

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5675840号  
(P5675840)

(45) 発行日 平成27年2月25日 (2015. 2. 25)

(24) 登録日 平成27年1月9日 (2015. 1. 9)

(51) Int. Cl.	F I
<b>G06F 17/30 (2006.01)</b>	G06F 17/30 110B
<b>G06F 12/00 (2006.01)</b>	G06F 17/30 330B
	G06F 17/30 340Z
	G06F 12/00 513D

請求項の数 60 (全 23 頁)

(21) 出願番号	特願2012-546227 (P2012-546227)	(73) 特許権者	509123208
(86) (22) 出願日	平成22年12月23日 (2010. 12. 23)		アビニシオ テクノロジー エルエルシー
(65) 公表番号	特表2013-516008 (P2013-516008A)		アメリカ合衆国 02421 マサチュー
(43) 公表日	平成25年5月9日 (2013. 5. 9)		セッツ州 レキシントン スプリング ス
(86) 国際出願番号	PCT/US2010/061979		トリート 201
(87) 国際公開番号	W02011/079251	(74) 代理人	100079108
(87) 国際公開日	平成23年6月30日 (2011. 6. 30)		弁理士 稲葉 良幸
審査請求日	平成25年12月20日 (2013. 12. 20)	(74) 代理人	100109346
(31) 優先権主張番号	61/289, 778		弁理士 大貫 敏史
(32) 優先日	平成21年12月23日 (2009. 12. 23)	(72) 発明者	スタンフィル, クレイグ, ダブリュー,
(33) 優先権主張国	米国 (US)		アメリカ合衆国, マサチューセッツ州 O
			1773, リンカーン, ハックルベリー
			ヒル ロード 43

最終頁に続く

(54) 【発明の名称】 クエリー管理

(57) 【特許請求の範囲】

【請求項 1】

1 つ又は複数のデータ・ソースについて実行されるクエリーを管理するための方法であって、

少なくとも第 1 のクエリーを記憶媒体に記憶することと、

処理のために前記第 1 のクエリーを選択することと、

第 1 のクエリー間隔に対して前記 1 つ又は複数のデータ・ソースにおけるデータの第 1 の部分についての前記第 1 のクエリーを処理するようにクエリーエンジンに指示することと、

前記データの第 1 の部分についての前記第 1 のクエリーの処理に基づいて前記クエリーエンジンから第 1 の結果データを受信することと、

前記第 1 のクエリー間隔後に前記第 1 のクエリーの状態を前記記憶媒体に保存することと、

前記第 1 のクエリー間隔後の第 2 のクエリー間隔の間に第 2 のクエリーを処理するように前記クエリーエンジンに指示することと、

前記第 1 のクエリーに関連付けられた優先順位を変更することと、

前記第 2 のクエリーの処理に基づいて前記クエリーエンジンから第 2 の結果データを受信することと、

前記第 2 のクエリーが中断され、前記クエリーエンジンによって処理されないように、前記第 1 のクエリーに関連付けられた前記優先順位に基づいて前記第 2 のクエリーを中断

10

20

することを決定することと、

前記第2のクエリー間隔後の第3のクエリー間隔の間に前記1つ又は複数のデータ・ソース内の第2のデータ部分についての前記第1のクエリーを処理するように前記クエリーエンジンに指示すること

を含む、方法。

【請求項2】

前記第1のクエリーに関係付けられた前記優先順位を前記記憶媒体に記憶することをさらに含み、

前記第1のクエリーに関係付けられた前記優先順位を変更することは、処理のために前記第1のクエリーを選択する前に行われ、

10

処理のために前記第1のクエリーを選択することが部分的に前記優先順位に基づいて前記クエリーを選択することを含む、請求項1記載の方法。

【請求項3】

前記第1のクエリー間隔が所定の時間量によって定義される、請求項1記載の方法。

【請求項4】

前記第1のクエリーに関連付けられた前記優先順位は、前記1つ又は複数のデータ・ソース内の前記データのうちのどのくらいの量が、前記第1のクエリー間隔に対して前記第1のクエリーが実行される前記データの第1の部分に含まれるか、に影響する、請求項3記載の方法。

【請求項5】

20

前記第1のクエリーを記憶することが、前記第1のクエリーを提供したリクエストに通知される前に使用可能になるべき前記第1の結果データの量の通知しきい値を記憶することを含む、請求項1記載の方法。

【請求項6】

前記第1の結果データの前記量が前記通知しきい値を超えた時に前記リクエストに通知することをさらに含み、前記第1のクエリーの前記状態を保存することが前記クエリーエンジンから受信した前記第1の結果データの前記量を記憶することを含む、請求項5記載の方法。

【請求項7】

前記リクエストからの要求による前記第1の結果データを返すことと、前記リクエストに返された前記第1の結果データの前記量を前記記憶媒体に記憶すること、をさらに含む、請求項6記載の方法。

30

【請求項8】

前記第1のクエリーを選択することが、前記クエリーエンジンから受信した前記第1の結果データの前記量と前記リクエストに返された前記第1の結果データの前記量に基づくものである、請求項7記載の方法。

【請求項9】

前記第1のクエリーの前記状態を保存することが、

前記第1のクエリーを中断するよう前記クエリーエンジンに指示することと、

前記第1のクエリーが中断された後に前記第1のクエリーの状態を保存すること、

40

を含む、請求項1記載の方法。

【請求項10】

前記第2のデータ部分に関する前記第1のクエリーを処理するように前記クエリーエンジンに指示することが、

前記第1のクエリーの前記保存状態をロードすることと、

前記第1のクエリーを再開するよう前記クエリーエンジンに指示すること、

を含む、請求項9記載の方法。

【請求項11】

前記第1のクエリーの前記状態を保存することが、二次索引へのオフセットを保存することを含む、請求項9記載の方法。

50

## 【請求項 1 2】

前記二次索引がブロック圧縮索引付きファイルの索引である、請求項 1 1 記載の方法。

## 【請求項 1 3】

前記第 1 のクエリーを複数の副クエリーに分割することと、前記副クエリーのうちの少なくともいくつかを同時に処理するよう前記クエリーエンジンに指示すること、をさらに含む、請求項 1 記載の方法。

## 【請求項 1 4】

前記第 1 のクエリー間隔が始まった後で前記第 2 のクエリーが受信され、前記記憶媒体に記憶される、請求項 1 記載の方法。

## 【請求項 1 5】

前記第 1 のクエリー間隔が始まる前に前記第 2 のクエリーが受信され、前記記憶媒体に記憶される、請求項 1 記載の方法。

## 【請求項 1 6】

1 つ又は複数のデータ・ソースについて実行されるクエリーを管理するためのコンピュータ・プログラムを記憶するコンピュータ可読記憶媒体であって、前記コンピュータ・プログラムが、

少なくとも第 1 のクエリーを記憶媒体に記憶することと、

処理のために前記第 1 のクエリーを選択することと、

第 1 のクエリー間隔に対して前記 1 つ又は複数のデータ・ソース内のデータの第 1 の部分についての前記第 1 のクエリーを処理するようクエリーエンジンに指示することと、

前記データの第 1 の部分についての前記第 1 のクエリーの処理に基づいて前記クエリーエンジンから第 1 の結果データを受信することと、

前記第 1 のクエリー間隔後に前記第 1 のクエリーの状態を前記記憶媒体に保存することと、

前記第 1 のクエリー間隔後の第 2 のクエリー間隔の間に第 2 のクエリーを処理するよう前記クエリーエンジンに指示することと、

前記第 1 のクエリーに関連付けられた優先順位を変更することと、

前記第 2 のクエリーの処理に基づいて前記クエリーエンジンから第 2 の結果データを受信することと、

前記第 2 のクエリーが中断され、前記クエリーエンジンによって処理されないように、前記第 1 のクエリーに関連付けられた前記優先順位に基づいて前記第 2 のクエリーを中断することを決定することと、

前記第 2 のクエリー間隔後の第 3 のクエリー間隔の間に前記 1 つ又は複数のデータ・ソース内の第 2 のデータ部分についての前記第 1 のクエリーを処理するよう前記クエリーエンジンに指示することと、

をコンピュータに実行させるための命令を含む、コンピュータ可読記憶媒体。

## 【請求項 1 7】

1 つ又は複数のデータ・ソースについて実行されるクエリーを管理するためのシステムであって、前記システムが、

少なくとも第 1 のクエリーを記憶する記憶媒体と、

前記 1 つ又は複数のデータ・ソース内のデータに関するクエリーを処理するよう構成されたクエリーエンジンと、

サーバであって、

処理のために前記第 1 のクエリーを選択し、

第 1 のクエリー間隔に対して前記 1 つ又は複数のデータ・ソース内のデータの第 1 の部分についての前記第 1 のクエリーを処理するよう前記クエリーエンジンに指示し、

