

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6107429号
(P6107429)

(45) 発行日 平成29年4月5日(2017.4.5)

(24) 登録日 平成29年3月17日(2017.3.17)

(51) Int.Cl.

F I

G O 6 F 17/30 (2006.01)

G O 6 F 12/00 (2006.01)

G O 6 F 13/00 (2006.01)

G O 6 F 17/30 3 3 O B

G O 6 F 17/30 1 1 O C

G O 6 F 12/00 5 1 2

G O 6 F 13/00 5 4 O E

G O 6 F 17/30 3 4 O D

請求項の数 9 (全 44 頁)

(21) 出願番号 特願2013-113984 (P2013-113984)
 (22) 出願日 平成25年5月30日(2013.5.30)
 (65) 公開番号 特開2014-232483 (P2014-232483A)
 (43) 公開日 平成26年12月11日(2014.12.11)
 審査請求日 平成28年2月26日(2016.2.26)

(73) 特許権者 000005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番
 1号
 (74) 代理人 100092152
 弁理士 服部 毅巖
 (72) 発明者 廣瀬 厚人
 神奈川県川崎市中原区上小田中4丁目1番
 1号 富士通株式会社内
 (72) 発明者 川上 俊弘
 神奈川県川崎市中原区上小田中4丁目1番
 1号 富士通株式会社内
 (72) 発明者 山本 明博
 神奈川県川崎市中原区上小田中4丁目1番
 1号 富士通株式会社内

最終頁に続く

(54) 【発明の名称】 データベースシステム、検索方法およびプログラム

(57) 【特許請求の範囲】

【請求項 1】

各々に同期されるデータを有する複数のデータベースに対応し、それぞれが前記複数のデータベースの何れかに接続された複数のサーバと、

前記複数のサーバに対して、同一の検索範囲を指定した検索要求を送信する検索要求装置と、を有し、

前記複数のサーバそれぞれは、1または2以上のキーの値を含む複数のノードが木構造に連結されたインデックス木を、検索を行うサーバの台数に応じた数のキーの値を前記複数のノードのうちのルートノードが含むように生成し、

前記複数のサーバそれぞれは、受信した前記検索要求が示す検索範囲のうち当該サーバが担当する部分検索範囲を前記インデックス木に基づいて算出し、前記部分検索範囲に限定して当該サーバに接続されたデータベースからデータを検索し、検索結果を前記検索要求装置に送信する、

データベースシステム。

【請求項 2】

前記部分検索範囲は、前記複数のサーバの間で互いに重複しないように算出する、

請求項 1 記載のデータベースシステム。

【請求項 3】

前記ルートノードは、前記サーバの台数以上の数のキーの値を含み、

前記サーバの台数が変化したとき、検索を行うサーバそれぞれは、前記サーバの台数が

10

20

変化する前に使用していた前記インデックス木の再構成を抑止する、

請求項 1 または 2 記載のデータベースシステム。

【請求項 4】

前記複数のサーバそれぞれは、当該サーバの優先順位を示す情報を有しており、

前記部分検索範囲は、前記サーバの台数と前記優先順位とに基づいて算出する、

請求項 1 乃至 3 の何れか一項に記載のデータベースシステム。

【請求項 5】

前記検索要求装置は、前記複数のサーバのうち一部のサーバから検索結果を受信し他のサーバから検索結果を受信していないときに、前記サーバの台数が変化したことを検出すると、前記一部のサーバの検索結果を破棄して前記検索要求を再送信する、

請求項 1 乃至 4 の何れか一項に記載のデータベースシステム。

【請求項 6】

前記複数のサーバのうちの第 1 のサーバは、前記検索要求が示す検索範囲のうち第 1 の部分検索範囲を算出し、前記第 1 のサーバに接続された第 1 のデータベースから、前記第 1 の部分検索範囲に対応する第 1 のデータを検索し、

前記複数のサーバのうちの第 2 のサーバは、前記検索要求が示す検索範囲のうち第 2 の部分検索範囲を算出し、前記第 2 のサーバに接続され前記第 1 のデータベースと同じデータを記憶する第 2 のデータベースから、前記第 1 の部分検索範囲に対応する前記第 1 のデータを無視して前記第 2 の部分検索範囲に対応する第 2 のデータを検索する、

請求項 1 乃至 5 の何れか一項に記載のデータベースシステム。

【請求項 7】

各々に同期されるデータを有する複数のデータベースに対応し、それぞれが前記複数のデータベースの何れかに接続された複数のサーバと、

前記複数のサーバに対して、同一の検索範囲を指定した検索要求を送信する検索要求装置と、を有し、

前記複数のサーバそれぞれは、1 または 2 以上のキーの値を含む複数のノードが木構造に連結されたインデックス木にアクセス可能であり、

前記複数のサーバそれぞれは、受信した前記検索要求が示す検索範囲のうち当該サーバが担当する部分検索範囲を、前記インデックス木のルートノードまたはルートノードから所定の深さ以内のブランチノードに登録されたキーの値を基準にして前記検索要求が示す検索範囲を分割することで算出し、前記部分検索範囲に限定して当該サーバに接続されたデータベースからデータを検索し、検索結果を前記検索要求装置に送信する、

データベースシステム。

【請求項 8】

データベースシステムが実行する検索方法であって、

各々に同期されるデータを有する複数のデータベースに対応しておりそれぞれが前記複数のデータベースの何れかに接続された複数のサーバそれぞれにおいて、1 または 2 以上のキーの値を含む複数のノードが木構造に連結されたインデックス木を、検索を行うサーバの台数に応じた数のキーの値を前記複数のノードのうちのルートノードが含むように生成し、

検索要求装置から前記複数のサーバに対して、同一の検索範囲を指定した検索要求を送信し、

前記複数のサーバそれぞれにおいて、受信された前記検索要求が示す検索範囲のうち当該サーバが担当する部分検索範囲を前記インデックス木に基づいて算出し、前記部分検索範囲に限定して当該サーバに接続されたデータベースからデータを検索し、

前記複数のサーバそれぞれから前記検索要求装置に検索結果を送信する、

検索方法。

【請求項 9】

各々に同期されるデータを有する複数のデータベースに対応しておりそれぞれが前記複数のデータベースの何れかに接続された複数のサーバの 1 つとして用いられるコンピュー

10

20

30

40

50

タに、

1 または 2 以上のキーの値を含む複数のノードが木構造に連結されたインデックス木を、検索を行うサーバの台数に応じた数のキーの値を前記複数のノードのうちのルートノードが含むように生成し、

同一の検索範囲を指定した検索要求を前記複数のサーバに送信する検索要求装置から、前記検索要求を受信し、

受信した前記検索要求が示す検索範囲のうち前記コンピュータが担当する部分検索範囲を前記インデックス木に基づいて算出し、前記部分検索範囲に限定して前記コンピュータに接続されたデータベースからデータを検索し、

前記部分検索範囲の検索結果を前記検索要求装置に送信する、
処理を実行させるプログラム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明はデータベースシステム、検索方法およびプログラムに関する。

【背景技術】

【0002】

現在、データベース管理システム（DBMS：Database Management System）と呼ばれるソフトウェアをサーバ上で動作させ、複数のクライアントからデータベースを利用できるようにしたクライアントサーバ型のデータベースシステムが広く利用されている。クライアントは検索条件を指定した検索要求をサーバに送信し、サーバは検索条件に該当するデータをデータベースから検索してクライアントに送信する。検索条件として、ある属性（テーブルのカラムに相当）の値の範囲が指定されることがある。

20

【0003】

ここで、検索要求に対するレスポンスを高速化するため、DBMSが動作するサーバを複数設けて並列に動作させる並列データベースシステムが考えられる。並列データベースシステムの構造としては、シェアードエブリシング（SE：Shared Everything）アーキテクチャとシェアードナッシング（SN：Shared Nothing）アーキテクチャとがある。

【0004】

SEアーキテクチャでは、複数のサーバが共通のデータベースに直接アクセスする。共通のデータベースは、例えば、複数のサーバからアクセス可能な共通の記憶装置上に実現される。一方、SNアーキテクチャでは、データベースのデータを予め複数のパーティションに分割しておき、各サーバは特定のパーティションにのみアクセスする。各パーティションは、例えば、他のパーティションとは異なる記憶装置上に実現される。他のパーティションのデータには、当該他のパーティションに対応する他のサーバを介して間接的にアクセスすることになる。負荷分散の観点では、SEアーキテクチャよりもSNアーキテクチャの方がアクセス競合しにくく、スループットを向上させやすい。

30

【0005】

なお、アプリケーションサーバ（APサーバ）とデータベースサーバ（DBサーバ）を含み、APサーバが複数のCPU（Central Processing Unit）コアを用いてデータを検索するデータ検索システムが提案されている。このAPサーバは、検索条件を複数の検索条件に分割し、複数のCPUコアそれぞれで動作する検索処理部に割当てて、各検索処理部は、分割された検索条件に対応するSQL文をDBサーバに送信する。そして、APサーバは、分割された検索条件に対応する部分的な検索結果のデータをマージする。

40

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開2012-59215号公報

【発明の概要】

50

【発明が解決しようとする課題】

【0007】

1つの検索要求に対して、複数のサーバが分担して検索を行い、当該検索要求に対するレスポンスを高速化することが考えられる。各サーバの検索範囲を限定する方法としては、S Nアーキテクチャのように、データベースのデータを予め複数のパーティションに分割しておく方法がある。しかし、検索を行うサーバの台数は、サーバを増設する場合やサーバが故障した場合等、並列データベースシステムの運用中に変化する可能性がある。データベースを予め分割する方法では、サーバ台数が変化したときにパーティションを再構成することになり、検索範囲の再設定の負担が大きいという問題がある。

【0008】

1つの側面では、本発明は、サーバ台数の変化が検索の並列化に与える影響を低減できるデータベースシステム、検索方法およびプログラムを提供することを目的とする。

【課題を解決するための手段】

【0009】

1つの態様では、各々に同期されるデータを有する複数のデータベースに対応し、それぞれが複数のデータベースの何れかに接続された複数のサーバと、複数のサーバに対して、同一の検索範囲を指定した検索要求を送信する検索要求装置と、を有するデータベースシステムが提供される。複数のサーバそれぞれは、受信した検索要求が示す検索範囲のうち当該サーバが担当する部分検索範囲を算出し、部分検索範囲に限定して当該サーバに接続されたデータベースからデータを検索し、検索結果を検索要求装置に送信する。

【0010】

また、1つの態様では、データベースシステムが実行する検索方法が提供される。この検索方法では、検索要求装置から、各々に同期されるデータを有する複数のデータベースに対応しておりそれぞれが複数のデータベースの何れかに接続された複数のサーバに対して、同一の検索範囲を指定した検索要求を送信する。複数のサーバそれぞれにおいて、受信された検索要求が示す検索範囲のうち当該サーバが担当する部分検索範囲を算出し、部分検索範囲に限定して当該サーバに接続されたデータベースからデータを検索する。複数のサーバそれぞれから検索要求装置に検索結果を送信する。

【0011】

また、1つの態様では、各々に同期されるデータを有する複数のデータベースに対応しておりそれぞれが複数のデータベースの何れかに接続された複数のサーバの1つとして用いられるコンピュータに実行させるプログラムが提供される。プログラムを実行するコンピュータは、同一の検索範囲を指定した検索要求を複数のサーバに送信する検索要求装置から、検索要求を受信する。受信した検索要求が示す検索範囲のうちコンピュータが担当する部分検索範囲を算出し、部分検索範囲に限定してコンピュータに接続されたデータベースからデータを検索する。部分検索範囲の検索結果を検索要求装置に送信する。

