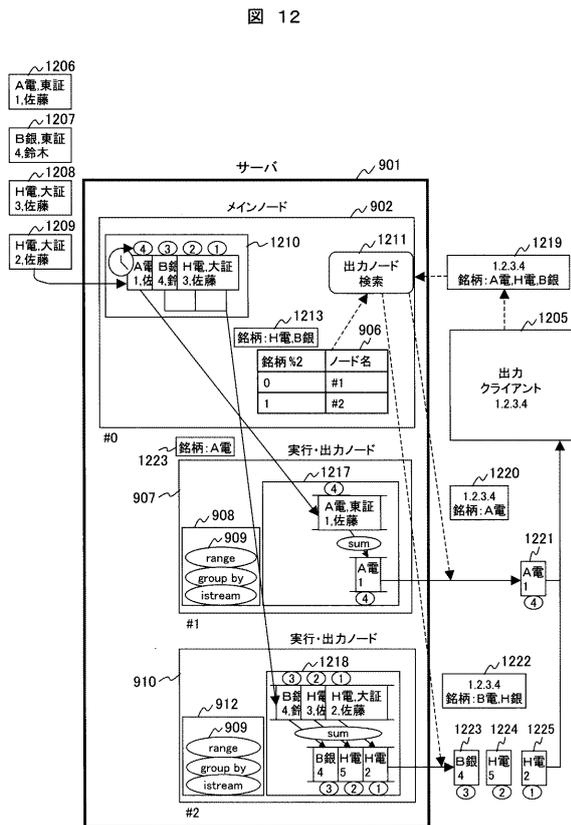
【 図 1 0 】



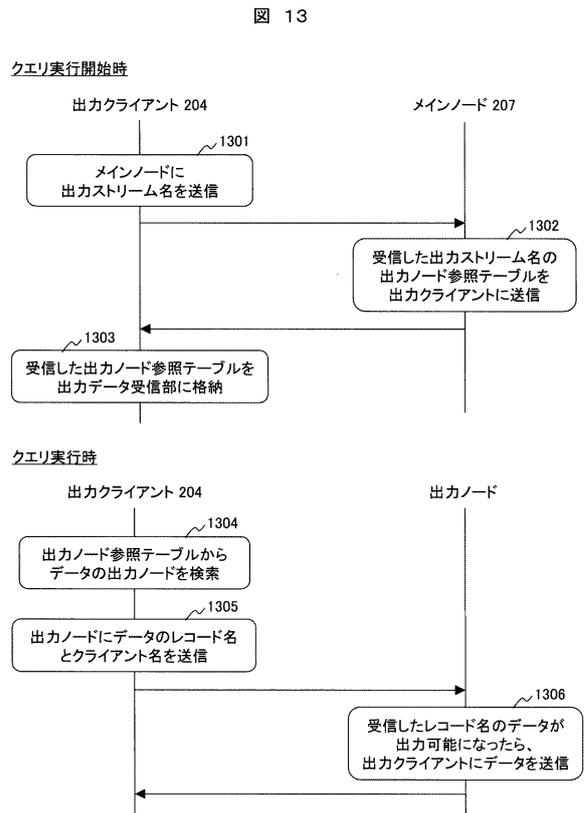
【 図 1 1 】



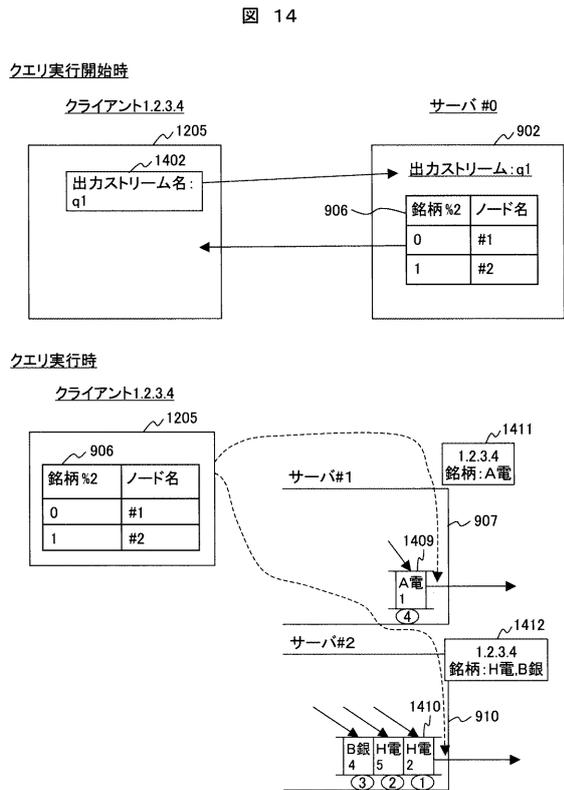
【 図 1 2 】



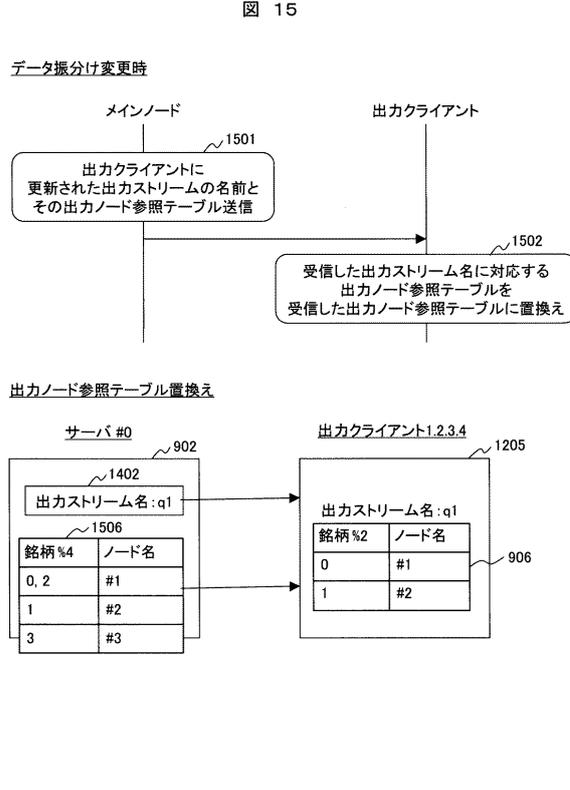