前記データの第 1 の部分についての前記第 1 のクエリーの処理に基づいて前記クエリーエンジンから第 1 の結果データを受信し、

前記第 1 のクエリー間隔後に前記第 1 のクエリーの状態を前記記憶媒体に保存し、

前記第 1 のクエリー間隔後の第 2 のクエリー間隔の間に第 2 のクエリーを処理するよ

10

20

30

40

50

うに前記クエリーエンジンに指示し、

前記第 1 のクエリーに関連付けられた優先順位を変更し、

前記第 2 のクエリーの処理に基づいて前記クエリーエンジンから第 2 の結果データを受信し、

前記第 2 のクエリーが中断され、前記クエリーエンジンによって処理されないように、前記第 1 のクエリーに関連付けられた前記優先順位に基づいて前記第 2 のクエリーを中断することを決定し、

前記第 2 のクエリー間隔後の第 3 のクエリー間隔の間に前記 1 つ又は複数のデータ・ソース内の第 2 のデータ部分についての前記第 1 のクエリーを処理するように前記クエリーエンジンに指示する、

10

ように構成されたサーバと、  
を含む、システム。

【請求項 18】

前記コンピュータ・プログラムが、

前記第 1 のクエリーに関係付けられた前記優先順位を前記記憶媒体に記憶すること、  
をコンピュータに実行させる命令をさらに含み、

前記第 1 のクエリーに関係付けられた前記優先順位を変更することは、処理のために前記第 1 のクエリーを選択する前に行われ、

処理のために前記第 1 のクエリーを選択することが部分的に前記優先順位に基づいて前記クエリーを選択することを含む、請求項 16 記載のコンピュータ可読記憶媒体。

20

【請求項 19】

前記第 1 のクエリー間隔が所定の時間量によって定義される、請求項 16 記載のコンピュータ可読記憶媒体。

【請求項 20】

前記第 1 のクエリーに関連付けられた前記優先順位は、前記 1 つ又は複数のデータ・ソース内の前記データのうちのどのくらいの量が、前記第 1 のクエリー間隔に対して前記第 1 のクエリーが実行される前記データの第 1 の部分に含まれるか、に影響する、請求項 19 記載のコンピュータ可読記憶媒体。

【請求項 21】

前記第 1 のクエリーを記憶することが、前記第 1 のクエリーを提供したリクエストに通知される前に使用可能になるべき前記第 1 の結果データの量の通知しきい値を記憶することを含む、請求項 16 記載のコンピュータ可読記憶媒体。

30

【請求項 22】

前記コンピュータ・プログラムが、前記第 1 の結果データの前記量が前記通知しきい値を超えた時に前記リクエストに通知することをコンピュータに実行させる命令をさらに含み、前記第 1 のクエリーの前記状態を保存することが前記クエリーエンジンから受信した前記第 1 の結果データの前記量を記憶することを含む、請求項 21 記載のコンピュータ可読記憶媒体。

【請求項 23】

前記コンピュータ・プログラムが、前記リクエストからの要求に対して前記第 1 の結果データを返すことと、前記リクエストに返された前記第 1 の結果データの前記量を前記記憶媒体に記憶することと、をコンピュータに実行させる命令を含む、請求項 22 記載のコンピュータ可読記憶媒体。

40

【請求項 24】

前記第 1 のクエリーを選択することが、前記クエリーエンジンから受信した前記第 1 の結果データの前記量と前記リクエストに返された前記第 1 の結果データの前記量とに基づく、請求項 23 記載のコンピュータ可読記憶媒体。

【請求項 25】

前記第 1 のクエリーの前記状態を保存することが、

前記第 1 のクエリーを中断するよう前記クエリーエンジンに指示することと、

50

前記第 1 のクエリーが中断された後に前記第 1 のクエリーの状態を保存すること、を含む、請求項 1 6 記載のコンピュータ可読記憶媒体。

【請求項 2 6】

前記第 2 のデータ部分に関する前記第 1 のクエリーを処理するように前記クエリーエンジンに指示することが、

前記第 1 のクエリーの前記保存状態をロードすることと、

前記第 1 のクエリーを再開するように前記クエリーエンジンに指示することと、を含む、請求項 2 5 記載のコンピュータ可読記憶媒体。

【請求項 2 7】

前記第 1 のクエリーの前記状態を保存することが、二次索引へのオフセットを保存することを含む、請求項 2 5 記載のコンピュータ可読記憶媒体。

10

【請求項 2 8】

前記二次索引がブロック圧縮索引付きファイルの索引である、請求項 2 7 記載のコンピュータ可読記憶媒体。

【請求項 2 9】

前記コンピュータ・プログラムが、前記第 1 のクエリーを複数の副クエリーに分割することと、前記副クエリーのうちの少なくともいくつかを同時に処理するよう前記クエリーエンジンに指示すること、をコンピュータに実行させる命令を含む、請求項 1 6 記載のコンピュータ可読記憶媒体。

【請求項 3 0】

20

前記第 1 のクエリー間隔が始まった後で前記第 2 のクエリーが受信され、前記記憶媒体に記憶される、請求項 1 6 記載のコンピュータ可読記憶媒体。

【請求項 3 1】

前記第 1 のクエリー間隔が始まる前に前記第 2 のクエリーが受信され、前記記憶媒体に記憶される、請求項 1 6 記載のコンピュータ可読記憶媒体。

【請求項 3 2】

前記サーバは、

前記第 1 のクエリーに関係付けられた前記優先順位を前記記憶媒体に記憶し、

処理のために前記第 1 のクエリーを選択する前に前記第 1 のクエリーに関係付けられた前記優先順位を変更し、

30

処理のために前記第 1 のクエリーを選択することが部分的に前記優先順位に基づいて前記クエリーを選択することを含む、

ようにさらに構成されている、請求項 1 7 記載のシステム。

【請求項 3 3】

前記第 1 のクエリー間隔が所定の時間量によって定義される、請求項 1 7 記載のシステム。

【請求項 3 4】

前記第 1 のクエリーに関連付けられた前記優先順位は、前記 1 つ又は複数のデータ・ソース内の前記データのうちのどのくらいの量が、前記第 1 のクエリー間隔に対して前記第 1 のクエリーが実行される前記データの第 1 の部分に含まれるか、に影響する、請求項 3 3 記載のシステム。

40

【請求項 3 5】

前記第 1 のクエリーを記憶することが、前記第 1 のクエリーを提供したリクエストに通知される前に使用可能になるべき前記第 1 の結果データの量の通知しきい値を記憶することを含む、請求項 1 7 記載のシステム。

【請求項 3 6】

前記サーバは、前記第 1 の結果データの前記量が前記通知しきい値を超えた時に前記リクエストに通知するようにさらに構成され、前記第 1 のクエリーの前記状態を保存することが前記クエリーエンジンから受信した前記第 1 の結果データの前記量を記憶することを含む、請求項 3 5 記載のシステム。

50

## 【請求項 37】

前記サーバは、前記リクエストからの要求に対して前記第1の結果データを返し、前記リクエストに返された前記第1の結果データの前記量を前記記憶媒体に記憶するようにさらに構成されている、請求項36記載のシステム。

## 【請求項 38】

前記第1のクエリーを選択することが、前記クエリーエンジンから受信した前記第1の結果データの前記量と前記リクエストに返された前記第1の結果データの前記量に基づくものである、請求項37記載のシステム。

## 【請求項 39】

前記第1のクエリーの前記状態を保存することが、  
前記第1のクエリーを中断するよう前記クエリーエンジンに指示することと、  
前記第1のクエリーが中断された後に前記第1のクエリーの状態を保存すること、  
を含む、請求項17記載のシステム。

10

## 【請求項 40】

前記第2のデータ部分に関する前記第1のクエリーを処理するように前記クエリーエンジンに指示することが、

前記第1のクエリーの前記保存状態をロードすることと、  
前記第1のクエリーを再開するよう前記クエリーエンジンに指示すること、  
を含む、請求項39記載のシステム。

## 【請求項 41】

前記第1のクエリーの前記状態を保存することが、二次索引へのオフセットを保存することを含む、請求項39記載のシステム。

20

## 【請求項 42】

前記二次索引がブロック圧縮索引付きファイルの索引である、請求項41記載のシステム。

## 【請求項 43】

前記サーバは、前記第1のクエリーを複数の副クエリーに分割し、前記副クエリーのうちの少なくともいくつかを同時に処理するよう前記クエリーエンジンに指示するようにさらに構成されている、請求項17記載のシステム。