【発明の効果】

【0012】

1つの側面では、サーバ台数の変化が検索の並列化に与える影響を低減できる。

【図面の簡単な説明】

【0013】

【図1】第1の実施の形態のデータベースシステムを示す図である。

【図2】第2の実施の形態のシステムを示す図である。

【図3】サーバ装置のハードウェア例を示すブロック図である。

【図4】システムによる検索の例を示す図である。

【図5】非インデックス検索の場合に部分検索範囲を算出する例を示す図である。

【図6】インデックスの例を示す図である。

【図7】インデックス検索の場合に部分検索範囲を算出する例を示す図である。

【図8】インデックス検索の場合に部分検索範囲を算出する例を示す図（続き）である。

【図9】システムによる更新の例を示す図である。

10

20

30

40

50

- 【図 1 0】サーバ装置の故障発生時の検索処理の例を示す図である。
- 【図 1 1】サーバ装置の故障発生時の検索処理の例を示す図（続き）である。
- 【図 1 2】サーバ装置の故障発生時のインデックスの更新の例を示す図である。
- 【図 1 3】システムの機能例を示すブロック図である。
- 【図 1 4】処理要求通知の例を示す図である。
- 【図 1 5】稼働サーバリストの例を示す図である。
- 【図 1 6】処理結果通知の例を示す図である。
- 【図 1 7】稼働サーバ数と応答サーバリストの例を示す図である。
- 【図 1 8】クライアント装置による検索制御の例を示すフローチャートである。
- 【図 1 9】クライアント装置による検索制御の例を示すフローチャート（続き）である。 10
- 【図 2 0】サーバ装置による検索制御の例を示すフローチャートである。
- 【図 2 1】サーバ装置による検索制御の例を示すフローチャート（続き）である。
- 【図 2 2】限定範囲の算出処理の例を示すフローチャートである。
- 【図 2 3】非インデックス検索の場合の限定範囲の算出処理の例を示すフローチャートである。
- 【図 2 4】クライアント装置による更新制御の例を示すフローチャートである。
- 【図 2 5】サーバ装置による更新制御の例を示すフローチャートである。
- 【図 2 6】非インデックス検索により検索されるデータの例を示す図である。
- 【図 2 7】インデックス検索により検索されるデータの例を示す図である。
- 【図 2 8】インデックス検索により検索されるデータの例を示す図（続き）である。 20
- 【図 2 9】システムによる検索時間の例を示す図である。
- 【図 3 0】インデックスの変形例を示す図である。
- 【図 3 1】限定範囲の算出処理の変形例を示すフローチャートである。
- 【図 3 2】サーバ装置の機能構成の変形例を示す図である。
- 【図 3 3】システムによる検索の変形例を示す図である。
- 【発明を実施するための形態】

【 0 0 1 4 】

以下、本実施の形態について図面を参照して説明する。

〔第 1 の実施の形態〕

図 1 は、第 1 の実施の形態のデータベースシステムを示す図である。 30

【 0 0 1 5 】

第 1 の実施の形態のデータベースシステムは、サーバ 1 0 , 1 0 a を含む複数のサーバと検索要求装置 2 0 とを有する。サーバ 1 0 , 1 0 a は、例えば、検索要求装置 2 0 からアクセスを受け付けるサーバ装置（サーバコンピュータ等）である。検索要求装置 2 0 は、例えば、サーバ 1 0 , 1 0 a にアクセスするクライアント装置（クライアントコンピュータ等）であり、ユーザが操作する端末装置であってもよい。ただし、検索要求装置 2 0 は、サーバ装置であってもよいし、クライアント装置からの要求に応じてサーバ 1 0 , 1 0 a にアクセスする装置（中継装置や中継サーバと呼ばれるもの）であってもよい。サーバ 1 0 , 1 0 a と検索要求装置 2 0 とは、例えば、ネットワークを介して通信する。ただし、検索要求装置 2 0 が、複数のサーバの何れかと同じ装置上に実現されてもよい。 40

【 0 0 1 6 】

以下の情報処理を行うため、サーバ 1 0 , 1 0 a および検索要求装置 2 0 は、プロセッサおよびメモリを有してもよい。例えば、プロセッサが、メモリに記憶されたプログラムを実行する。プロセッサには、CPU、DSP (Digital Signal Processor)、ASIC (Application Specific Integrated Circuit)、FPGA (Field Programmable Gate Array) 等が含まれ得る。複数のプロセッサの集合（マルチプロセッサ）を、「プロセッサ」と呼ぶこともある。メモリには、RAM (Random Access Memory) が含まれ得る。

【 0 0 1 7 】

複数のサーバは、データベース 1 1 , 1 1 a を含む複数のデータベースに対応して設けられる。複数のデータベースは、各々に同期されるデータを有し、原則として同じ内容の 50

データを含むことになる。各サーバは、それら複数のデータベースの何れかに接続される。サーバ10はデータベース11に接続され、サーバ10aはデータベース11aに接続されている。各サーバは、当該サーバに接続されたデータベースにアクセスする一方、他のサーバに接続されたデータベースにはアクセスしない。よって、サーバ10はデータベース11にはアクセスするが、データベース11aにはアクセスしない。また、サーバ10aはデータベース11aにはアクセスするが、データベース11にはアクセスしない。

【0018】

各データベースは、例えば、他のデータベースとは異なる記憶装置上に実現される。データベースのデータを記憶する記憶装置は、HDD (Hard Disk Drive) 等の不揮発性の記憶装置でもよいし、RAM等の揮発性の記憶装置でもよい。1つのサーバと1つのデータベースとが「接続している」とき、当該データベースは当該サーバの筐体内に存在してもよいし、当該サーバの筐体外に存在してもよい。後者の場合、サーバとデータベースとが、スイッチ等の通信装置を介して接続されていてもよい。

【0019】

検索要求装置20は、サーバ10, 10aを含む複数のサーバに対して、検索条件として同一の検索範囲を指定した検索要求を送信する。例えば、検索要求装置20は、サーバ10, 10aに対して同じ検索要求21を送信する。検索要求21の送信は、マルチキャストやブロードキャストであってもよい。図1に例示した検索範囲 $X1 < C < X2$ は、属性Cの値がX1以上且つX2未満であるデータを検索することを意味する。

【0020】

複数のサーバそれぞれは、検索要求装置20から検索要求を受信すると、受信した検索要求が示す検索範囲のうち当該サーバが担当する部分検索範囲を算出する。例えば、サーバ10は、検索要求21が示す検索範囲 $X1 < C < X2$ から、部分検索範囲 $X1 < C < T$ を算出する。この部分検索範囲は、属性Cの値がX1以上且つT未満であるデータを検索することを意味する。一方、サーバ10aは、検索要求21が示す検索範囲 $X1 < C < X2$ から、サーバ10とは異なる部分検索範囲 $T < C < X2$ を算出する。この部分検索範囲は、属性Cの値がT以上且つX2未満であるデータを検索することを意味する。

【0021】

複数のサーバが算出する部分検索範囲は、好ましくは、互いに重複がないようにする。部分検索範囲の重複をなくすことで、無駄な検索処理を削減できる。このような部分検索範囲は、好ましくは、サーバ間で通信を行わず、所定のアルゴリズムに従って各サーバで独立に算出できるようにする。各サーバが他のサーバとは独立に部分検索範囲を算出することで、検索開始までのオーバーヘッドを小さくできる。例えば、検索を行うサーバそれぞれに対して、複数のサーバの中での優先順位を付与しておく。各サーバは、検索を行うサーバの台数に応じて、検索要求21が示す検索範囲を複数の部分検索範囲に分割し、当該サーバの優先順位に応じて、複数の部分検索範囲の1つを選択する。

【0022】

そして、複数のサーバそれぞれは、算出した部分検索範囲に限定して、当該サーバに接続されたデータベースからデータを検索し、部分検索範囲に対応する検索結果を検索要求装置20に送信する。例えば、サーバ10は、部分検索範囲 $X1 < C < T$ に該当するデータをデータベース11から検索し、部分検索範囲 $X1 < C < T$ の検索結果12を検索要求装置20に送信する。サーバ10aは、部分検索範囲 $T < C < X2$ に該当するデータをデータベース11aから検索し、部分検索範囲 $T < C < X2$ の検索結果12aを検索要求装置20に送信する。検索要求装置20は、例えば、サーバ10から受信した検索結果12とサーバ10aから受信した検索結果12aとをマージする。

【0023】

上記のような並列検索では、データ全体が仮想的に複数のパーティションに分割されているということもできる。例えば、データベース11, 11aに同じデータが記憶されていても、サーバ10からは一部の仮想パーティションのデータのみがデータベース11に記憶されているように見える。同様に、サーバ10aからは、データベース11と異なる

10

20

30

40

50

一部の仮想パーティションのデータのみがデータベース 11a に記憶されているように見える。これにより、サーバ 10, 10a が検索を分担することができる。

【0024】

一方で、パーティションが仮想的であるため、サーバの増設や故障等により検索を行うサーバの台数が変化しても、データベース間でデータを移動しなくてもよく、パーティションを再構成する負荷を抑制できる。例えば、各サーバがサーバ台数をパラメータとするアルゴリズムに従って部分検索範囲を算出する場合、サーバ台数の変化に応じて、当該サーバが認識する仮想パーティションが自動的に調整されることになる。

【0025】

第1の実施の形態のデータベースシステムによれば、データを検索するときサーバ 10, 10a が異なるデータベースにアクセスするため、共通のデータベースにアクセスする場合と比べて、アクセス競合を抑制できスループットが向上しやすくなる。また、データベース 11, 11a のデータが同期されるため、あるサーバが故障しても、残りのサーバを用いてデータ全体へのアクセスを確保でき、耐故障性が向上する。

【0026】

また、サーバ 10, 10a それぞれが、検索要求装置 20 から指定された検索範囲のうち担当する部分検索範囲を算出し、部分検索範囲に限定してデータを検索するため、検索処理をサーバ 10, 10a が効率的に分担することができる。また、検索を行うサーバの台数が変化しても、データベース 11, 11a の間でデータを移動しなくてもよく、サーバ 10, 10a に検索処理を分散させるための再設定の負担が軽減される。

【0027】

[第2の実施の形態]

図2は、第2の実施の形態のシステムを示す図である。

第2の実施の形態のシステムは、サーバ装置 100, 100a, 100b およびクライアント装置 200 を有する。サーバ装置 100, 100a, 100b は、第1の実施の形態のサーバ 10, 10a の一例である。クライアント装置 200 は、第1の実施の形態の検索要求装置 20 の一例である。クライアント装置 200 は、ネットワーク 30 を経由してサーバ装置 100, 100a, 100b と接続している。なお、サーバ装置の数は、2台でもよいし、4台以上でもよい。また、各サーバ装置はクラスタ化されていてもよい。

【0028】

クライアント装置 200 は、例えば、デスクトップ型コンピュータやノート型コンピュータ等、ユーザが操作するコンピュータである。クライアント装置 200 は、アプリケーションプログラムを実行する。クライアント装置 200 が実行するアプリケーションプログラムは、サーバ装置 100, 100a, 100b の有するデータベースのデータを用いてアプリケーションプログラムを実行する。クライアント装置 200 は、例えば、データベースのデータの更新や検索を要求する通知をサーバ装置 100, 100a, 100b へ一斉送信する。なお、クライアント装置 200 は、他の端末装置にサービスを提供するアプリケーションサーバとして用いられてもよい。

【0029】

サーバ装置 100, 100a, 100b は、データベースを管理するサーバコンピュータである。サーバ装置 100, 100a, 100b それぞれは、個別にデータベースを有する。また、各サーバ装置のデータベースのデータは、サーバ装置間で同期がとられた状態である。また、サーバ装置 100, 100a, 100b それぞれは、クライアント装置 200 の要求に応じて、データベースのデータへのアクセス処理を実行し、実行結果をクライアント装置 200 に応答する。