【 図 1 3 】



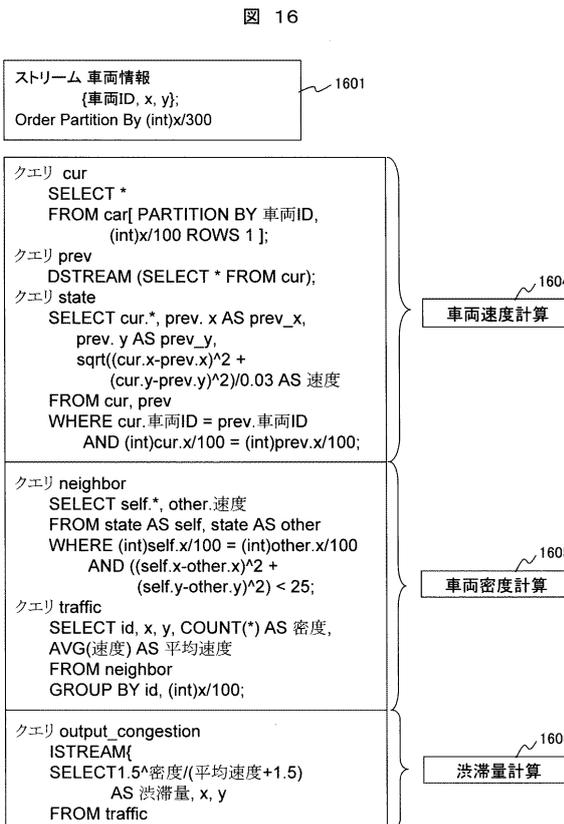
【 図 1 4 】



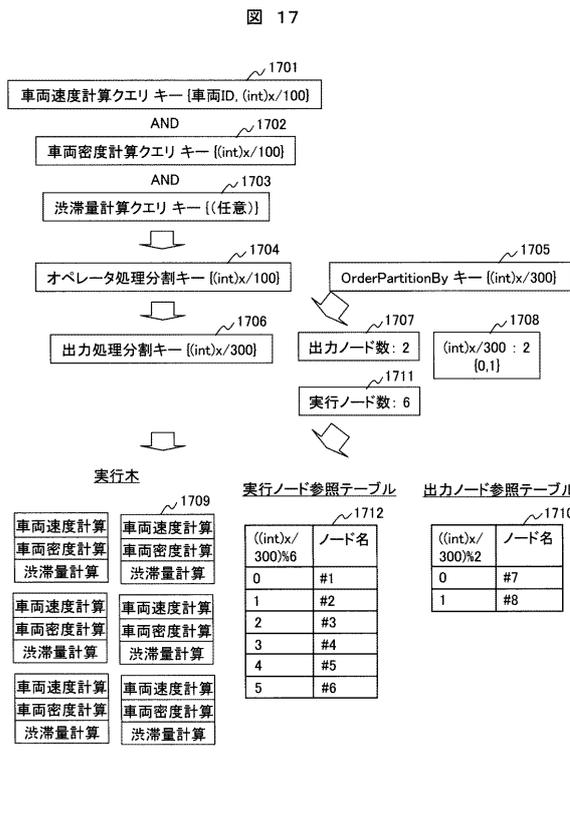
【 図 1 5 】



【 図 1 6 】

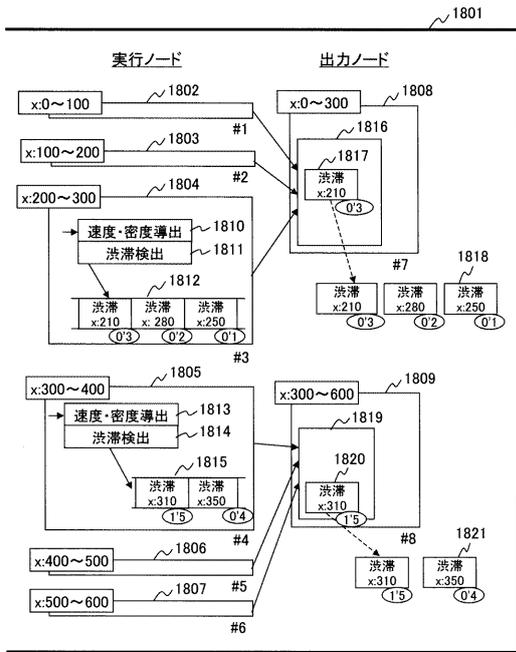


【 図 1 7 】