## 【請求項 44】

前記第1のクエリー間隔が始まった後で前記第2のクエリーが受信され、前記記憶媒体に記憶される、請求項17記載のシステム。

30

## 【請求項 45】

前記第1のクエリー間隔が始まる前に前記第2のクエリーが受信され、前記記憶媒体に記憶される、請求項17記載のシステム。

## 【請求項 46】

1つ又は複数のデータ・ソースについて実行されるクエリーを管理するためのシステムであって、

少なくとも第1のクエリーを記憶する手段と、  
処理のために前記第1のクエリーを選択する手段と、  
第1のクエリー間隔に対して前記1つ又は複数のデータ・ソース内のデータの第1の部分についての前記第1のクエリーを処理するようクエリーエンジンに指示する手段と、  
前記データの第1の部分についての前記第1のクエリーの処理に基づいて前記クエリーエンジンから第1の結果データを受信する手段と、  
前記第1のクエリー間隔後に前記第1のクエリーの状態を前記記憶媒体に保存する手段と、

40

前記第1のクエリー間隔後の第2のクエリー間隔の間に第2のクエリーを処理するよう前記クエリーエンジンに指示する手段と、

前記第1のクエリーに関連付けられた優先順位を変更する手段と、  
前記第2のクエリーの処理に基づいて前記クエリーエンジンから第2の結果データを受

50

信する手段と、

前記第 2 のクエリーが中断され、前記クエリーエンジンによって処理されないように、前記第 1 のクエリーに関連付けられた前記優先順位に基づいて前記第 2 のクエリーを中断することを決定する手段と、

前記第 2 のクエリー間隔後の第 3 のクエリー間隔の間に前記 1 つ又は複数のデータ・ソース内の第 2 のデータ部分についての前記第 1 のクエリーを処理するように前記クエリーエンジンに指示する手段と

を含む、システム。

【請求項 4 7】

前記第 1 のクエリーに関係付けられた前記優先順位を前記記憶媒体に記憶する手段をさらに含み、

前記第 1 のクエリーに関係付けられた前記優先順位の前記変更は、処理のために前記第 1 のクエリーを選択する前に行われ、

処理のために前記第 1 のクエリーを選択することが部分的に前記優先順位に基づいて前記クエリーを選択することを含む、請求項 4 6 記載のシステム。

【請求項 4 8】

前記第 1 のクエリー間隔が所定の時間量によって定義される、請求項 4 6 記載のシステム。

【請求項 4 9】

前記第 1 のクエリーに関連付けられた前記優先順位は、前記 1 つ又は複数のデータ・ソース内の前記データのうちのどのくらいの量が、前記第 1 のクエリー間隔に対して前記第 1 のクエリーが実行される前記データの第 1 の部分に含まれるか、に影響する、請求項 4 8 記載のシステム。

【請求項 5 0】

前記第 1 のクエリーを記憶することが、前記第 1 のクエリーを提供したリクエストに通知される前に使用可能になるべき前記第 1 の結果データの量の通知しきい値を記憶することを含む、請求項 4 6 記載のシステム。

【請求項 5 1】

前記第 1 の結果データの前記量が前記通知しきい値を超えた時に前記リクエストに通知する手段をさらに含み、

前記第 1 のクエリーの前記状態を保存することが前記クエリーエンジンから受信した前記第 1 の結果データの前記量を記憶することを含む、請求項 5 0 記載のシステム。

【請求項 5 2】

前記リクエストからの要求に対して前記第 1 の結果データを返し、前記リクエストに返された前記第 1 の結果データの前記量を前記記憶媒体に記憶する手段をさらに含む、請求項 5 1 記載のシステム。

【請求項 5 3】

前記第 1 のクエリーを選択することが、前記クエリーエンジンから受信した前記第 1 の結果データの前記量と前記リクエストに返された前記第 1 の結果データの前記量に基づくものである、請求項 5 2 記載のシステム。

【請求項 5 4】

前記第 1 のクエリーの前記状態を保存することが、  
前記第 1 のクエリーを中断するよう前記クエリーエンジンに指示することと、  
前記第 1 のクエリーが中断された後に前記第 1 のクエリーの状態を保存すること、  
を含む、請求項 4 6 記載のシステム。

【請求項 5 5】

前記第 2 のデータ部分に関する前記第 1 のクエリーを処理するように前記クエリーエンジンに指示することが、

前記第 1 のクエリーの前記保存状態をロードすることと、

前記第 1 のクエリーを再開するよう前記クエリーエンジンに指示すること、

10

20

30

40

50

を含む、請求項 5 4 記載のシステム。

【請求項 5 6】

前記第 1 のクエリーの前記状態を保存することが、二次索引へのオフセットを保存することを含む、請求項 5 4 記載のシステム。

【請求項 5 7】

前記二次索引がブロック圧縮索引付きファイルの索引である、請求項 5 6 記載のシステム。

【請求項 5 8】

前記第 1 のクエリーを複数の副クエリーに分割し、前記副クエリーのうちの少なくともいくつかを同時に処理するよう前記クエリーエンジンに指示する手段をさらに含む、請求項 4 6 記載のシステム。

10

【請求項 5 9】

前記第 1 のクエリー間隔が始まった後で前記第 2 のクエリーが受信され、前記記憶媒体に記憶される、請求項 4 6 記載のシステム。

【請求項 6 0】

前記第 1 のクエリー間隔が始まる前に前記第 2 のクエリーが受信され、前記記憶媒体に記憶される、請求項 4 6 記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

20

関連出願の相互参照

本出願は、参照により本明細書に組み込まれる 2009 年 1 月 23 日出願の米国特許出願第 61 / 289778 号に対する優先権を主張するものである。

【0002】

本明細書はクエリー (queries) の管理に関する。

【背景技術】

【0003】

いくつかのデータ記憶システム (例えば、データベース) は、多数のクエリーの処理をサポートするような何らかの方法で、記憶された大量のデータを記憶する。例えば、いくつかのシステムは、並列記憶装置、並列クエリー処理エンジン、又はその両方の使用による並列処理能力を備える。

30

【発明の概要】

【課題を解決するための手段】

【0004】

ある一つの態様では、一般に、1つ又は複数のデータ・ソースについて実行されるクエリーを管理するための方法は、少なくとも第 1 のクエリーを記憶媒体に記憶することと、処理のために第 1 のクエリーを選択することと、第 1 のクエリー間隔に対して 1つ又は複数のデータ・ソース内のデータの第 1 の部分についての第 1 のクエリーを処理するようにクエリーエンジンに指示することと、データの第 1 の部分についての第 1 のクエリーの処理に基づいてクエリーエンジンから結果データを受信することと、第 1 のクエリー間隔後に第 1 のクエリーの状態を記憶媒体に保存することと、第 1 のクエリー間隔後の第 2 のクエリー間隔の間に第 2 のクエリーを処理するようにクエリーエンジンに指示することと、第 2 のクエリー間隔後の第 3 のクエリー間隔の間に 1つ又は複数のデータ・ソース内のデータの第 2 の部分についての第 1 のクエリーを処理するようにクエリーエンジンに指示することを含む。

40

【0005】

諸態様は以下の特徴のうちの 1つ又は複数を含むことができる。

【0006】

この方法は、第 1 のクエリーに関連する優先順位を記憶媒体に記憶することと、処理のために第 1 のクエリーを選択する前に第 1 のクエリーに関連する優先順位を変更すること

50

をさらに含み、処理のために第1のクエリーを選択することが部分的に優先順位に基づいてクエリーを選択することを含む。

【0007】

第1のクエリー間隔は所定の時間量によって定義される。

【0008】

第1のクエリーの優先順位は、1つ又は複数のデータ・ソース内のデータのうちのどのくらいの量が、第1のクエリー間隔に対して第1のクエリーが実行されるデータの第1の部分に含まれるか、に影響する。

【0009】

第1のクエリーを記憶することは、第1のクエリーを提供したリクエストに通知される前に使用可能になるべき結果データの数量の通知しきい値を記憶することを含む。

10

【0010】

この方法は、結果データの数量が通知しきい値を超えた時にリクエストに通知することをさらに含み、第1のクエリーの状態を保存することがクエリーエンジンから受信した結果データの数量を記憶することを含む。

【0011】

この方法は、リクエストからの要求次第で結果データを返すことと、リクエストに返された結果データの数量を記憶媒体に記憶することをさらに含む。

【0012】

クエリーを選択することは、クエリーエンジンから受信した結果データの数量とリクエストに返された結果データの数量に基づくものである。

20

【0013】

第1のクエリーの状態を保存することは、第1のクエリーを中断するようクエリーエンジンに指示することと、第1のクエリーが中断された後に第1のクエリーの状態を保存することを含む。

【0014】

第2のデータ部分に関する第1のクエリーを処理するようクエリーエンジンに指示することは、第1のクエリーの保存状態をロードすることと、第1のクエリーを再開するようクエリーエンジンに指示することを含む。

【0015】

第1のクエリーの状態を保存することは、二次索引 (secondary index) へのオフセットを保存することを含む。

30

【0016】

二次索引はブロック圧縮索引付きファイル (block compressed indexed file) である。

【0017】

この方法は、第1のクエリーを複数の副クエリーに分割することと、その副クエリーのうちの少なくともいくつかを同時に処理するようクエリーエンジンに指示することをさらに含む。