【0030】

なお、各サーバ装置のデータベースは、HDD にデータを配置するディスク型データベースでもよいし、主記憶上にデータを配置して高速にアクセスを行うインメモリ型データベースでもよい。また、同期がとられたデータベースそれぞれは、各サーバ装置に対応付けられて外部に接続されたデータベースサーバ装置に記憶されていてもよい。

【 0 0 3 1 】

全サーバ装置のうち、1台のサーバ装置の種別には“現用”が割当てられ、他のサーバ装置の種別には“待機1”、“待機2”等が割当てられる。以下、種別が“現用”であるサーバ装置を“現用系のサーバ装置”と記載し、種別が“現用”以外であるサーバ装置を“待機系のサーバ装置”と記載する場合がある。

【 0 0 3 2 】

現用系のサーバ装置は、更新処理および検索処理を行う。更新処理には、例えば、データベースのデータの登録や削除等、データベースのデータを更新する処理を言う。検索処理は、データベースからデータを検索する処理を言う。また、現用系のサーバ装置は、更新処理の実行後、他のサーバ装置に自己のデータベースの更新後のデータと同期をとるよう要求する。また、各サーバ装置は、並列にデータベースからデータを検索する。

10

【 0 0 3 3 】

現用系のサーバ装置が故障した場合、待機系のサーバ装置の何れか1台のサーバ装置（例えば、種別が“待機1”であるサーバ装置）が現用系のサーバ装置となる。

図3は、サーバ装置のハードウェア例を示すブロック図である。サーバ装置100は、プロセッサ101、RAM102、HDD103、画像信号処理部104、入力信号処理部105、ディスクドライブ106および通信インタフェース107を有する。これらのユニットは、サーバ装置100内でバス108に接続されている。

【 0 0 3 4 】

プロセッサ101は、プログラムの命令を実行する演算器を含むプロセッサであり、例えばCPUである。プロセッサ101は、HDD103に記憶されているプログラムやデータの少なくとも一部をRAM102にロードしてプログラムを実行する。なお、プロセッサ101は複数のプロセッサコアを備えてもよい。また、サーバ装置100は、複数のプロセッサを備えてもよい。また、サーバ装置100は、複数のプロセッサまたは複数のプロセッサコアを用いて並列処理を行ってもよい。また、2以上のプロセッサの集合、FPGAやASIC等の専用回路、2以上の専用回路の集合、プロセッサと専用回路の組み合わせ等を「プロセッサ」と呼んでもよい。

20

【 0 0 3 5 】

RAM102は、プロセッサ101が実行するプログラムやプログラムから参照されるデータを一時的に記憶する揮発性メモリである。なお、サーバ装置100は、RAM以外の種類のメモリを備えてもよく、複数個の揮発性メモリを備えてもよい。

30

【 0 0 3 6 】

HDD103は、OS（Operating System）やファームウェアやアプリケーションソフトウェア等のソフトウェアのプログラムおよびデータを記憶する不揮発性の記憶装置である。なお、サーバ装置100は、フラッシュメモリ等の他の種類の記憶装置を備えてもよく、複数個の不揮発性の記憶装置を備えてもよい。

【 0 0 3 7 】

画像信号処理部104は、プロセッサ101からの命令に従って、サーバ装置100に接続されたディスプレイ41に画像を出力する。ディスプレイ41としては、CRT（Cathode Ray Tube）ディスプレイや液晶ディスプレイ等を用いることができる。

40

【 0 0 3 8 】

入力信号処理部105は、サーバ装置100に接続された入力デバイス42から入力信号を取得し、プロセッサ101に通知する。入力デバイス42としては、マウスやタッチパネル等のポインティングデバイス、キーボード等を用いることができる。

【 0 0 3 9 】

ディスクドライブ106は、記録媒体43に記録されたプログラムやデータを読み取る駆動装置である。記録媒体43として、例えば、フレキシブルディスク（FD：Flexible Disk）やHDD等の磁気ディスク、CD（Compact Disc）やDVD（Digital Versatile Disc）等の光ディスク、光磁気ディスク（MO：Magneto-Optical disk）を使用できる。ディスクドライブ106は、プロセッサ101からの命令に従って、記録媒体43から

50

読み取ったプログラムやデータをRAM 102またはHDD 103に格納する。

【0040】

通信インタフェース107は、ネットワーク30等のネットワークを介して他の情報処理装置（例えば、クライアント装置200等）と通信を行う。

なお、サーバ装置100はディスクドライブ106を備えていなくてもよく、専ら他の端末装置から制御される場合には、画像信号処理部104や入力信号処理部105を備えていなくてもよい。また、ディスプレイ41や入力デバイス42は、サーバ装置100の筐体と一体に形成されていてもよい。

【0041】

なお、サーバ装置100a、100bおよびクライアント装置200も、サーバ装置100と同様のハードウェアを用いて実現できる。

図4は、システムによる検索の例を示す図である。サーバ装置100は、サーバ制御部110およびデータベース120を有する。サーバ装置100aは、サーバ制御部110aおよびデータベース120aを有する。サーバ装置100bは、サーバ制御部110bおよびデータベース120bを有する。サーバ装置100の種別に“現用”が割当てられ、サーバ装置100aの種別に“待機1”が割当てられ、サーバ装置100bの種別に“待機2”が割当てられている。

【0042】

データベース120、120a、120bは、クライアント装置200が用いるデータを記憶する。データベース120、120a、120bのデータは、同期がとられた状態である。データベース120、120a、120bではそれぞれ、データが1または2以上のテーブルにより管理されている。

【0043】

サーバ制御部110、110a、110bは、テーブルに記憶されているデータへのアクセス処理を実行し、処理結果をクライアント装置200に送信する。検索処理の処理結果には、検索されたデータや検索件数を示す情報が含まれる。更新処理の処理結果には、更新処理の成否を示す情報が含まれる。

【0044】

また、サーバ制御部110、110a、110bは、クライアント装置200から検索処理を要求されたとき、まず、自サーバ装置の種別や検索に参加するサーバ装置の数や検索条件の示す検索範囲等に基づいて部分検索範囲を算出する。部分検索範囲は、検索条件の示す検索範囲のうち各サーバ装置が検索を担当する検索範囲である。また、部分検索範囲は、サーバ装置間で範囲が重複しないように算出される。そして、サーバ制御部110、110a、110bは、算出した部分検索範囲について検索処理を実行する。

【0045】

また、サーバ制御部110、110a、110bは、クライアント装置200から更新処理を要求されたとき、自サーバ装置の種別に基づいてデータの更新を行うか否か判定する。

【0046】

クライアント装置200は、アプリケーションソフトウェア210およびクライアント制御部220を有する。アプリケーションソフトウェア210は、データベース120、120a、120bのデータを用いて所定の情報処理を実行する。

【0047】

クライアント制御部220は、サーバ装置100、100a、100bにデータの検索や更新等を要求し、各サーバ装置から処理結果を受信する。そして、受信した処理結果をマージし、アプリケーションソフトウェア210に出力する。

【0048】

以下、第2の実施の形態のシステムによる検索の例をステップ番号に沿って説明する。

まず、アプリケーションソフトウェア210は、クライアント制御部220にデータベースからのデータの検索を要求する（S1）。次に、クライアント装置200は、サーバ

10

20

30

40

50

装置 100, 100a, 100b に、データベースからのデータの検索要求を一斉送信する。検索要求の一斉送信は、ブロードキャストやマルチキャストとして実行してもよい。例えば、テーブル名が“T01”であるテーブルについてデータの検索を要求する(S2, S2a, S2b)。以下、テーブル名が“T01”であるテーブルを“テーブル(T01)”と記載する場合がある。

【0049】

次に、サーバ制御部 110, 110a, 110b それぞれは、自己のサーバ装置の種別と検索に参加するサーバ装置の数に基づき部分検索範囲を算出する。例えば、サーバ制御部 110 は部分検索範囲として“部分検索範囲 # 1”を算出し、サーバ制御部 110a は部分検索範囲として“部分検索範囲 # 2”を算出し、サーバ制御部 110b は部分検索範囲として“部分検索範囲 # 3”を算出する。そして、サーバ制御部 110, 110a, 110b それぞれは、自己のテーブルから算出された部分検索範囲についてデータを検索する(S3, S3a, S3b)。

10

【0050】

次に、サーバ制御部 110, 110a, 110b それぞれは、検索結果をクライアント装置 200 に送信する(S4, S4a, S4b)。そして、クライアント制御部 220 は、サーバ装置 100, 100a, 100b から受信した検索結果をマージし、マージした検索結果をアプリケーションソフトウェア 210 に通知する(S5)。

【0051】

次に、図 5 ~ 8 を用いて、インデックスを用いて検索する場合とインデックスを用いずに検索する場合において、部分検索範囲を算出する例について説明する。以下、インデックスを用いて検索することをインデックス検索と記載する場合がある。また、インデックスを用いずに検索することを非インデックス検索と記載する場合がある。

20

【0052】

サーバ制御部 110 は、インデックス検索により検索処理を実行する場合と、非インデックス検索により検索処理を実行する場合とがある。インデックス検索により実行する場合は、例えば、検索条件に含まれているテーブルの特定の列に対応するインデックスが存在する場合は挙げられる。この際、インデックスには、B - T r e e インデックス等の木構造のインデックスが用いられる。非インデックス検索により実行する場合は、例えば、検索条件に含まれているテーブルの特定の列に対応するインデックスが存在しない場合は挙げられる。

30

【0053】

図 5 は、非インデックス検索の場合に部分検索範囲を算出する例を示す図である。サーバ制御部 110 は、非インデックス検索の場合、クライアント装置 200 が指定した検索条件の示す検索範囲を、検索に参加するサーバ装置の数で分割することで部分検索範囲を算出する。

【0054】

例えば、クライアント装置 200 に、サーバ装置 100, 100a が接続されている。この状態で、クライアント制御部 220 が、テーブル(T01)から“ $10 < C01 < 100$ ”の検索条件を満たすデータの検索を、サーバ装置 100, 100a それぞれに要求したとする。“C01”は、テーブル(T01)の有する列である。

40

【0055】

このとき、検索条件の下限値である“10”と上限値である“100”の中間の値は“55”であるため、“ $10 < C01 < 55$ ”が、部分検索範囲 # 1 としてサーバ制御部 110 により算出される。また、“ $55 < C01 < 100$ ”が、部分検索範囲 # 2 としてサーバ制御部 110a により算出される。

【0056】

図 6 は、インデックスの例を示す図である。各サーバ装置は、データベースのテーブルからデータを検索する際、検索条件に含まれる列に対応するインデックスを用いる。テーブル 121 は、データベース 120 のデータを格納する。テーブル 121 は、列名が“C

50

“ 0 1 ” である列および列名が “ C 0 2 ” である列を有する。以下、列名が “ C 0 1 ” である列を “ 列 (C 0 1) ” と記載する場合がある。

【 0 0 5 7 】

サーバ装置 1 0 0 は、例えば、テーブル 1 2 1 の列 (C 0 1) が特定の検索条件を満たすデータをテーブル 1 2 1 から検索するとき、列 (C 0 1) の値をキーに設定したインデックス 1 3 1 を用いる。なお、各キーには、複数の列の値の組み合わせが設定されてもよい。

【 0 0 5 8 】

インデックス 1 3 1 は、データベース 1 2 0 に記憶されている。インデックス 1 3 1 は、テーブル 1 2 1 に格納された列 (C 0 1) の値の集合に対応して一意に生成される。そのため、同期がとられた複数のデータベースからは、同一のインデックスが生成される。すなわち、サーバ装置 1 0 0 以外の各サーバ装置も、テーブル 1 2 1 と同期がとられたテーブルを有するため、インデックス 1 3 1 と同一のインデックスを有することになる。