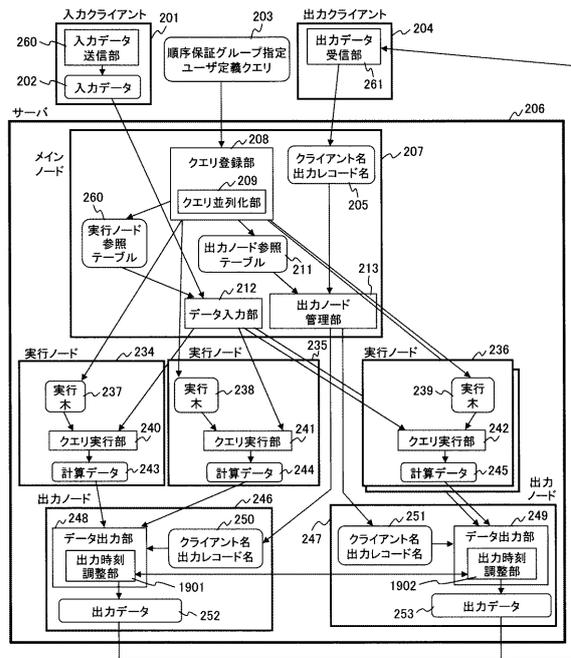
【図 18】

図 18



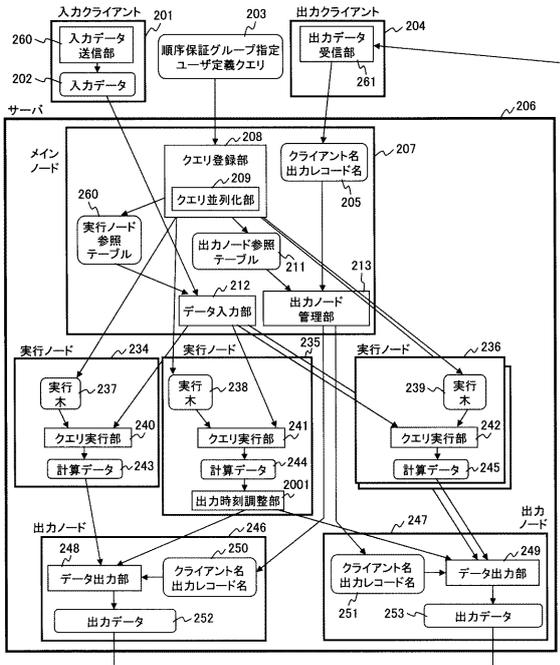
【図 19】

図 19



【図 20】

図 20



【図 21】

図 21

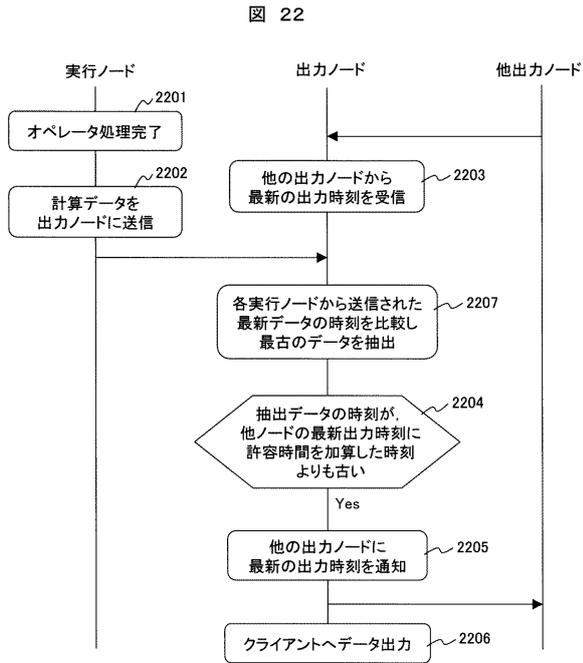
```

2101