【0018】

第1のクエリー間隔が始まった後で第2のクエリーが受信され、記憶媒体に記憶される。

40

【0019】

第1のクエリー間隔が始まる前に第2のクエリーが受信され、記憶媒体に記憶される。

【0020】

他の態様では、一般に、コンピュータ可読媒体は、1つ又は複数のデータ・ソースについて実行されるクエリーを管理するためのコンピュータ・プログラムを記憶する。このコンピュータ・プログラムは、少なくとも第1のクエリーを記憶媒体に記憶することと、処理のために第1のクエリーを選択することと、第1のクエリー間隔に対して1つ又は複数のデータ・ソース内のデータの第1の部分についての第1のクエリーを処理するようク

50

エリーエンジンに指示することと、データの第1の部分についての第1のクエリーの処理に基づいてクエリーエンジンから結果データを受信することと、第1のクエリー間隔後に第1のクエリーの状態を記憶媒体に保存することと、第1のクエリー間隔後の第2のクエリー間隔の間に第2のクエリーを処理するようにクエリーエンジンに指示することと、第2のクエリー間隔後の第3のクエリー間隔の間に1つ又は複数のデータ・ソース内のデータの第2の部分についての第1のクエリーを処理するようにクエリーエンジンに指示することをコンピュータに実行させるための命令を含む。

【0021】

他の態様では、一般に、1つ又は複数のデータ・ソースについて実行されるクエリーを管理するためのシステムが提供される。このシステムは、少なくとも第1のクエリーを記憶する記憶媒体を含む。このシステムは、1つ又は複数のデータ・ソース内のデータに関するクエリーを処理するように構成されたクエリーエンジンを含む。また、このシステムは、処理のために第1のクエリーを選択し、第1のクエリー間隔に対して1つ又は複数のデータ・ソース内のデータの第1の部分についての第1のクエリーを処理するようにクエリーエンジンに指示し、第1のデータ部分についての第1のクエリーの処理に基づいてクエリーエンジンから結果データを受信し、第1のクエリー間隔後に第1のクエリーの状態を記憶媒体に保存し、第1のクエリー間隔後の第2のクエリー間隔の間に第2のクエリーを処理するようにクエリーエンジンに指示し、第2のクエリー間隔後の第3のクエリー間隔の間に1つ又は複数のデータ・ソース内のデータの第2の部分についての第1のクエリーを処理するようにクエリーエンジンに指示するように構成されたサーバも含む。

【0022】

他の態様では、一般に、1つ又は複数のデータ・ソースについて実行されるクエリーを管理するためのシステムが提供される。このシステムは、少なくとも第1のクエリーを記憶する記憶媒体を含む。このシステムは、1つ又は複数のデータ・ソース内のデータについてのクエリーを処理するように構成されたクエリーエンジンを含む。このシステムは記憶媒体内のクエリーを管理するための手段を含み、その管理は、第1のクエリー間隔に対して1つ又は複数のデータ・ソース内のデータの第1の部分についての第1のクエリーを処理するようにクエリーエンジンに指示することと、データの第1の部分についての第1のクエリーの処理に基づいてクエリーエンジンから結果データを受信することと、第1のクエリー間隔後に第1のクエリーを記憶媒体に保存することと、第1のクエリー間隔後の第2のクエリー間隔の間に第2のクエリーを処理するようにクエリーエンジンに指示することと、第2のクエリー間隔後の第3のクエリー間隔の間に1つ又は複数のデータ・ソース内のデータの第2の部分についての第1のクエリーを処理するようにクエリーエンジンに指示することを含む。

【発明の効果】

【0023】

諸態様は以下の利点のうちの1つ又は複数を含むことができる。

【0024】

部分的にクエリーに関連する優先順位に基づいてクエリーを選択すると、並列クエリー処理システムにおいて効率的な処理が可能になる。クエリーの各部分を部分的に処理し、次に中断することができる複数の間隔に時間をスライスすると、いくつかのクエリーをより速やかに処理することができ、特に優先順位の高いクエリーの場合、システム内の潜在的なバックログが低減される。

【0025】

本発明のその他の特徴及び利点は、以下の説明並びに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】

【0026】

【図1】クエリー処理を描写する概略図である。

【図2】クエリー処理を描写する概略図である。

10

20

30

40

50

- 【図 3】データ記憶システムのブロック図である。  
 【図 4】索引付き圧縮データ記憶装置の概略図である。  
 【図 5 A】クエリーの処理に関連する時間間隔を示す図である。  
 【図 5 B】クエリーの処理に関連する時間間隔を示す図である。  
 【図 6 A】クエリーの処理に関連する時間間隔を示す図である。  
 【図 6 B】クエリーの処理に関連する時間間隔を示す図である。  
 【図 7 A】クエリーの処理に関連する時間間隔を示す図である。  
 【図 7 B】クエリーの処理に関連する時間間隔を示す図である。  
 【図 7 C】クエリーの処理に関連する時間間隔を示す図である。  
 【図 7 D】クエリーの処理に関連する時間間隔を示す図である。  
 【図 8】スライスされたクエリー処理の概略図である。  
 【図 9】索引付き圧縮データ記憶装置のクエリー処理を示す概略図である。  
 【図 10】クエリーを管理するためのプロセスのフローチャートである。  
 【図 11】クエリーを管理するためのプロセスのフローチャートである。  
 【発明を実施するための形態】

10

【0027】

## 1 概要

図 1 を参照すると、いくつかの問題が分散クエリー管理において発生し得る。例えば、先入れ先出し法でデータ記憶システムのクエリーエンジンにクエリーが引き渡されると、システムはバックログになる可能性がある。いくつかのケースでは、引き渡されたクエリーは、あまりリソースを必要とせずに迅速に実行される短いクエリー 102、104、108、112、118 と、実行するのにより長い時間を必要とし、大量のシステム・リソースを使用する長いクエリー 110、114、116 と、短いクエリーと長いクエリーの間どこかに入るクエリーとを含む可能性がある。特定のクエリーが実行される前にそのクエリーが要求するシステム・リソースの量をあらかじめ決めることは実用的ではない可能性がある。図 1 は、複数のクエリーエンジンを使用してクエリーを処理するためのシステムの一例を示している。クエリーは、非同期的に受信されて、待ち行列 101 に記憶され、データ記憶システムのクエリーサーバ 100 上で実行されるクエリーエンジンによって処理される機会を待つ。この例では、初めに、長いクエリー 116 が処理のために第 1 のクエリーエンジン 120 に割り当てられ、短いクエリー 118 が処理のために第 2 のクエリーエンジン 122 に割り当てられる。図 2 を参照すると、短い時間の後に短いクエリー 118 は完了した可能性があり、次に並んでいるクエリーである長いクエリー 114 が空いているクエリーエンジン 122 に割り当てられる。この時点で残りのクエリー 102、104、108、110、112 は、長いクエリー 116、114 のうちの一方が処理を完了し、クエリーエンジン内の処理リソースを解放するまで待つ。この現象は短いクエリーの待ち時間を増やすものであり、迅速応答が期待されているクエリーにおいて受け入れがたい遅延を引き起こす可能性がある。

20

30

【0028】

図 3 を参照すると、データ記憶システム 300 は、クエリーを実行するための要求を受信するために、フロントエンド・サービス 302、例えば、ウェブ・サービスを提供するように構成される。仲介サーバ 304 は複数のクエリーエンジン 312 によるクエリー実行をスケジュールする。各クエリーは割り当てられた期間に対して実行が許可され、その期間は、例えば、時間（例えば、CPU クロックによって測定されたもの）、持続時間、処理された行数、又は検索された行数によって測定され得る。クエリーエンジン 312 は、1 つ又は複数のデータ・ソース 310 A、310 B、310 C からのデータにアクセスし、処理セット 314 を生成するためにクエリーを処理する。データ・ソースを提供する記憶装置は、例えば、仲介サーバ 304 を実現するコンピュータに接続された記憶媒体上に記憶されていて、システム 300 にとってローカルなもの（例えば、ハードディスク・ドライブ）である場合もあれば、リモート接続により通信接続しているリモート・システム上でホストとして処理され、仲介サーバ 304 にとってリモートなもの（例えば、メイ

40

50

ンフレーム)である場合もある。

【0029】

仲介サーバ304は結果セット314を管理する。仲介サーバ304は、クエリーに関する追加情報、例えば、クエリーの優先順位、要求された行数、クエリーから返された行数、リクエストに返された行数、クエリーがどのように使用されるかを示す表示要素(indication)、一度に必要な行数、クエリーがクエリーエンジンによって最後に実行された時間、及び状況標識を記憶することができる。状況標識は、クエリーが待っていること、実行していること、中断されていること、割り込まれていること、又は完了したことを示すことができる。いくつかの編成(arrangements)では、クエリー状態は、クエリー処理中に発生したエラー条件の存在も示すことができる。待ち状態にあるクエリーは、現在実行されていないが、実行する資格のあるものである。実行状態にあるクエリーは、クエリーエンジンによって現在処理されている。中断状態にあるクエリーは、クライアントによって現在要求されている行数を仲介サーバがすでに返したので、実行する資格のないものである。割り込み状態にあるクエリーは、優先順位がより高いクエリーによって先取りされたので、実行する資格のないものである。完了状態にあるクエリーは実行を完了したものである。いくつかの編成では、追加の状況がサポートされる。