【 0 0 5 9 】

インデックス 1 3 1 には、例えば、B - T r e e インデックスが用いられる。B - T r e e インデックスは、木構造のデータ構造であるB木を用いたインデックスである。

インデックス 1 3 1 は、キー # 1 , # 2 を含むルートノード、キー # 3 を含むブランチノード、キー # 4 を含むブランチノード、キー # 5 を含むリーフノード、キー # 6 を含むリーフノード、キー # 7 を含むリーフノード、キー # 8 を含むリーフノードを有する。

【 0 0 6 0 】

ルートノードは、木構造の頂点にあたるノードである。ルートノードは、1または2以上のキーおよび、2以上のポイントを有する。ルートノードのポイント1つは、ブランチノード1つを指し示す。ルートノードはルートブロックと呼ばれることもある。ブランチノードは、ルートノードとリーフノードとの間にある中間のノードである。ブランチノードは、1または2以上のキー、および2以上のポイントを有する。ブランチノードのポイント1つは、他のブランチノードまたはリーフノード1つを指し示す。ブランチノードはブランチブロックと呼ばれることもある。リーフノードは、木構造の終端にあたるノードである。リーフノードは、1または2以上のキーおよび木構造の終端を示す情報（例えば、“ N U L L ”）が設定された2以上のポイントを有する。リーフノードはリーフブロックと呼ばれることもある。

【 0 0 6 1 】

各ノードには、最大でk個のキーおよびk + 1個のポイントが、ポイントを先頭に交互に配置される。例えば、インデックス 1 3 1 において、ルートノードには、ポイント # 1 1、キー # 1、ポイント # 1 2 およびキー # 2 が順番に配置されている。また、キー # 3 を含むブランチノードは、ポイント # 1 3、キー # 3 およびポイント # 1 4 が順番に配置されている。また、キー # 4 を含むブランチノードには、ポイント # 1 5、キー # 4 およびポイント # 1 6 が順番に配置されている。

【 0 0 6 2 】

各ノードのキーの値は、昇順に並んでいる。例えば、インデックス 1 3 1 において、キー # 1 には “ 2 4 ” が設定され、キー # 2 には “ 4 6 ” が設定されている。

各キーは、キーに対応するレコードへのポイントを有する。例えば、キー # 1 はレコード (C 0 1 = 2 4) へのポイントを有し、キー # 4 はレコード (C 0 1 = 3 6) へのポイントを有し、キー # 7 はレコード (C 0 1 = 2 9) へのポイントを有する。

【 0 0 6 3 】

ポイント # 1 1 , # 1 2 , # 1 3 , # 1 4 は、1つ下の階層のノードを指し示す。各ポイントの指し示す1つ下のノードには、ポイントの前のキーより大きい値、かつ、ポイントの後ろのキーより小さい値のキーが配置される。

【 0 0 6 4 】

例えば、ポイント # 1 1 は、キー # 3 を含むブランチノードを指し示し、キー # 3 にはキー # 1 の値より小さい “ 1 3 ” が設定されている。ポイント # 1 2 は、キー # 4 を含む

10

20

30

40

50

ブランチノードを指し示し、キー # 4 にはキー # 1 の値より大きくキー # 2 の値より小さい “ 3 6 ” が設定されている。また、ポインタ # 1 3 はキー # 5 を含むリーフノードを指し示し、キー # 5 にはキー # 3 の値よりも小さい “ 6 ” が設定されている。ポインタ # 1 4 はキー # 6 を含むリーフノードを指し示し、キー # 6 にはキー # 3 の値よりも大きい “ 1 7 ” が設定されている。ポインタ # 1 5 はキー # 7 を含むリーフノードを指し示し、キー # 7 にはキー # 4 の値よりも小さい “ 2 9 ” が設定されている。ポインタ # 1 6 はキー # 8 を含むリーフノードを指し示し、キー # 8 にはキー # 4 の値よりも大きい “ 4 0 ” が設定されている。

【 0 0 6 5 】

このように、B - T r e e インデックスでは、キーの値の範囲が、ルートノードからリーフノードに向かって階層的に分割される。1つ下の階層のノードを指し示すポインタは、キーの値の範囲1つに対応していると言える。図6の例の場合、ルートノードにおいて、列 (C 0 1) の値が “ C 0 1 < 2 4 ” , “ 2 4 ” , “ 2 4 < C 0 1 < 4 6 ” , “ 4 6 ” を含む複数の範囲に分割されていると言える。また、ブランチノードにおいて、“ C 0 1 < 2 4 ” の範囲が更に “ C 0 1 < 1 3 ” , “ 1 3 ” , “ 1 3 < C 0 1 < 2 4 ” に分割され、“ 2 4 < C 0 1 < 4 6 ” の範囲が更に “ 2 4 < C 0 1 < 3 6 ” , “ 3 6 ” , “ 3 6 < C 0 1 < 4 6 ” に分割されていると言える。インデックス 1 3 1 のルートノードのキーを部分検索範囲の境界値とすることで、各サーバ装置が担当する部分検索範囲の検索負荷 (例えば、検索されるレコードの数) が、ほぼ均等になるものと期待できる。

【 0 0 6 6 】

1つの値を指定したインデックス検索は、例えば、次のように行う。

まず、ルートノードについて、検索する値と一致するキーがルートノードに含まれているか判定する。検索する値と一致するキーがある場合、当該キーのポインタの指し示すレコードを抽出する。

【 0 0 6 7 】

検索する値と一致するキーがない場合、ルートノードについて、検索する値より大きいキーのうち最も左側 (最も手前) にあるキーを特定する。検索する値より大きいキーが特定された場合、特定されたキーの1つ前のポインタを選択する。検索する値より大きいキーが特定されなかった場合、末尾のポインタを選択する。

【 0 0 6 8 】

次に、選択したポインタが指し示すブランチノードについて、ルートノードと同様に、検索する値と一致するキーが当該ブランチノードに含まれているか判定する。検索する値と一致するキーがある場合、当該キーのポインタの指し示すレコードを抽出する。検索する値と一致するキーがない場合、ルートノードと同様、着目しているブランチノードについて、検索する値より大きいキーのうち最も左側にあるキーを特定する。検索する値より大きいキーが特定された場合、特定されたキーの1つ前のポインタを選択する。検索する値より大きいキーが特定されなかった場合、末尾のポインタを選択する。次に、選択したポインタの指し示すノードが1つ下の階層のブランチノードの場合、1つ上の階層のブランチノードと同様に、レコード抽出し、あるいは、ポインタを選択する。

【 0 0 6 9 】

選択したポインタが指し示すノードがリーフノードの場合、そのリーフノードについて、検索する値と一致するキーがリーフノードに含まれているか判定する。検索する値と一致するキーがある場合、当該キーのポインタの指し示すレコードを抽出する。検索する値と一致するキーがない場合、例えば、検索条件に該当するレコードが検索されなかったとして終了する。

【 0 0 7 0 】

例えば、レコード (C 0 1 = 2 9) をインデックス検索する場合、まず、ルートノードについて、“ 2 9 ” より大きい最も左側にあるキー # 2 が特定され、そのキーの1つ前のポインタ # 1 2 が選択される。これは、検索範囲を “ 2 4 < C 0 1 < 4 6 ” に絞ることを意味する。次に、ポインタ # 1 2 が指し示すブランチノードについて、“ 2 9 ” より大き

10

20

30

40

50

いキー # 4 が特定され、キー # 4 の 1 つ前のポインタ # 1 5 が選択される。これは、検索範囲を “ 2 4 < C 0 1 < 3 6 ” に絞ることを意味する。次に、ポインタ # 1 5 が指し示すリーフノードについて、値が “ 2 9 ” であるキー # 7 が特定され、特定されたキー # 7 のポインタが指し示すレコード (C 0 1 = 2 9) が抽出される。

【 0 0 7 1 】

インデックスを更新する場合、ルートノードからリーフノードに向かって、削除するキーを含むノードまたは追加するキーを格納すべきノードを探す。例えば、レコード (C 0 1 = 2 9) が削除された場合、キー # 1 を含むルートノード、キー # 4 を含むブランチノードを経由して検索されたキー # 7 を削除する。また、例えば、その後、レコード (C 0 1 = 3 1) が登録されたときに “ 3 1 ” であるキーを何れかのノードに追加する場合、キー # 1 を含むルートノード、キー # 4 を含むブランチノードを経由して、ポインタ # 1 5 が指し示すリーフノードに値が “ 3 1 ” であるキーを追加する。

10

【 0 0 7 2 】

B - T r e e インデックスは、1 つのノードが有するポインタの数の最大値を j とすると、ブランチノードの有するポインタは少なくとも $j / 2$ となる特徴を有する。例えば、1 つのノードが有するポインタの数を最大 “ 2 ” とすると、各ブランチノードは少なくとも 1 つのポインタを有することとなる。これにより、リーフノードの階層の深さが均等化され、インデックス 1 3 1 を用いたデータの検索回数が最大でも、対数オーダーの階層の深さ以内の回数となる。よって、データの検索速度が安定する。

【 0 0 7 3 】

20

なお、インデックス 1 3 1 には、例えば、B * T r e e インデックスや B + T r e e インデックス等、他の木構造のインデックスを用いてもよい。

図 7 は、インデックス検索の場合に部分検索範囲を算出する例を示す図である。各サーバ装置は、検索に参加するサーバ装置の数が 2 以上である場合は、ルートノードのキーの数が “ 検索に参加するサーバ装置の数 - 1 ” になるようにインデックスを生成する。検索に参加するサーバ装置の数が 1 である場合、上記の方法でインデックスを生成するとルートノードのキーの数が 0 になるため、ルートノードのキーの数が 1 以上の任意の数となるようにインデックスを生成する。また、各サーバ装置のサーバ制御部は、次のように、部分検索範囲を算出する。

【 0 0 7 4 】

30

まず、各サーバ装置のサーバ制御部は、自己のサーバ装置の種別に対応するルートノードのキーが境界値となるように、仮想パーティションに含まれるデータの範囲を算出する。仮想パーティションは、検索に参加するサーバ装置の数でデータベースのデータを分割したときのデータの集合である。データベースのデータは複数のサーバ装置の間で同期されている (複数のデータベースが同じデータを含んでいる) ことから、このパーティションはデータベースを物理的に分割したものではなく、仮想的に分割したものであると言える。各サーバ装置は、データベースのデータを検索する際、仮想パーティションに含まれるデータに限定して検索する。すなわち、検索処理の際、各サーバ装置のデータベースには、仮想パーティションの範囲に限定してデータが格納されているように見える。以下、仮想パーティションに含まれるデータの範囲を “ 限定範囲 ” と記載する場合がある。そして、各サーバ装置のサーバ制御部は、算出した限定範囲と、検索条件の示す検索範囲との重複する範囲を部分検索範囲として算出する。

40

【 0 0 7 5 】

例えば、図 7 では、クライアント装置 2 0 0 に、サーバ装置 1 0 0 , 1 0 0 a が接続されている。サーバ装置 1 0 0 , 1 0 0 a それぞれは、同期がとられたテーブル (T 0 1) を有する。また、インデックス 1 3 1 はサーバ装置 1 0 0 に記憶され、インデックス 1 3 1 a はサーバ装置 1 0 0 a に記憶されている。

【 0 0 7 6 】

また、インデックス 1 3 1 , 1 3 1 a それぞれのルートノードは、 “ 2 4 ” が設定されているキーを有する。検索に参加するサーバ装置の数は 2 であるため、インデックス 1 3