ストリーム 株注文
(銘柄,取引所,注文数,口座,業種);
order partition by 銘柄
limit 1 second

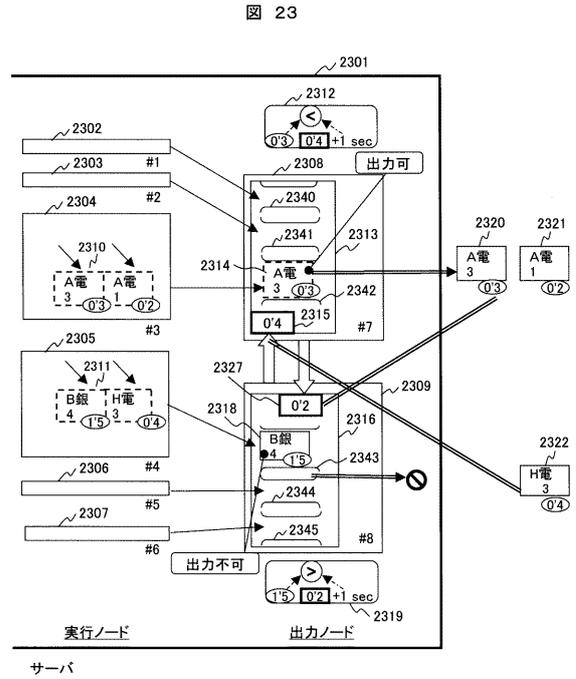
クエリ q1
istream(
SELECT 銘柄, 取引所, sum(注文数)
FROM 株注文 [range 1 minute]
GROUP BY 銘柄, 取引所);