10

【0030】

仲介サーバ304は、クエリー結果が処理の準備ができていることをいつ、どのようにリクエストに通知しなければならないかを同定する情報も記憶し得る。いくつかの編成では、複数のクエリーエンジンが単一クエリーの種々の部分について独立して動作し得る。各クエリーエンジンは仲介データベースを独立に更新する。トリガ・イベントが発生した時、例えば、結合されたクエリーエンジンが要求された行数を返した時に、通知イベントがトリガされる。

20

【0031】

いくつかの実装例では、仲介サーバは仲介論理モジュール306と仲介データベース308とを含む。いくつかの実装例では、仲介論理モジュール306は、個々のクエリーエンジン312、フロントエンド・サービス302に組み込まれる場合もあれば、複数のサービスに分割される場合もある。クエリーエンジン312は、結果セット314で使用可能な行数について仲介データベース308を更新し得る。

30

【0032】

いくつかの実装例では、図4を参照すると、データ・ソース310は、索引付き圧縮データ記憶装置402を含む。索引付き圧縮データ記憶装置は、例えば、ファイル内に記憶された複数の圧縮データ・ブロック404を含む。それぞれの圧縮データ・ブロックは、その圧縮データ・ブロック内のデータの位置特定を可能にする、少なくとも1つの索引406に関連付けられている。いくつかの実装例では、第1のキー(例えば、一次キー)に基づいてサーチ可能な一次索引が提供され、その他のキー(例えば、外部キー)に基づいてサーチ可能な1つ又は複数の二次索引が提供される。索引のうちいくつかは、それぞれのキー値が固有のものであるサロゲート・キーで構成することができ、その他の索引は、キーの値がデータ・セット内で固有のものではない可能性がある自然キーに基づくものである。いくつかの実装例では、自然索引を結合して、単一連結索引を作成することができる。索引付き圧縮データ記憶技法及びシステムについては、参照により本明細書に組み込まれる米国特許出願公報第2008/0104149A1号により詳細に記載されている。

40

【0033】

2 クエリースライシング

図5Aを参照すると、一連のクエリーA502、B504、C506、及びD508が、異なるクエリーに関係付けられた時間間隔を示す図に示されている。クエリーが引き渡された順序で実行される場合、クエリーAは間隔502の間に実行されて完了し、次にクエリーBは間隔504の間に実行されて完了し、次に間隔506の間のクエリーCと間隔508の間のクエリーDが続く。これらの条件下では、クエリーAは時間510で完了す

50

るまで結果を返さず、クエリー B は時間 5 1 2 で完了するまで結果を返さず、クエリー C は時間 5 1 4 で完了するまで結果を返さず、クエリー D は時間 5 1 6 で完了するまで結果を返さないであろう。クエリー D は短いクエリーであるが、たまたま他の長いクエリーの後ろに位置していたので結果を返すのに長い時間を要する。

#### 【 0 0 3 4 】

仲介サーバ 3 0 4 のいくつかの実装例では、クエリーを完了まで必然的に順次実行する代わりに、仲介サーバは 1 つのクエリーを複数の異なる小さい部分に分割する。クエリーエンジン 3 0 4 は、特定の間隔の間にクエリーを実行するよう指示される。この間隔は、期間、返すべき行数、処理された行数、又はその他の何らかの基準に基づいて定義することができる。この手法を使用して、図 5 B を参照すると、クエリー A は間隔 5 2 8 の間に実行され、クエリー B は間隔 5 3 0 の間に実行され、クエリー C は間隔 5 3 2 の間に実行され、クエリー D は間隔 5 3 4 の間に実行され（完了する）、次にクエリー A は第 2 の間隔の間にもう一度実行される。いくつかのケースでは、1 つのクエリーが処理されるそれぞれの間隔後にそのクエリーをサブミットしたプロセスに対し、そのクエリーによるいくつかの結果が返される可能性がある。例えば、クエリー A によるいくつかの結果は時間 5 2 0 後に返され、クエリー B、C、及び D によるいくつかの結果は、それぞれ時間 5 2 2、5 2 4、5 2 6 後に返される可能性がある。これらのクエリーを小さい実行間隔に分割することにより、システム 3 0 0 は、他のクエリーが実行される前にクエリーが完了するのを待たなければならない場合より速やかにより多くのクエリーに関するいくつかの結果を生成することができる。さらに、他のクエリーを遅延させるというトレードオフにより、いくつかのクエリーは、本来完了したと思われる時期より速やかに完了することができる。この例では、クエリー D は時間 5 2 6 で完了し、クエリー C は時間 5 4 0 で完了し、クエリー A は時間 5 4 2 で完了し、クエリー B は時間 5 4 4 で完了する。従って、この例では、長いクエリー A 及び B を遅延させるという犠牲を払って、短いクエリー C 及び D の方がより速やかに完了する。

#### 【 0 0 3 5 】

クエリーの分割方法の決定は、システムに望ましい動作特性に依存しうる。例えば、時間に基づいてクエリーを分割することは、それぞれのクエリーが特定の量の作業を実行できることを保証する可能性があるが、その作業がどのくらいの長さの暦時間を費やせるかという保証はなく、1 つの実行間隔でどのくらいの行数が返されるかについても保証はない。対照的に、いくつかの行が返されるまでクエリーを実行できるようにすることにより、いくつかの結果を生成するためにどのくらいの数の実行間隔が必要になるかが決定されるが、1 つの間隔がどのくらい長く持続するかについての保証はない。いくつかの行が処理されるまでクエリーを実行できるようにすることは、システムに、クエリーを完了するのにどのくらいの数の実行間隔が必要になるかを同定することを可能にするが、特定の数の行を返すのにどのくらいのサイクル数が必要であるか又は具体的に特定の実行サイクルがどのくらいの長さの時間を要するかを知らせることはない。

#### 【 0 0 3 6 】

クエリーを処理するための時間は、単一クエリーのみが処理されている場合でも複数の実行間隔（又は「クエリー間隔」）に分割することができる。あるクエリー間隔の終わりに、新しいクエリーが到着している場合、処理されているクエリーは中断され、次のクエリー間隔を使用して新しいクエリーを処理する。代わって、そのクエリー間隔の終わりに新しいクエリーが到着していない場合、処理されているクエリーは追加のクエリー時間の間に処理を続行することができる。例えば、図 6 A の例では、クエリー B はクエリー A の処理中の時間 6 1 0 に到着し、図 6 B の例では、クエリー A 又はクエリー B のいずれかの処理が始まる前に両方のクエリー A 及び B が到着する。

#### 【 0 0 3 7 】

図 6 A の例では、クエリー A は間隔 6 0 2 の間に実行され、クエリー A が間隔 6 0 2 の終わりに完了していない場合、システムは、追加のクエリー間隔の間にクエリー A を処理しなければならないかどうか又は他のクエリーが処理を待っているかどうかを判断するた

10

20

30

40

50

めにチェックする。クエリー B は間隔 6 0 2 の終わりにまだ到着していないので、クエリー A はクエリー間隔 6 0 4 の間に処理される。同様に、クエリー A は次のクエリー間隔 6 0 6 の間にも処理される。しかし、クエリー間隔 6 0 6 の終わりに、システムは、時間 6 1 0 に到着したクエリー B を間隔 6 0 8 の間に処理しなければならないと判断する。次に、それぞれが完了するまで（この例では、クエリー A は時間 6 1 2 に完了し、クエリー B は時間 6 1 4 に完了する）交互の間隔においてクエリー A 及び B が処理される。図 6 B の例では、クエリー A は間隔 6 2 0 の間に実行され、クエリー A が間隔 6 2 0 の終わりに完了していない場合、システムは、追加のクエリー間隔の間にクエリー A を処理しなければならないかどうか又は他のクエリーが処理を待っているかどうかを判断するためにチェックする。クエリー B は間隔 6 2 0 の終わり以前にすでに到着しているので、クエリー B はクエリー間隔 6 2 2 の間に処理される。次に、それぞれが完了するまで交互の間隔においてクエリー A 及び B が処理される。

10

**【 0 0 3 8 】**

クエリー間隔の終わりにクエリーを中断することは、クエリーの状態を仲介データベースに保存することを含む。ある編成では、1つの間隔後に仲介データベース内でクエリー状態を「中断」又は他の状態に更新し、そのクエリーが実行する資格のないものであることを示すことができる。所定の間隔後に、クエリーの状況を「待ち」に更新し、そのクエリーをもう一度実行できるようにすることができる。他の編成では、仲介サーバは自動的に所定の間隔後に直ちにクエリーをスケジュールする。