50

1, 131 aそれぞれのルートノードのキーの数は、“ $2 - 1 = 1$ ”となる。

【0077】

この状態で、テーブル(T01)について、“ $10 < C01 < 100$ ”の検索条件を満たすデータの検索を、クライアント制御部220がサーバ装置100, 100 aへ要求したとする。すると、サーバ制御部110, 110 aにより、“ $C01 < 24$ ”および“ $24 < C01$ ”の限定範囲が算出される。例えば、“ $C01 < 24$ ”である限定範囲は、サーバ制御部110により算出され、“ $24 < C01$ ”である限定範囲はサーバ制御部110 aにより算出される。

【0078】

そして、算出された限定範囲と要求された検索条件の示す検索範囲との重複する範囲が、部分検索範囲としてサーバ制御部110, 110 aにより算出される。例えば、要求された検索条件の示す検索範囲と、算出された“ $C01 < 24$ ”の重複する“ $10 < C01 < 24$ ”が、部分検索範囲#1としてサーバ制御部110により算出される。また、要求された検索条件の示す検索範囲と、算出された“ $24 < C01$ ”の重複する“ $24 < C01 < 100$ ”が、部分検索範囲#2としてサーバ制御部110 aにより算出される。

【0079】

なお、インデックス検索の場合における検索範囲の算出方法では、クライアント装置200から要求された検索条件の示す検索範囲が、ルートノードのキーで分割された検索範囲に含まれない場合がある。例えば、図7のインデックス131を用いて検索条件の示す“ $30 < C01 < 100$ ”の検索範囲を分割する場合、検索条件の示す検索範囲と、サーバ装置100が算出した限定範囲(“ $C01 < 24$ ”)とが重複しないため、サーバ装置100の担当する部分検索範囲が存在しないことになる。

【0080】

この場合、少なくとも1つのキーが検索条件の示す検索範囲に含まれるブランチノードを選択し、選択したブランチノードのキーを境界値として部分検索範囲を算出してもよい。このブランチノードとしては、検索条件に該当するキーを有するブランチノードのうち、最もルートノードに近いもの(ルートノードからの深さが最も小さいもの)を選択することが好ましい。例えば、“ $30 < C01 < 100$ ”の検索条件を指定した場合、ルートノードのキーは検索条件に含まれないため、少なくとも1つのキーが検索条件に含まれているブランチノードを選択する。その結果、“36”が設定されたキーを含むブランチノードが選択される。そして、“36”を境界値とした“ $30 < C01 < 36$ ”および“ $36 < C01 < 100$ ”の部分検索範囲が算出される。

【0081】

なお、選択するブランチノードについて、検索条件の示す検索範囲に含まれるキーの数が多ければ、多くのサーバ装置に空でない部分検索範囲が割当てられることになる。そのため、検索に参加するサーバ装置が3以上あるとき、選択するブランチノードには、検索条件の示す検索範囲内のキーが多く含まれることが好ましい。ただし、検索範囲が過度に小さい部分検索範囲に分割されないように、選択するブランチノードを、ルートノードから所定の深さ以内にあるものに限定してもよい。

【0082】

図8は、インデックス検索の場合に部分検索範囲を算出する例を示す図(続き)である。

例えば、図8では、クライアント装置200に、サーバ装置100, 100 a, 100 bが接続されている。サーバ装置100, 100 a, 100 bそれぞれは、同期がとられたテーブル(T01)を有する。インデックス131, 131 a, 131 bそれぞれは、テーブル(T01)の列に対応するインデックスである。また、インデックス131はサーバ装置100に記憶され、インデックス131 aはサーバ装置100 aに記憶され、インデックス131 bはサーバ装置100 bに記憶されている。

【0083】

インデックス131, 131 a, 131 bそれぞれのルートノードは、“24”が設定

10

20

30

40

50

されているキーと、“46”が設定されているキーとを有する。検索に参加するサーバ装置の数は3であるため、インデックス131, 131a, 131bのルートノードのキーの数は、“3 - 1 = 2”となる。

【0084】

この状態で、テーブル(T01)について、“10 < C01 < 100”の検索条件を満たすデータの検索を、クライアント制御部220がサーバ装置100, 100a, 100bそれぞれへ要求したとする。すると、サーバ制御部110, 110a, 110bにより、“C01 < 24”、“24 C01 < 46”および“46 C01”の限定範囲が算出される。例えば、“C01 < 24”の限定範囲はサーバ制御部110により算出され、“24 C01 < 46”である限定範囲はサーバ制御部110aにより算出され、“46 C01”である限定範囲はサーバ制御部110bにより算出される。

10

【0085】

そして、算出された限定範囲と要求された検索条件の示す検索範囲との重複する範囲が、部分検索範囲としてサーバ制御部110, 110a, 110bにより算出される。例えば、要求された検索条件の示す検索範囲と、算出された“C01 < 24”の重複する“10 < C01 < 24”が、部分検索範囲#1としてサーバ制御部110により算出される。また、要求された検索条件の示す検索範囲と、算出された“24 C01 < 46”の重複する“24 C01 < 46”が、部分検索範囲#2としてサーバ制御部110aにより算出される。また、要求された検索条件の示す検索範囲と、算出された“46 C01”の重複する“46 C01 < 100”が、部分検索範囲#3としてサーバ制御部110bにより算出される。

20

【0086】

図5で説明した方法により、部分検索範囲を算出すると、特定の範囲にデータが偏って存在する場合、部分検索範囲の間で検索負荷が均等にならない場合がある。図6～8で説明したように、各サーバ装置は、B - T r e eインデックスのルートノードのキーを境界値とすることで、部分検索範囲の間の検索負荷をほぼ均等化できる。

【0087】

図9は、システムによる更新の例を示す図である。サーバ装置100のデータベース120は、インデックス131を記憶している。同様に、サーバ装置100aのデータベース120aは、インデックス131aを記憶し、サーバ装置100bのデータベース120bは、インデックス131bを記憶している。ここでは、現用系のサーバ装置のみがクライアント装置200からの更新要求に応答し、待機系のサーバ装置は更新要求に応答しない。以下、第2の実施の形態のシステムによる更新の例をステップ番号に沿って説明する。

30

【0088】

まず、アプリケーションソフトウェア210は、クライアント制御部220にデータの更新を要求する(S11)。クライアント制御部220は、サーバ装置100, 100a, 100bにデータベースのデータの更新要求を一斉送信する(S12, S12a, S12b)。

【0089】

次に、サーバ制御部110, 110a, 110bそれぞれは、自サーバ装置の種別に基づき更新処理を行うか判定する。具体的には、自サーバ装置が現用系である場合、更新処理を行うと判定し、自サーバ装置が待機系である場合、更新処理を行わないと判定する。図9では、現用系のサーバ装置であるサーバ装置100がデータベース120のデータについて更新処理を実行する(S13)。

40

【0090】

次に、サーバ制御部110は、待機系の各サーバ装置にデータベース120のデータと同期をとるように要求する。同期の要求を受信した各サーバ装置は、自己のデータベースのデータについて、データベース120のデータと同期をとる。例えば、サーバ装置100aはデータベース120aのデータについて同期をとり(S14)、サーバ装置100

50

b は、データベース 120 b のデータについて同期をとる (S14a)。

【0091】

次に、各サーバ装置は、同期がとられたデータベースのテーブルの特定の列に対応するインデックスを更新する。具体的には、サーバ装置 100 はデータベース 120 のテーブルの特定の列に対応するインデックス 131 を更新し (S15)、サーバ装置 100 a はデータベース 120 a のテーブルの特定の列に対応するインデックス 131 a を更新し (S15a)、サーバ装置 100 b はデータベース 120 b のテーブルの特定の列に対応するインデックス 131 b を更新する (S15b)。

【0092】

次に、サーバ制御部 110 は、更新結果をクライアント装置 200 に送信する (S16) 。そして、クライアント制御部 220 は、受信した更新結果をアプリケーションソフトウェア 210 に出力する (S17) 。

【0093】

なお、図 9 では、サーバ装置 100 は、データベース 120、120 a、120 b の同期の後に更新結果をクライアント装置 200 に送信したが、データベース 120、120 a、120 b の同期の前に送信してもよい。また、インデックス 131、131 a、131 b の更新は、サーバ装置 100 がクライアント装置 200 に更新結果を送信した後に行われてもよい。

【0094】

次に、図 10 ~ 12 を用いて、現用系のサーバ装置が故障した場合の動作について説明する。

図 10 は、サーバ装置の故障発生時の検索処理の例を示す図である。図 10 では、図 4 と同様のシステム構成である場合について説明する。この場合、サーバ装置 100、100 a、100 b それぞれでは、自サーバ装置の種別に応じて部分検索範囲 #1、#2、#3 の何れかが算出される。この状態で、サーバ装置 100 が故障したとする。

【0095】

図 11 は、サーバ装置の故障発生時の検索処理の例を示す図 (続き) である。サーバ装置 100 が故障すると、サーバ装置 100、100 a、100 b の集合がサーバ装置 100 a、100 b の集合に縮退して、クライアント装置 200 からの要求を処理することになる。このとき、種別が “待機 1” であるサーバ装置 100 a の種別が “現用” に変更され、種別が “待機 2” であるサーバ装置 100 b の種別が “待機 1” に変更される。また、サーバ装置 100 a、100 b それぞれには、自サーバ装置の種別に応じて部分検索範囲 #1、#2 の何れかが算出される。検索を行うサーバ装置の数が変わったことに伴って、各サーバ装置が認識する仮想パーティションが変わり、限定範囲が変わる。よって、同じ検索範囲を指定した検索要求が受信された場合であっても、各サーバ装置は、サーバ装置 100 が故障する前とは異なる部分検索範囲を算出することになる。

【0096】

図 12 は、サーバ装置の故障発生時のインデックスの更新の例を示す図である。図 12 では、図 10 と同様のシステム構成である場合について説明する。この場合、検索に参加するサーバ装置の数は 3 台であるため、各サーバ装置のインデックスのルートノードのキーの数は、“ $3 - 1 = 2$ ” である。また、各サーバ装置のインデックスは、図 12 の上側のインデックス 131 a のような状態である。

【0097】

この状態で、サーバ装置 100 が故障した場合、クライアント装置 200 と接続されている検索に参加するサーバ装置の数は、サーバ装置 100 a、100 b の 2 台となる。そのため、各サーバ装置は、ルートノードのキーの数が “2” から “ $2 - 1 = 1$ ” となるようインデックスを更新する。その結果、各サーバ装置のインデックスは、例えば、図 12 の下側のインデックス 131 a のように更新される。

【0098】

また、故障したサーバ装置 100 が復旧した場合、検索に参加するサーバ装置の数は 2

10

20

30

40

50

台から3台に戻るため、各サーバ装置は、ルートノードのキーの数が“1”から“2”となるようインデックスを更新する。その結果、各サーバ装置のインデックスは、例えば、図12の下側のインデックス131aから、図12の上側のインデックス131aのように更新される。

【0099】

なお、インデックスの更新は、ノード間のキーの配置を組み替えることで実現されてもよいし、インデックスに対応するデータベース上のデータに基づいて再構築することで実現されてもよい。

【0100】

図13は、システムの機能例を示すブロック図である。クライアント装置200は、アプリケーションソフトウェア210、クライアント制御部220および稼働サーバ情報記憶部230を有する。稼働サーバ情報記憶部230は、クライアント装置200が備えるRAMやHDDに確保された記憶領域として実現できる。アプリケーションソフトウェア210およびクライアント制御部220は、クライアント装置200が備えるプロセッサが実行するプログラムのモジュールとして実現できる。