```

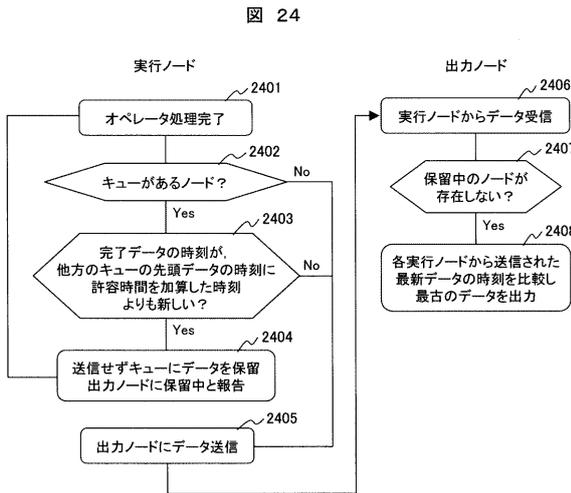
【 図 2 2 】



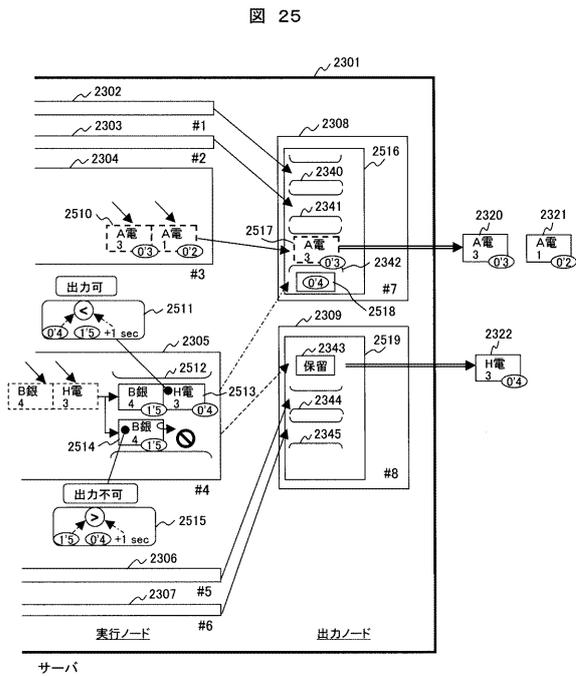
【 図 2 3 】



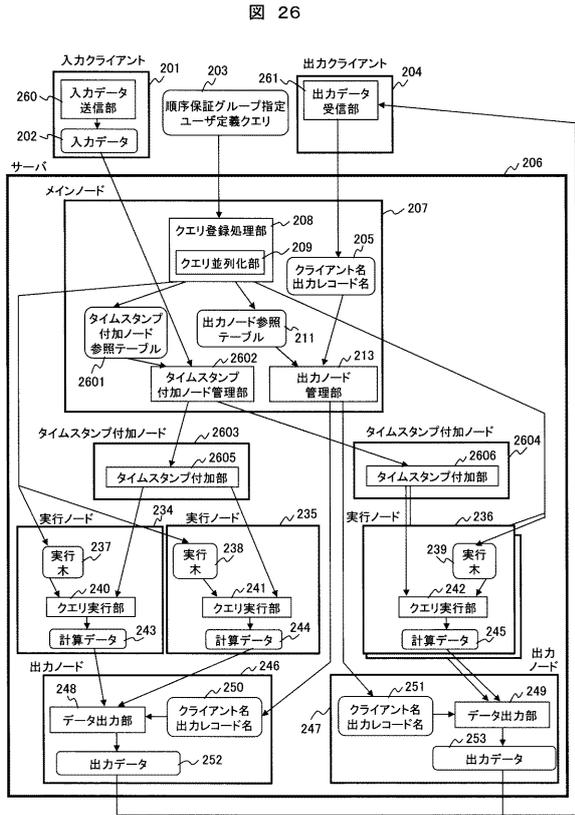
【 図 2 4 】



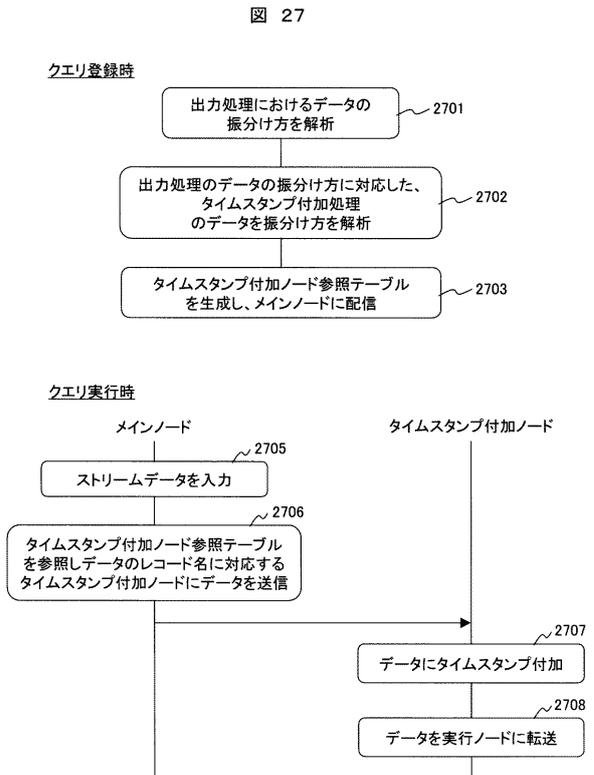