**【 0 0 3 9 】**

20

**3 クエリーの優先順位付け及び再優先順位付け**

仲介データベースは個々のクエリーに関係付けられた優先順位を記憶し得る。この優先順位は、クエリーが実行される頻度及び方法に影響し得る。図 7 A を参照すると、優先順位の高いクエリー A には、クエリー B（間隔 7 0 4 の間に処理される）又は優先順位が低いクエリー C（間隔 7 0 6 の間に処理される）より大きい実行間隔 7 0 2 が提供され得る。この例では、優先順位の高いクエリー A にはクエリー B に提供される実行間隔 7 0 4 より大きい実行間隔 7 0 2 が提供され、クエリー B には優先順位の低いクエリー C に提供される実行間隔 7 0 6 より大きい実行間隔 7 0 4 が提供される。或いは、図 7 B を参照すると、優先順位の高いクエリー A には標準的な優先順位のクエリー B（間隔 7 1 0 の間に処理される）より高い頻度の実行間隔 7 0 8 が提供され、標準的な優先順位のクエリー B には優先順位の低いクエリー C（間隔 7 1 2 の間に処理される）より高い頻度の実行間隔が提供され得る。図 7 C を参照すると、ある状況では、クエリー A には、クエリー A が（間隔 7 1 4 後に）実行を完了するまで他のクエリー B 及び C の処理が中断されるように十分な高い優先順位が提供され、実行を完了した時点で、それぞれ間隔 7 1 6 と 7 1 8 の間で交互に、中断されたクエリー B 及び C の実行が再開される。

30

**【 0 0 4 0 】**

また、仲介データベースは、クエリーが実行している間、クエリーに再優先順位付けを可能にする。例えば、図 7 D を参照すると、優先順位の高いクエリー A は（間隔 7 2 0 の間）、通常の優先順位のクエリー B（間隔 7 2 2 の間）及び優先順位の低いクエリー C（間隔 7 2 4 の間）とともに、仲介データベースによってスケジュールされる。時間 7 2 6 で、優先順位の高いクエリー A は通常の優先順位レベルに再優先順位付けされる。その時点で、仲介データベースは、新しい優先順位付けに基づいてクエリーのスケジューリングを調整する。再優先順位付け後に前進すると、次に、通常の優先順位のクエリー A には通常の優先順位のクエリー B に提供される間隔 7 2 2 と同様のサイズの実行間隔 7 2 8 が提供される。

40

**【 0 0 4 1 】**

再優先順位付けは、要求しているプロセスによってなされた判断により生じる場合もあれば、それ自体の基準に基づいて仲介サーバ内で生じる場合もある。例えば、仲介サーバにはクエリーを完了するための期限が設けられる可能性があり、期限が近づくにつれて、サーバは適時完了を保証するためにクエリーの優先順位を高め得る。いくつかのケースで

50

は、仲介サーバが優先順位のより高いトラフィックをチェックできるようにするために、単一のより大きい実行間隔の代わりに、複数のより小さい実行間隔がクエリーに提供され得る。その他のケースでは、仲介サーバは、優先順位のより高いクエリーが実行できるようにするために、実行中のクエリーの実行間隔に割り込むことができる。

#### 【 0 0 4 2 】

いくつかのケースでは、クエリーは、前のクエリーの実行前又は実行中のいずれかの実行に先立って次の間隔に入る次のクエリーとともに、スケジューリングされ得る。いくつかのケースでは、実行のためにスケジューリングされるべき次のクエリーは、選択基準に基づいて実行直前に選択され得る。

#### 【 0 0 4 3 】

### 4 並列クエリー処理

多くのシステムにとって、複数のクエリーを一度に実行することが有利であり得る。例えば、単一システム上で実行されている2つのクエリーは、1つのシステム上で実行される単一クエリーよりも改善されたパフォーマンスを実現 (experience) する。これは、例えば、第2のクエリーが異なるリソースを使用している間に一方のクエリーが1つのコンピューティング・リソースを使用できるように生じ得る。一度に両方のクエリーを実行することによって、スループットは改善される。いくつかの実装例では、図8を参照すると、優先順位の高いクエリー802は、複数のクエリースライス804、806、808、810、812に分割される。各スライスは、個別のクエリーエンジン814、816、818、820、822によって処理され得る。

#### 【 0 0 4 4 】

優先順位の高いクエリー802は、上記のように処理すべき行数に基づいてスライスされ得る。クエリーを完了するためにどのくらいの数の実行間隔が必要になるかを判断するために、パーティション化情報は、クエリーのターゲットである索引付き圧縮データ記憶装置の二次索引と比較され得る。これは、それぞれのクエリースライスによって索引付き圧縮データ記憶装置のどの部分が処理されるかを同定することにもなる。例えば、図9を参照すると、索引付き圧縮ファイル902は複数のデータ・ブロック904、906、908、910を含み、各データ・ブロックは複数のデータ・レコードを含む。索引付き圧縮ファイル902は、データ・ブロックを参照する索引912に関係付けられる。いくつかの編成では、この索引は、それぞれのデータ・ブロックに関する1つの索引レコード922を含む可能性があり、他の編成では、索引912はデータ・ブロックより少ない索引レコード922を含む得る。いくつかの編成では、各索引レコード922はデータ・ブロック904、906、908、910を参照し、他の編成では、各索引レコード922はデータ・ブロックの第1のデータ・レコードを参照する。仲介サーバは、索引912を検討し、索引レコードに基づいてクエリー実行間隔 (又は「クエリースライス」) を決定する。この例では、クエリーエンジンは、索引912に基づいて4つのクエリースライス914、916、918、920を作成することを選択する。1つのクエリースライス914は、ブロック1:904から始まるデータ・レコードを処理し、ブロック10 (図示せず) の終わりで終了し、クエリースライス916は、ブロック11:906から始まるデータを処理し、ブロック20 (図示せず) の終わりで終了し、クエリースライス918は、ブロック21:908から処理を開始し、ブロック30 (図示せず) の終わりで終了し、最後にクエリースライス920は、ブロック31:910から処理を開始し、索引付き圧縮ファイル902の終わりで処理を終了する。この例では、仲介サーバは、索引912内の索引レコード922の数によってのみ制限される、任意の数のクエリースライスを作成することを選ぶことができる。

#### 【 0 0 4 5 】

図8を参照すると、クエリーの各スライスは、クエリーエンジン814、816、818、820、822のそれぞれ異なる1つによって同時に処理され得る。例えば、クエリースライス804はクエリーエンジン814によって処理され、クエリースライス806は実質的に同時にクエリーエンジン816によって処理される。同時に、クエリースライ

10

20

30

40

50

ス 8 0 8 はクエリーエンジン 8 1 8 によって処理され、クエリースライス 8 1 0 はクエリーエンジン 8 2 0 によって処理され、クエリー 8 1 2 はクエリーエンジン 8 2 2 によって処理される。各クエリーエンジンは、そのクエリーパーティションに関する結果セットを生成する。すべての結果セットが生成されると、結果セットは、クエリー全体に対する完全な結果セットを形成するように、結合されうる。この方法を使用すると、優先順位の高いクエリーは、通常、その動作を完了するのに要する時間の一部で完了することができる。

#### 【 0 0 4 6 】

##### 5 コールバック

あらかじめ指定された基準によって定義されたトリガが満たされると、システムは通知を行う。図 3 を参照すると、新しいクエリーがフロントエンド・サービス 3 0 2 にサブミットされる場合、このサブミットは、条件が満たされた時に（フロントエンド・サービス 3 0 2 を介して）リクエストに通知するよう仲介サーバ 3 0 4 に要求する情報を含み得る。ある編成では、この条件は、特定の数の結果データ要素がリクエストによってアクセスできる状態になった時の通知であり得る。例えば、リクエストは、1 0 0 個の結果レコードの準備ができた時に通知され得る。いくつかのケースでは、リクエストは、通知の前に準備ができなければならない結果データ要素の数を指定し得る。他のケースでは、リクエストは、リクエストが通知を受ける前に満たさなければならない他の基準を提供し得る。例えば、リクエストは、クエリーが中断された時又はすべての処理が完了した時に通知を受けたいと希望する可能性がある。いくつかのケースでは、トリガ基準は、仲介データベース 3 0 8 内で追跡された状態情報に制限されるかも知れない。他のケースでは、トリガ基準は制限がないかも知れない。トリガは、いくつかの異なる方法で仲介サーバ 3 0 4 に提供され得る。例えば、トリガは、それぞれのクエリー間隔後に仲介サーバ 3 0 4 が実行するスクリプトとして、又は所定のアプリケーション・プログラミング・インターフェース（API）に適合するコンパイル済みクラスとして、提供され得る。いくつかのケースでは、例えば、1 0 0 個の結果レコードが発見されたという条件など、その条件が一度しか発生しない可能性がある。他の編成では、例えば、1 0 0 個の追加の結果レコードが発見されるごとに通知を求める要求など、その条件が再発する可能性もある。