【0101】

アプリケーションソフトウェア210については、図4で説明したため、説明を省略する。稼働サーバ情報記憶部230は、各サーバ装置へ送信した検索要求に対して、応答があったサーバ装置の一覧情報が格納される応答サーバリストを記憶する。また、稼働サーバ情報記憶部230は、検索に参加するサーバ装置の数を示す情報を記憶する。

【0102】

クライアント制御部220の例について、図4で説明していない点について説明し、図4で説明した点については説明を省略する。クライアント制御部220は、処理要求部221および実行結果制御部222を有する。処理要求部221は、1または2以上のテーブルについて、データの更新や検索の要求をアプリケーションソフトウェア210から取得する。また、処理要求部221は、取得した要求に基づいて処理要求通知を生成する。処理要求通知は、データベースのデータへのアクセス処理（例えば、更新処理や検索処理）の要求を示す通知である。処理要求通知には、データベースのデータへのアクセス処理に関する情報が設定される。処理要求部221は、生成した処理要求通知を複数のサーバ装置（サーバ装置100等）に一斉送信する。

【0103】

実行結果制御部222は、処理結果通知をサーバ装置から受信する。処理結果通知は、処理要求部221が送信した処理要求通知に対する処理結果を示す通知である。受信した処理結果通知が更新処理の処理要求通知に対応するものである場合、実行結果制御部222は、1台のサーバ装置のみ（例えば、現用系のサーバ装置）から処理結果通知を受信する。そして、処理結果をアプリケーションソフトウェア210に出力する。

【0104】

また、受信した処理結果通知が検索処理の処理要求通知に対応するものである場合、実行結果制御部222は、1または2以上のサーバ装置から処理要求通知を受信する。そして、実行結果制御部222は、受信した処理要求通知に含まれる検索結果をマージし、マージされた検索結果をアプリケーションソフトウェア210に出力する。

【0105】

サーバ装置100は、サーバ制御部110、データベース120および稼働サーバリスト140を有する。サーバ制御部110およびデータベース120の例について、図4で説明していない点について説明し、図4で説明した点については説明を省略する。データベース120および稼働サーバリスト140は、サーバ装置100が備えるRAM102やHDD103に確保された記憶領域として実現できる。サーバ制御部110は、サーバ装置100が備えるプロセッサ101が実行するプログラムのモジュールとして実現できる。データベース120は、第1の実施の形態のデータベース11、11aの一例である。

。

10

20

30

40

50

【 0 1 0 6 】

データベース 1 2 0 は、インデックス情報記憶部 1 3 0 を有する。インデックス情報記憶部 1 3 0 は、データベース 1 2 0 に記憶されているテーブル（例えば、テーブル 1 2 1 等）の特定の列に対応するインデックス（例えば、インデックス 1 3 1 等）を記憶する。

【 0 1 0 7 】

稼働サーバリスト 1 4 0 は、稼働中のサーバ装置の識別情報のリストである。稼働サーバリスト 1 4 0 には、サーバ装置の識別情報が稼働優先度の高い順に記憶される。稼働優先度とは、予め定義されたサーバ装置間の優先順位であり、異なるサーバ装置には異なる稼働優先度が設定される。第 2 の実施の形態のシステムにおいて、稼働優先度が 1 位のサーバ装置が現用系のサーバ装置となる。また、稼働優先度が 2 位以下のサーバ装置が待機系のサーバ装置となる。待機系のサーバ装置が 2 以上ある場合、“待機系 1 ”，“待機系 2 ”，・・・のように、それら待機系のサーバ装置の間でも優先順位が定義されることになる。現用系のサーバ装置が故障した場合は、待機系のサーバ装置のうち稼働優先度が最も高いサーバ装置が現用系のサーバ装置となる。また、待機系 1 のサーバ装置が故障した場合は、待機系 2 のサーバ装置が待機系 1 に繰り上がる。すなわち、ある優先順位のサーバ装置が故障すると、それ以下のサーバ装置の優先順位が 1 つずつ繰り上がる。

10

【 0 1 0 8 】

サーバ制御部 1 1 0 は、処理内容判定部 1 1 1、データベース制御部 1 1 2 およびシステム管理部 1 1 3 を有する。

処理内容判定部 1 1 1 は、クライアント装置 2 0 0 から処理要求通知を受信する。受信した処理要求通知が更新処理を要求するものである場合、次の処理を実行する。

20

【 0 1 0 9 】

処理内容判定部 1 1 1 は、自己が現用系のサーバ装置であるとき、処理要求通知に含まれる更新処理をデータベース制御部 1 1 2 に要求する。処理内容判定部 1 1 1 は、自己が待機系のサーバ装置であるとき、更新処理を要求しない。

【 0 1 1 0 】

また、受信した処理要求通知が検索処理を要求するものである場合、処理内容判定部 1 1 1 は、処理要求通知に含まれる検索条件、検索に参加するサーバ装置の数、自己の稼働優先度、および、インデックス情報記憶部 1 3 0 に記憶されているインデックスのルートノードに基づき、自サーバ装置が担当する部分検索範囲を算出する。そして、算出された部分検索範囲について検索処理をデータベース制御部 1 1 2 に要求する。

30

【 0 1 1 1 】

データベース制御部 1 1 2 は、処理内容判定部 1 1 1 からの要求に応じて、データベース 1 2 0 のデータに対し、更新処理や検索処理を行う。また、データベース制御部 1 1 2 は、他のサーバ装置からテーブルの同期要求通知を受信する。同期要求通知には、同期をとるテーブルを示す情報が含まれる。データベース制御部 1 1 2 は、受信した同期要求通知に含まれるテーブルについて、同期要求通知の送信元のサーバ装置の有する同じテーブル名のテーブルと同期をとる。そして、データベース制御部 1 1 2 は、同期がとられたテーブルの特定の列に対応するインデックス（例えば、インデックス 1 3 1 ）を更新する。

40

【 0 1 1 2 】

システム管理部 1 1 3 は、各サーバ装置が正常に稼働しているか（故障していないか）確認するため、各サーバ装置に対し応答を要求する応答要求通知を送信する。システム管理部 1 1 3 は、応答要求通知に対する応答がない場合、その応答要求通知の送信先のサーバ装置を、故障したものと判断して稼働サーバリスト 1 4 0 から削除する。そして、システム管理部 1 1 3 は、削除後の検索に参加するサーバ装置の数に基づいて、インデックス情報記憶部 1 3 0 に格納されているインデックスを更新する。

【 0 1 1 3 】

また、システム管理部 1 1 3 は、故障していたサーバ装置が復旧した場合、復旧したサーバ装置を稼働サーバリスト 1 4 0 に登録する。そして、システム管理部 1 1 3 は、登録後の検索に参加するサーバ装置の数に基づいて、インデックス情報記憶部 1 3 0 に格納さ

50

れているインデックスを更新する。なお、システム管理部 113 は、サーバ装置が増設された（新しいサーバ装置が追加された）とき、故障していたサーバ装置が復旧したときと同様の処理を行ってもよい。

【0114】

次に、図 14 ～ 17 を用いて、第 2 の実施の形態によるシステムが用いるテーブルまたは通知について説明する。

図 14 は、処理要求通知の例を示す図である。処理要求通知 51 は、データベースのデータへのアクセス処理の要求を示す通知である。処理要求通知 51 は、クライアント装置 200 から、複数のサーバ装置に一斉送信される。

【0115】

処理要求通知 51 は、制御情報、並列フラグ、種別、テーブル、列および条件の項目を有する。制御情報の項目には、処理要求通知 51 に含まれる文字の文字数や文字コード等、処理要求通知 51 に対する処理を制御するための通知制御情報が設定される。

【0116】

並列フラグの項目には、検索処理を 2 以上のサーバ装置が並列に実行するか否かを示す情報が設定される。例えば、図 4 のようにサーバ装置 100, 100a, 100b が並列に検索処理を実行する場合、並列フラグの項目には“TRUE”が設定される。一方、1 つのサーバ装置（例えば、現用系のサーバ装置）のみが検索処理を実行する場合は、“FALSE”が設定される。

【0117】

種別の項目には、クライアント装置 200 が宛先のサーバ装置に要求する処理の種別を示す情報が設定される。例えば、検索処理を要求する場合、種別の項目には、“検索”が設定される。また、データの追加を要求する場合、種別の項目には、“追加”が設定される。また、データの削除を要求する場合、種別の項目には、“削除”が設定される。また、データの更新を要求する場合、種別の項目には、“更新”が設定される。

【0118】

テーブルの項目には、検索処理や更新処理の対象となるテーブルを識別する情報が設定される。列の項目には、テーブルの項目に設定されたテーブルについて、データを抽出する列（“カラム”という場合がある）または書換える列を識別するための情報が設定される。条件の項目には、検索処理や更新処理の対象とするレコードを限定するための条件が設定される。

【0119】

例えば、テーブル（T01）において“10 < C01 < 100”の検索条件を満たすレコードの列（C01）および列（C02）の値の検索をクライアント装置 200 が要求する場合、テーブルの項目には“T01”が設定される。列の項目には“C01, C02”が設定され、条件の項目には“C01 > 10 AND C01 < 100”が設定される。

【0120】

また、例えば、テーブル（T01）において列（C01）の値を“20”から“10”に更新する処理をクライアント装置 200 が要求する場合、テーブルの項目には“T01”が設定される。列の項目には“C01 = 10”が設定され、条件の項目には“C01 = 20”が設定される。

【0121】

なお、検索条件を満たすレコードについて、全ての列を取得するときは、列の項目に“*”を設定するようにしてもよい。

また、処理要求通知は、種別、テーブル、列および条件の項目の代わりに、検索処理や更新処理を示す SQL 文を含んでもよい。

【0122】

図 15 は、稼働サーバリストの例を示す図である。稼働サーバリスト 140 は、稼働中のサーバ装置の識別情報を格納するリストである。

稼働サーバリスト 140 は、サーバの項目を有する。サーバの項目には、サーバ装置を

10

20

30

40

50

識別するための識別子が設定される。サーバの項目に設定される識別子は、サーバ装置の稼働優先度が高いほど上に設定される。稼働優先度に基づいてサーバ装置の種別が判定される。例えば、サーバ装置の種別は、稼働優先度が一番高い順に“現用”、“待機1”、“待機2”と判定される。以下、識別子が“SV#A”であるサーバ装置を“サーバ装置(SV#A)”と記載する。

【0123】

例えば、稼働サーバリスト140に“SV#A”、“SV#B”、“SV#C”の順に上から登録されているとする。この場合、サーバ装置(SV#A)の稼働優先度は“1”となり、サーバ装置(SV#B)の稼働優先度は“2”となり、サーバ装置(SV#C)の稼働優先度は“3”となる。そして、サーバ装置(SV#A)の種別は“現用”となり、サーバ装置(SV#B)の種別は“待機1”となり、サーバ装置(SV#C)の種別は“待機2”となる。

10

【0124】

なお、稼働サーバリスト140は、稼働優先度を示す情報を設定する項目やサーバ装置の種別を示す情報を設定する項目を、サーバの項目の他に有してもよい。また、サーバ装置の識別子として、サーバ名を用いてもよいし、IP(Internet Protocol)アドレス等のネットワークアドレスを用いてもよい。

【0125】

図16は、処理結果通知の例を示す図である。処理結果通知52は、クライアント装置200が送信した処理要求通知51に対する処理結果を示す通知である。

20

処理結果通知52は、制御情報、直接指定フラグおよび処理結果の項目を有する。

【0126】

制御情報の項目には、検索に参加するサーバ装置の数や、処理結果通知52の送信元のサーバ装置の稼働優先度や、他の制御情報(例えば、通知の文字数や文字コード等)が設定される。サーバ装置の数や稼働優先度は、処理結果通知52を送信したサーバ装置が送信時点で認識しているものである。例えば、検索に参加するサーバ装置の数が“3”であり、処理結果通知52の送信元のサーバ装置の稼働優先度が“1”であり、他の制御情報が“情報A”である場合、制御情報の項目には、“3:1:情報A”が設定される。