【 図 2 5 】



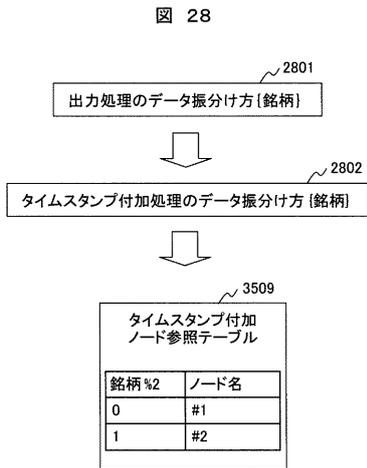
【図 26】



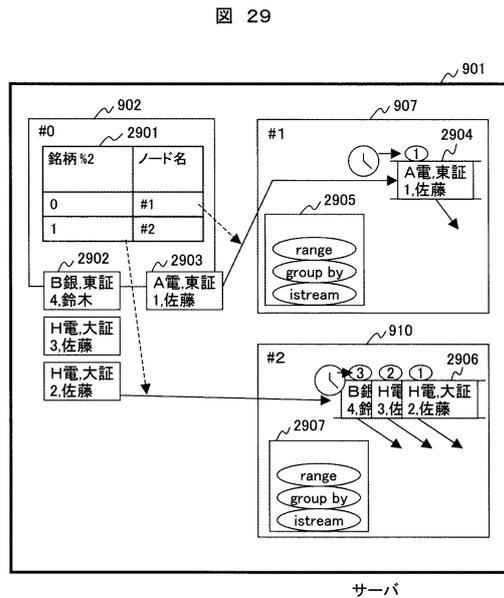
【図 27】



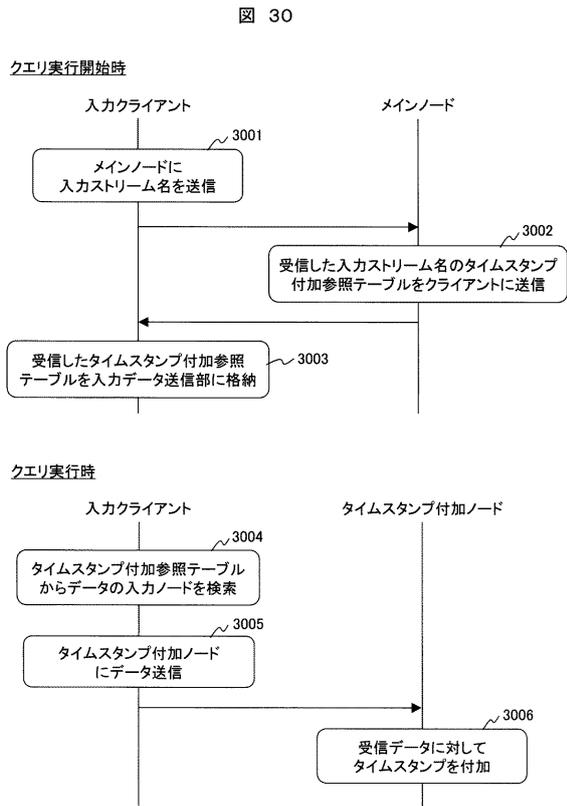
【図 28】



【図 29】



【図30】



【図31】