#### 【 0 0 4 7 】

いくつかのケースでは、トリガ条件のサブミットはアクション定義も含み得る。このアクションは、トリガとともに仲介データベース 3 0 8 に記憶され得る。アクションは、条件が満たされた時に仲介サーバ 3 0 4 がどのように応答するかを定義するものである。アクションは、例えば、通知、要約など、所定の 1 組の可能なアクションのうちの 1 つにすることができる。アクションは、仲介サーバ 3 0 4 上で実行されるスクリプトにすることができる。例えば、1 つのアクションは、返された結果をクエリーパラメータとして使用して、追加のクエリーをシステムにサブミットすることができる。また、アクションは、事前確立された API に適合するコンパイル済みクラスとして提供することができる。

#### 【 0 0 4 8 】

##### 6 クエリーの中断

いくつかの実装例では、仲介サーバ 3 0 4 はクエリーの処理を中断することができる。仲介サーバは、そのクエリーに中断というマークを付けることができ、そのクエリーが再開されるまでいかなる処理も行われぬ。クエリーの中断とともに、クエリーサーバはクエリーの状態を保存し得る。この状態はクエリーのプロセスの表示要素である。例えば、この状態は索引付き圧縮データ・ストアへのオフセットになる場合もあれば、B 木内で最後に処理されたノードを含む場合もある。

#### 【 0 0 4 9 】

いくつかのケースでは、仲介サーバはそれ自体のイニシアチブについてのクエリーを中断することを選択しうる。これは、例えば、1 つのクエリーがいくつかのレコードを結果セット内に生成し、リクエストに引き渡されるのを待っている結果セット内の行数がしきい値を超えた時に生じ得る。

10

20

30

40

50

## 【 0 0 5 0 】

例えば、クエリーをサブミットするリクエストは、固定数の行の引き渡しを後で要求し得る（例えば、ユーザ・インターフェースが画面を埋め尽くす（populate）ために25行分のデータの「ページ」を要求するならば、システムはクエリーから25行分のデータを要求する。）。その後、ユーザがより多くのクエリー結果を見ることを希望していることを示した場合、システムはクエリー結果の次の「ページ」又は26～50の結果を要求することができる。仲介データベースは、クエリーから返された結果の数と、ユーザに返された結果の数を追跡する。例えば、クエリーは300行を返した可能性があるが、25行がリクエストに送信された可能性がある。クエリーから返された行数がマージン（例えば、25、50、75、又は100行）によってリクエストに送信された行数を超える場合、仲介データベースはそのクエリーの処理を中断し得る。これは、仲介データベース内でそのクエリーに中断というマークを付けることによるか、又はそのクエリーの次の実行をスケジュールする前のチェックを介して、成し遂げられる。

10

## 【 0 0 5 1 】

いくつかのケースでは、しきい値はシステム300によって定義され、他のケースでは、そのクエリーがどのように使用されるかに依存して各クエリーについて個別にしきい値が定義され得る。例えば、固定数の項目を有するWebページ上にデータ・リストを表示するためにその結果が使用されるクエリーは、4ページ分のデータが待っている時に、中断し得る。対照的に、そのクエリーによって返されるすべてのデータの要約レポート、例えば、月末レポートを作成するためにその結果が使用されるクエリーは、けっして中断し得ない。いくつかのケースでは、しきい値は、リクエストに通知する前に収集するための行数から推測され得る。

20

## 【 0 0 5 2 】

いくつかのケースでは、クエリーは、仲介データベース内のそのクエリーの状態情報を更新することによって、明示的に中断され得る。例えば、優先順位のより高いクエリーを実行できるようにするために、クエリーに中断というマークが付され得る。他のケースでは、仲介サーバのスケジューリング・アルゴリズムは、中断されたクエリーの状態を有するクエリーがスケジュールされないようになっているので、クエリーは、暗黙的に中断され得る。クエリーの中断は、クエリーがサブミットされ、その後、そのクエリーが完了する前に呼び出しプログラムが終了する時にリソースの浪費を最小限にするという追加の利点を有する。仲介サーバは、リクエストが定義済み期間の間に結果にアクセスしなかった場合にクエリー及び結果セットを削除することを選ぶことができる。

30

## 【 0 0 5 3 】

## 7 仲介サーバの処理

図10を参照すると、フローチャート1000は、外部要求なしにクエリーの処理を中断すべきかどうかに関する判断を含む、仲介サーバ304の動作の模範的な編成を表している。

## 【 0 0 5 4 】

動作は、実行のためにクエリーを選択すること1002を含む。一例では、クエリーは、仲介サーバによって確立された所定のスケジュールの一部として選択され得る。他の例では、クエリーは、クエリーの優先順位及びクエリーが最後に実行された時間を含み得る何らかの基準に基づいて選択され得る。ある編成では、仲介サーバは実行を待っている（例えば、待ち状態にある）クエリーについて繰り返す。クエリーのそれぞれはクエリーエンジン上で実行するようにスケジュールされる。待っているすべてのクエリーが実行されると、仲介サーバは、依然として実行を待っているクエリーについてプロセスを繰り返す。他の編成では、仲介サーバは、実行のために最も長い間隔の間、待っていたクエリーを選択する。他の編成では、仲介サーバは、実行のために優先順位が最も高いクエリーを選択する。

40

## 【 0 0 5 5 】

動作はまた、クエリーエンジン上でクエリーを実行すること1004を含む。一例では

50

、選択されたクエリーはクエリーエンジンに割り当てることができ、そのクエリーエンジンはデータ・レコードに対してクエリーを実行し、結果セットを更新し、返された行数を仲介サーバに通知する。

【 0 0 5 6 】

動作はまた、引き渡されるのを待っている行数のチェック 1 0 0 6 を含む。引き渡されるのを待っている行数が通知しきい値を超える場合、仲介サーバはリクエストへのコールバック 1 0 0 8 を実行する。

【 0 0 5 7 】

動作はまた、リクエストがアクセスするのを待っている行数が中断しきい値を超えるかどうかのチェック 1 0 1 0 を含み、その場合、1 0 1 2 でクエリーが中断される。クエリーが中断されるかどうかにかかわらず、仲介サーバは次に処理するクエリーの選択に移行する。

10

【 0 0 5 8 】

図 1 1 を参照すると、フローチャート 1 1 0 0 は、例えば、結果がアクセスできる状態になっていることをコールバックがリクエストに通知した後に、クエリーによって返された結果セットの一部にリクエストがアクセスしたことに応じて、仲介サーバ 3 0 4 の動作の模範的な編成を表している。

【 0 0 5 9 】

動作は、リクエストがクエリーからの結果を要求すること 1 1 0 2 を含む。いくつかの編成では、リクエストは、返すべき行数の指示、例えば、2 5 行を返すよう求める要求を送信し得る。他の編成では、リクエストは特定の範囲の結果を返すよう要求し得る。例えば、リクエストは、5 0 ~ 1 2 6 の結果を返すよう要求し得る。さらに他の編成では、リクエストは収集したすべての結果を返すことを要求し得る。

20

【 0 0 6 0 】

動作はまた、結果を返し、レコードを更新すること 1 1 0 4 を含む。要求に応答して、仲介サーバは要求された行へのアクセスを提供し得る。いくつかの編成では、仲介サーバはまた、そのクエリーが依然として追加の結果を処理しているという表示要素をリクエストに送信し得る。他の編成では、仲介サーバは、追加の結果が即時引き渡しに使用可能であるという表示要素も提供し得る。

【 0 0 6 1 】

動作はまた、そのクエリーが現在中断されているかどうかを判断するためのチェック 1 1 0 6 を含む。クエリーが中断されている場合、制御は次の動作 1 1 0 8 に移行する。そうではない場合、プロセスは完了する。

30

【 0 0 6 2 】

動作はまた、引き渡しを待っている行数が中断しきい値未満であるかどうかを判断するためのチェック 1 1 0 8 を含む。そうである場合、クエリーは 1 1 1 0 で再開され、仲介サーバによる処理のためにスケジュールされ得る。

【 0 0 6 3 】

上記のクエリー管理手法は、コンピュータ上で実行するためにソフトウェアを使用して実現することができる。例えば、このソフトウェアは、1 つ又は複数のプログラム式又はプログラム可能コンピュータ・システム（分散、クライアント/サーバ、又はグリッドなどの様々なアーキテクチャのものにすることができる）上で実行される 1 つ又は複数のコンピュータ・プログラム内の手順を形成し、それぞれのコンピュータ・システムは少なくとも 1 つのプロセッサと、少なくとも 1 つのデータ記憶システム（揮発性及び不揮発性メモリ及び/又は記憶素子を含む）、少なくとも 1 つの入力装置又はポート、並びに少なくとも 1 つの出力装置又はポートを含む。このソフトウェアは、例えば、計算グラフの設計及び構成に関連するその他のサービスを提供する、より大きいプログラムの 1 つ又は複数のモジュールを形成することができる。グラフのノード及び要素は、コンピュータ可読媒体に記憶されたデータ構造又はデータ・リポジトリに記憶されたデータ・モデルに適合するその他の組織化されたデータとして実現することができる。