【0127】

直接指定フラグの項目には、インデックス検索を行う場合において、直接指定により検索されたか否かを示す情報が設定される。例えば、データベースのデータが直接指定により検索された場合、“TRUE”が設定される。また、範囲指定により検索された場合、直接指定フラグの項目には“FALSE”が設定される。直接指定とは、“C01=10”等のように、検索条件において、列の値の範囲ではなく列の特定の値を直接指定することである。一方、“ $1 < C01 < 10$ ”のように、検索条件において列の値の範囲を指定することを範囲指定という。

30

【0128】

処理結果の項目には、処理結果通知52の送信元のサーバ装置による処理結果が設定される。処理結果通知52の送信元のサーバ装置が検索処理を実行した場合、処理結果の項目には、検索された列のデータが設定される。例えば、列(C01)に“20”が設定され、列(C02)に“aa”が設定されたレコードと、列(C01)に“25”が設定され、列(C02)に“bb”が設定されたレコードが検索されたとき、検索結果が検索結果の項目には、“(20, aa), (25, bb)”が設定される。データが検索されなかった場合は、処理結果の項目は空になるか、データが検索されなかったことを示す情報(例えば、“データ無し”)を含む。また、処理結果が異常である場合、処理結果の項目には、処理結果が異常であることを示す情報(例えば、“異常終了”)が設定される。

40

【0129】

また、処理結果通知52の送信元のサーバ装置が更新処理を実行した場合、処理結果の項目には、更新処理の成否を示す情報が設定される。

図17は、稼働サーバ数と応答サーバリストの例を示す図である。稼働サーバ数231

50

は、検索に参加するサーバ装置の数を示す情報である。応答サーバリスト 2 3 2 は、処理結果通知 5 2 の送信元のサーバ装置を示すリストである。稼働サーバ数 2 3 1 および応答サーバリスト 2 3 2 は、稼働サーバ情報記憶部 2 3 0 に記憶される。

【 0 1 3 0 】

稼働サーバ数 2 3 1 には、1 個目の処理結果通知 5 2 が示す検索に参加するサーバ装置の数が登録される。クライアント装置 2 0 0 は、2 個目以降の処理結果通知 5 2 を受信したとき、2 個目以降の処理結果通知 5 2 の制御情報の項目に含まれる検索に参加するサーバ装置の数を示す情報と、稼働サーバ数 2 3 1 とを比較する。比較結果が一致しない場合、クライアント装置 2 0 0 は、同じ処理要求通知 5 1 に対する全サーバ装置からの処理結果通知 5 2 が揃う前に、サーバ装置の故障や復旧により検索に参加するサーバ装置の数が
10 変化したと判定する。クライアント装置 2 0 0 は、サーバ装置の数が変化したと判定したとき、受信済みの処理結果を破棄し、処理要求通知 5 1 を再度複数のサーバ装置に一斉送信する。

【 0 1 3 1 】

また、クライアント装置 2 0 0 は、検索の処理要求通知 5 1 に対する処理結果通知 5 2 を受信したとき、受信した処理結果通知 5 2 の送信元のサーバ装置の稼働優先度を、応答サーバリスト 2 3 2 に登録する。サーバ装置の稼働優先度は、処理結果通知 5 2 の制御情報の項目に含まれる。

【 0 1 3 2 】

次に、図 1 8 ~ 2 3 を用いて、データベースのデータの検索処理について説明する。
20

図 1 8 は、クライアント装置による検索制御の例を示すフローチャートである。以下、図 1 8 に示す処理をステップ番号に沿って説明する。

【 0 1 3 3 】

(ステップ S 2 1) 処理要求部 2 2 1 は、アプリケーションソフトウェア 2 1 0 から検索を要求される。検索の要求には、検索対象となるテーブルを示す情報、検索対象となるテーブルの列を示す情報および検索条件を示す情報等が含まれる。アプリケーションソフトウェア 2 1 0 が出力する上記の情報は、SQL 文として記述されていてもよい。

【 0 1 3 4 】

(ステップ S 2 2) 処理要求部 2 2 1 は、アプリケーションソフトウェア 2 1 0 からの検索の要求に基づいて、次のように処理要求通知 5 1 を生成する。
30

処理要求部 2 2 1 は、処理要求通知 5 1 の制御情報の項目に通知制御情報を設定する。また、処理要求部 2 2 1 は、並列フラグの項目に、複数のサーバ装置に並列して検索処理をさせるかを示す情報を設定する。並列して検索処理をさせるかを示す情報は、稼働サーバ情報記憶部 2 3 0 に予め記憶されていてもよいし、アプリケーションソフトウェア 2 1 0 から指定されてもよい。また、処理要求部 2 2 1 は、種別の項目に“検索”を設定する。また、処理要求部 2 2 1 は、テーブルの項目に、検索の要求に含まれるテーブルを示す情報を設定する。また、処理要求部 2 2 1 は、列の項目に、検索の要求に含まれるテーブルの列を示す情報を設定する。また、処理要求部 2 2 1 は、条件の項目に、検索の要求に含まれる検索条件を示す情報を設定する。

【 0 1 3 5 】

そして、処理要求部 2 2 1 は、生成した処理要求通知 5 1 を複数のサーバ装置（例えば、サーバ装置 1 0 0 , 1 0 0 a , 1 0 0 b）に一斉送信する。
40

(ステップ S 2 3) 実行結果制御部 2 2 2 は、現在時刻を処理開始時刻として一時的に格納する。処理開始時刻の格納場所には、例えば、稼働サーバ情報記憶部 2 3 0 等の記憶領域を用いる。

【 0 1 3 6 】

(ステップ S 2 4) 実行結果制御部 2 2 2 は、ステップ S 2 2 で送信した処理要求通知 5 1 に対する 1 つの処理結果通知 5 2 を、1 つのサーバ装置から受信したか判定する。処理結果通知 5 2 を受信した場合は、処理をステップ S 2 5 へ進める。処理結果通知 5 2 を受信していない場合は、処理をステップ S 3 2 へ進める。
50

【 0 1 3 7 】

(ステップ S 2 5) 実行結果制御部 2 2 2 は、処理結果通知 5 2 が異常終了を示すか判定する。異常終了を示すか否かは、例えば、処理結果通知 5 2 の処理結果の項目に“異常終了”が設定されているかで判断する。処理結果通知 5 2 が異常終了を示す場合、処理をステップ S 3 6 へ進める。処理結果通知 5 2 が異常終了を示さない場合、処理をステップ S 2 6 へ進める。

【 0 1 3 8 】

(ステップ S 2 6) 実行結果制御部 2 2 2 は、処理結果通知 5 2 の処理結果の項目に設定されている検索されたデータをクライアント装置 2 0 0 の備える R A M 等の記憶領域に一時的に格納する。

10

【 0 1 3 9 】

(ステップ S 2 7) 実行結果制御部 2 2 2 は、受信した処理結果通知 5 2 が処理要求通知 5 1 に対する 1 つ目の応答である場合、その処理結果通知 5 2 が示すサーバ装置の数を稼働サーバ数 2 3 1 として登録する。稼働サーバ数 2 3 1 は、1 つ目の処理結果通知 5 2 を送信したサーバ装置が送信時点で認識していた、検索に参加するサーバ装置の数を意味する。また、実行結果制御部 2 2 2 は、処理結果通知 5 2 の送信元のサーバ装置の稼働優先度を示す情報を応答サーバリスト 2 3 2 に登録する。送信元のサーバ装置の稼働優先度を示す情報は、処理結果通知 5 2 の制御情報の項目に含まれる。

【 0 1 4 0 】

(ステップ S 2 8) 実行結果制御部 2 2 2 は、処理結果通知 5 2 の直接指定フラグの項目が“T R U E”であるか判定する。直接指定フラグが“T R U E”である場合、処理をステップ S 3 9 へ進める。直接指定フラグが“F A L S E”である場合、処理をステップ S 2 9 へ進める。

20

【 0 1 4 1 】

(ステップ S 2 9) 実行結果制御部 2 2 2 は、ステップ S 2 2 で一斉送信した処理要求通知 5 1 において、並列フラグの項目が“T R U E”であるか判定する。並列フラグが“T R U E”である場合、処理をステップ S 3 0 へ進める。並列フラグが“F A L S E”である場合、処理をステップ S 3 9 へ進める。

【 0 1 4 2 】

(ステップ S 3 0) 実行結果制御部 2 2 2 は、検索に参加する全サーバ装置から処理結果通知 5 2 を受信済みであるか判定する。全サーバ装置から処理結果通知 5 2 を受信したかは、例えば、応答サーバリスト 2 3 2 に登録されているサーバ装置の数が、稼働サーバ数 2 3 1 と一致するかにより判断する。全サーバ装置から処理結果通知 5 2 を受信した場合、処理をステップ S 3 8 へ進める。少なくとも 1 つのサーバ装置から処理結果通知 5 2 が受信されていない場合、処理をステップ S 3 1 へ進める。

30

【 0 1 4 3 】

(ステップ S 3 1) 実行結果制御部 2 2 2 は、ステップ S 2 3 で設定した処理開始時刻と現在時刻との差(処理要求通知 5 1 を送信してからの経過時間)が閾値未満か判定する。閾値は、例えば、クライアント装置 2 0 0 の備える R A M 等の記憶領域に格納されている。時刻差が閾値未満である場合、処理をステップ S 3 2 へ進める。時刻差が閾値以上である場合、処理をステップ S 3 7 へ進める。

40

【 0 1 4 4 】

(ステップ S 3 2) 実行結果制御部 2 2 2 は、処理結果通知 5 2 の受信待ちの状態であるか判定する。処理結果通知 5 2 の受信待ちの状態であるかは、ステップ S 2 4 で処理結果通知 5 2 を受信したかにより判断する。処理結果通知 5 2 の受信待ちの状態である(ステップ S 2 4 で処理結果通知 5 2 を受信しなかった)場合、処理をステップ S 3 4 へ進める。処理結果通知 5 2 の受信待ちの状態でない(ステップ S 2 4 で処理結果通知 5 2 を受信した)場合、処理をステップ S 3 3 へ進める。

【 0 1 4 5 】

(ステップ S 3 3) 実行結果制御部 2 2 2 は、処理要求通知 5 1 に応じて行われる並列

50

検索の途中で、検索に参加するサーバ装置の数が増減したか判定する。

検索に参加するサーバ装置の数は、サーバ装置の故障、サーバ装置の復旧、サーバ装置の追加等によって変化する。実行結果制御部 222 は、今回受信した処理結果通知 52 が示すサーバ装置の数が、稼働サーバ数 231 と一致しないとき、検索に参加するサーバ装置の数が増減したと判定する。これは、今回受信した処理結果通知 52 の送信時点で認識されていたサーバ装置の数が、1 つ目の処理結果通知 52 の送信時点で認識されていたサーバ装置の数から増減したことを示している。ただし、実行結果制御部 222 は、復旧したサーバ装置を示す情報を含むコマンドを、第 2 の実施の形態のシステムの管理者から受信したとき、サーバ装置の数が増減したと判定してもよい。

【0146】

検索に参加するサーバ装置の数が増減した場合、処理をステップ S35 へ進める。検索に参加するサーバ装置の数が増減していない場合、処理をステップ S34 へ進める。