40

50

## 【 0 0 6 4 】

このソフトウェアは、汎用又は特殊目的プログラム可能コンピュータによって読み取り可能なCD-ROMなどの記憶媒体上で提供するか、或いはそれが実行されるコンピュータへのネットワークの通信媒体により配布する（伝搬信号にコード化する）ことができる。すべての機能は、特殊目的コンピュータ上で又はコプロセッサなどの特殊目的ハードウェアを使用して実行することができる。このソフトウェアは、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タによって実行されるという分散方法で実現することができる。それぞれのこのようなコンピュータ・プログラムは、好ましくは、汎用又は特殊目的プログラム可能コンピュータによって読み取り可能なストレージ・メディア又はデバイス（例えば、ソリッドステート・メモリ又はメディア或いは磁気又は光メディア）上に記憶されるか又はそれにダウンロードされ、本明細書に記載された手順を実行するためにそのストレージ・メディア又はデバイスがコンピュータ・システムによって読み取られた時にそのコンピュータを構成し操作する。また、本発明のシステムは、コンピュータ・プログラムとともに構成されたコンピュータ可読記憶媒体として実現されるものと見なすことができ、このように構成された記憶媒体は本明細書に記載された機能を実行するためにコンピュータ・システムを具体的かつ定義済みの方法で動作させる。

10

## 【 0 0 6 5 】

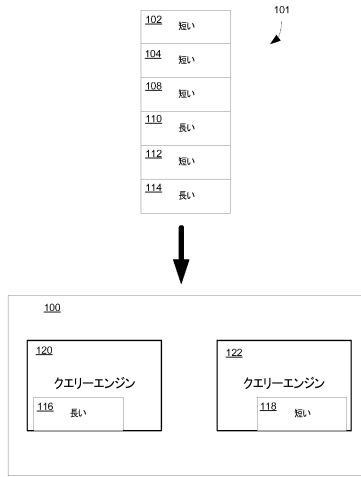
本発明のいくつかの実施形態について説明してきた。それにもかかわらず、本発明の精神及び範囲を逸脱せずに様々な変更が可能であることが理解されるであろう。例えば、上記の諸ステップのうちのいくつかは順序とは無関係なものにすることができ、従って、上記のものとは異なる順序で実行することができる。

20

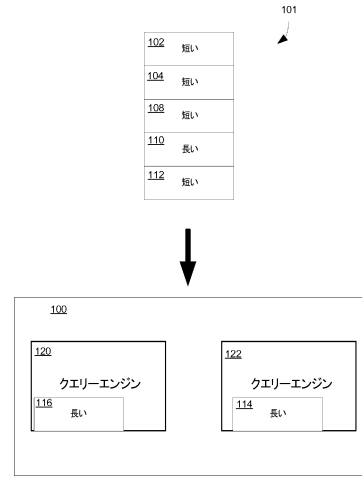
## 【 0 0 6 6 】

上記の説明は本発明を例示するためのものであって、本発明の範囲を限定するためのものではなく、本発明は特許請求の範囲の範囲によって定義されることを理解されたい。例えば、上記の機能ステップのうちのいくつかは、処理全体に実質的に影響せずに異なる順序で実行することができる。その他の諸実施形態は特許請求の範囲の範囲内である。

【図 1】



【図 2】



【図 3】

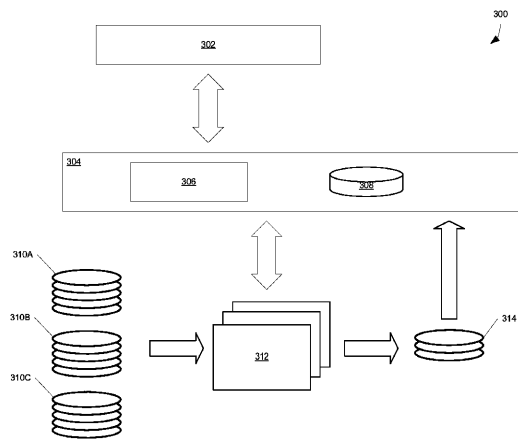


FIG. 3

【図 4】

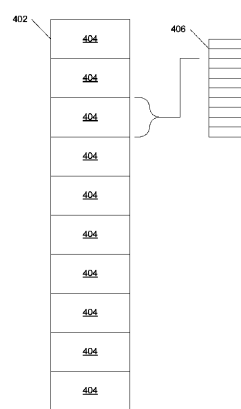


FIG. 4

【図 5 A】

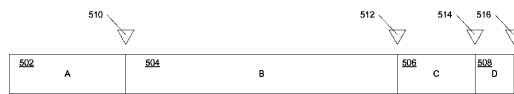


FIG. 5A

【 図 5 B 】

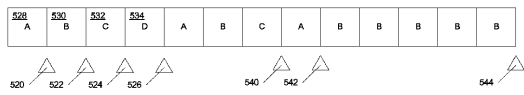


FIG. 5B

【 図 6 A 】

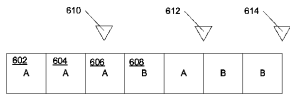


FIG. 6A

【 図 6 B 】



FIG. 6B

【 図 7 A 】

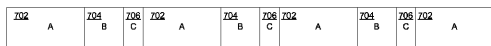


FIG. 7A

【 図 7 B 】

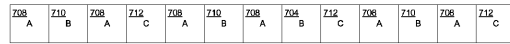


FIG. 7B

【 図 7 C 】



FIG. 7C

【 図 7 D 】

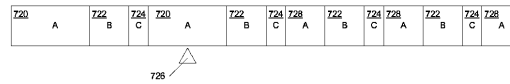


FIG. 7D

【 図 8 】

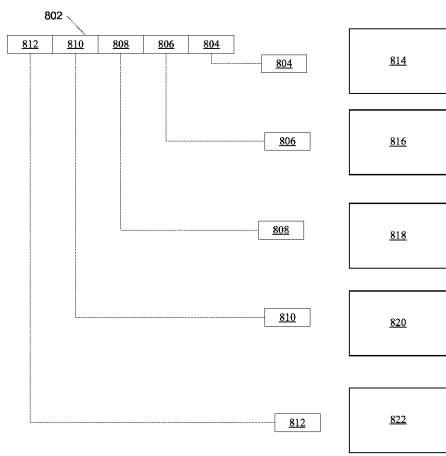
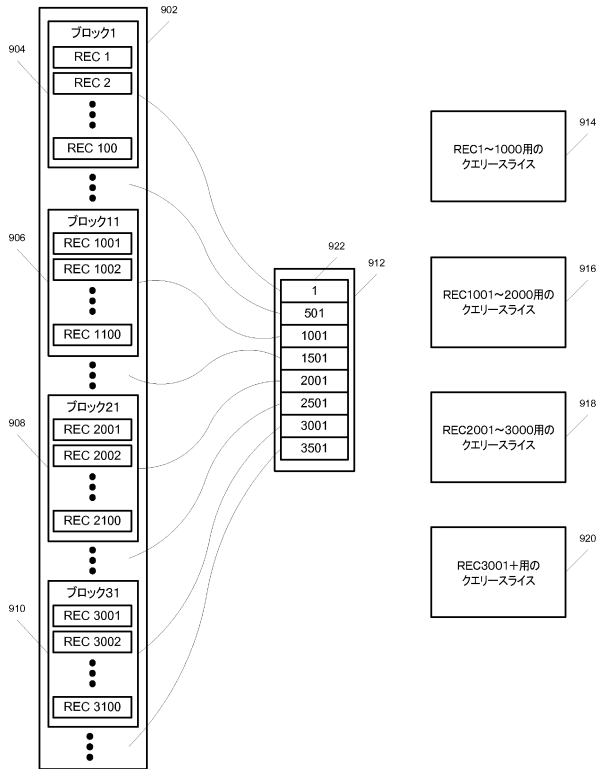


FIG. 8

【 図 9 】





---

フロントページの続き

(72)発明者 マクリーン, ジョン  
イギリス国, スコットランド ピーエー5 9 ビージェイ, エルダースライ, エルダール グローブ  
, ニューランドクレイグス アベニュー 34

審査官 早川 学

(56)参考文献 米国特許出願公開第2008/0033920 (US, A1)  
特開平08-314965 (JP, A)  
特開2006-260511 (JP, A)  
米国特許出願公開第2008/0104149 (US, A1)  
特開平8-292956 (JP, A)

(58)調査した分野(Int.Cl., DB名)  
G06F 17/30  
G06F 12/00