(ステップ S34) 実行結果制御部 222 は、一定時間(例えば、10 ミリ秒または 100 ミリ秒)経過するのを待つ。一定時間は、第 2 の実施の形態のシステムのユーザに設定されてもよいし、クライアント装置 200 の備える HDD 等の記憶領域に予め記憶されていてもよい。そして、処理をステップ S24 へ進め、次の処理結果通知 52 を受信したか判定する。

【0147】

(ステップ S35) 実行結果制御部 222 は、ステップ S26 で格納した受信済みのデータをクリアする。また、実行結果制御部 222 は、稼働サーバ数 231 および応答サーバリスト 232 に登録されている情報を全てクリアする。そして、実行結果制御部 222 は、処理をステップ S22 へ進める。これにより、処理要求通知 51 が再送される。なお、実行結果制御部 222 は、処理要求通知 51 の再送後に古い処理要求通知 51 に対応する処理結果通知 52 を受信した場合は、その処理結果通知 52 を無視する。

【0148】

図 19 は、クライアント装置による検索制御の例を示すフローチャート(続き)である。以下、図 19 に示す処理をステップ番号に沿って説明する。

(ステップ S36) 実行結果制御部 222 は、要求された検索が異常終了した旨をアプリケーションソフトウェア 210 に通知する。なお、通知後に処理要求通知 51 に対応する処理結果通知 52 (未受信分の処理結果通知)を受信した場合は、その処理結果通知 52 を無視する。また、実行結果制御部 222 は、ステップ S26 で格納した受信済みのデータをクリアする。また、実行結果制御部 222 は、稼働サーバ数 231 および応答サーバリスト 232 に登録されている情報を全てクリアする。アプリケーションソフトウェア 210 は、異常終了に対応する処理を実行する。そして、クライアントの検索制御を終了する。

【0149】

(ステップ S37) 実行結果制御部 222 は、要求された検索がタイムアウトした旨をアプリケーションソフトウェア 210 に通知する。なお、通知後に処理要求通知 51 に対応する処理結果通知 52 (未受信分の処理結果通知)を受信した場合は、その処理結果通知 52 を無視する。また、実行結果制御部 222 は、ステップ S26 で格納した受信済みのデータをクリアする。また、実行結果制御部 222 は、稼働サーバ数 231 および応答サーバリスト 232 に登録されている情報を全てクリアする。アプリケーションソフトウェア 210 は、タイムアウトに対応する処理を実行する。そして、クライアントの検索制御を終了する。

【0150】

(ステップ S38) 実行結果制御部 222 は、ステップ S26 で格納した複数のサーバ装置からのデータをマージする。例えば、実行結果制御部 222 は、複数のサーバ装置から受信したレコードのリストを連結する。

【0151】

(ステップ S39) 実行結果制御部 222 は、マージしたデータをアプリケーションソ

10

20

30

40

50

フトウェア 210 に通知する。また、実行結果制御部 222 は、稼働サーバ数 231 および応答サーバリスト 232 に登録されている情報を全てクリアする。アプリケーションソフトウェア 210 は、通知されたデータを用いて処理を実行する。

【0152】

図 20 は、サーバ装置による検索制御の例を示すフローチャートである。図 20 ~ 23 で説明する処理は、サーバ装置 100 が実行しているものとする。また、図 20 ~ 23 では、直接指定による検索条件で非インデックス検索の要求の場合、処理要求通知 51 の並列フラグの項目には、“FALSE” が設定されているものとする。以下、図 20 に示す処理をステップ番号に沿って説明する。

【0153】

(ステップ S41) 処理内容判定部 111 は、クライアント装置 200 から処理要求通知 51 を受信する。処理内容判定部 111 は、処理要求通知 51 の種別の項目が“検索”であることを確認する。

【0154】

(ステップ S42) 処理内容判定部 111 は、受信した処理要求通知 51 の条件の項目を参照し、クライアント装置 200 から要求された検索条件の示す検索範囲を確認する。

(ステップ S43) 処理内容判定部 111 は、現在時刻を処理開始時刻に設定する。処理開始時刻は、例えば、サーバ装置 100 が備える RAM 102 に記憶される。

【0155】

(ステップ S44) 処理内容判定部 111 は、検索に参加するサーバ装置の数(以下、検索サーバ装置数 a)を取得する。検索に参加するサーバ装置の数は、例えば、稼働サーバリスト 140 に登録されているサーバ装置の数をカウントすることで取得する。

【0156】

(ステップ S45) 処理内容判定部 111 は、サーバ装置 100 の稼働優先度を取得する。サーバ装置 100 の稼働優先度(以下、稼働優先度 b)は、例えば、図 15 のように、稼働サーバリスト 140 に記憶されているサーバ装置 100 の順位により取得できる。

【0157】

(ステップ S46) 処理内容判定部 111 は、検索サーバ装置数 a および稼働優先度 b に基づいて、処理結果通知 52 に設定する制御情報を生成する。ただし、制御情報の生成は、処理結果通知 52 を送信する直前に行ってもよい。

【0158】

(ステップ S47) 処理内容判定部 111 は、受信した処理要求通知 51 の並列フラグの項目が“TRUE”であるか判定する。並列フラグが“TRUE”である場合、処理をステップ S49 へ進める。並列フラグが“FALSE”である場合、処理をステップ S48 へ進める。

【0159】

(ステップ S48) 処理内容判定部 111 は、稼働優先度 b が“1”であるか判定する。稼働優先度 b が“1”である場合(サーバ装置 100 が現用系である場合)、サーバ装置 100 のみで検索を行うと判定し、処理をステップ S52 へ進める。稼働優先度 b が“1”以外である場合(サーバ装置 100 が待機系である場合)、他の 1 つのサーバ装置のみで検索を行うと判定し、処理をステップ S65 へ進める。

【0160】

(ステップ S49) 処理内容判定部 111 は、インデックス情報記憶部 130 に記憶されているインデックス、ステップ S42 で確認した検索範囲、検索サーバ装置数 a、および稼働優先度 b に基づいて、クライアント装置 200 から検索条件として指定された列の値について限定範囲を算出する。詳細は、図 22 で説明する。

【0161】

(ステップ S50) 処理内容判定部 111 は、ステップ S42 で確認した検索範囲と、算出された限定範囲の重複する範囲が存在するか判定する。重複する範囲が存在する場合、処理をステップ S51 へ進める。重複する範囲が存在しない場合、処理をステップ S6

10

20

30

40

50

5へ進める。

【0162】

(ステップS51) 処理内容判定部111は、図7～8で説明したように、ステップS42で確認した検索範囲と、限定範囲とが重複する範囲を部分検索範囲として算出する。

なお、インデックス検索と非インデックス検索のうち、非インデックス検索においては、図23で説明するように、実質的に、処理要求通知51で指定される検索範囲を分割することで限定範囲を算出している。この場合、算出される限定範囲は指定される検索範囲と明らかに重複しているため、処理内容判定部111は、ステップS50の処理を実行しなくてもよい。また、サーバ装置100の稼働優先度が1位でなく末尾でもない場合、図23の処理に従って算出される限定範囲は、そのまま部分検索範囲とみなすことができる。その場合、処理内容判定部111は、ステップS51の処理を実行しなくてもよい。

10

【0163】

(ステップS52) データベース制御部112は、並列検索が行われる場合(ステップS47のYES)、ステップS51で算出した部分検索範囲に限定して、データベース120からデータを検索する。また、データベース制御部112は、サーバ装置100のみが検索を行う場合(ステップS48のYES)、クライアント装置200から受信された処理要求通知51に従って、データベース120からデータを検索する。検索するテーブルは、処理要求通知51のテーブルの項目を確認する。データを抽出する列は、処理要求通知51の列の項目を確認する。そして、処理をステップS61へ進める。

【0164】

20

図21は、サーバ装置による検索制御の例を示すフローチャート(続き)である。以下、図21に示す処理をステップ番号に沿って説明する。

(ステップS61) システム管理部113は、検索に参加するサーバ装置の数が増減したか判定する。検索に参加するサーバ装置の数は、サーバ装置の故障、サーバ装置の復旧、サーバ装置の追加等により変化する。

【0165】

システム管理部113は、サーバ装置の故障を次のように判定する。システム管理部113は、定期的に各サーバ装置に対し応答要求通知を送信する。送信先のサーバ装置のアドレスは、例えば、サーバ装置100の備えるHDD103等の記憶領域に格納されている。システム管理部113は、応答要求通知に対する応答がない場合、送信先のサーバ装置が故障したと判定する。

30

【0166】

また、システム管理部113は、復旧したサーバ装置を示す情報を含むコマンドを、第2の実施の形態のシステムの管理者から受信したとき、サーバ装置が復旧したと判定する。コマンドには、復旧したサーバ装置の識別子が含まれる。ただし、システム管理部113は、故障したサーバ装置にも応答要求通知を定期的に送信し、応答が得られたときに当該サーバ装置が復旧したと判定してもよい。

【0167】

検索に参加するサーバ装置の数が増減していない場合、処理をステップS65へ進める。検索に参加するサーバ装置の数が増減した場合、処理をステップS62へ進める。

40

(ステップS62) システム管理部113は、稼働サーバリスト140を更新する。例えば、ステップS61でサーバ装置の故障を検出したとき、システム管理部113は、故障と判定されたサーバ装置を稼働サーバリスト140から削除する。また、例えば、ステップS61でコマンドを受け付けたとき、システム管理部113は、当該コマンドから復旧したサーバ装置の識別子を確認し、確認した識別子を稼働サーバリスト140の末尾に登録する。

【0168】

(ステップS63) データベース制御部112は、図20のステップS43で設定した処理開始時刻と、現在時刻との差(処理要求通知51が受信されてからの経過時間)が閾値より大きいか判定する。時刻差が閾値より大きい場合、処理をステップS64へ進める

50

。時刻差が閾値以下である場合、処理をステップS 4 4へ進める。なお、このサーバ側の閾値は、図1 8のステップS 3 1のクライアント側の閾値とは異なってもよい。このサーバ側の閾値は、ステップS 3 1のクライアント側の閾値より小さい（クライアント装置2 0 0よりもサーバ装置1 0 0の方が先にタイムアウトする）ことが望ましい。

【0 1 6 9】

（ステップS 6 4）データベース制御部1 1 2は、タイムアウト通知をクライアント装置2 0 0に送信する。そして、サーバの検索制御を終了する。

（ステップS 6 5）データベース制御部1 1 2は、要求された検索が直接指定による検索処理か判定する。直接指定による検索処理である場合、処理をステップS 6 6へ進める。範囲指定による検索処理である場合、処理をステップS 6 8へ進める。

10

【0 1 7 0】

（ステップS 6 6）データベース制御部1 1 2は、図2 0のステップS 5 2で検索処理を実行したか判定する。検索処理を実行した場合、処理をステップS 6 7へ進める。検索処理を実行していない場合、クライアント装置2 0 0へ応答せずに（処理結果通知5 2を送信せずに）、サーバの検索制御を終了する。

【0 1 7 1】

（ステップS 6 7）データベース制御部1 1 2は、処理結果通知5 2の直接指定フラグの項目に“TRUE”を設定する。そして、処理をステップS 6 9へ進める。

（ステップS 6 8）データベース制御部1 1 2は、処理結果通知5 2の直接指定フラグの項目に“FALSE”を設定する。

20

【0 1 7 2】

（ステップS 6 9）データベース制御部1 1 2は、処理結果通知5 2の制御情報の項目に、図2 0のステップS 4 6で生成した制御情報を設定する。また、データベース制御部1 1 2は、処理結果通知5 2の処理結果の項目に、図2 0のステップS 5 2の検索処理により検索されたデータの集合を設定する。検索条件または部分検索条件に該当するレコードが1件もなかった場合、このデータの集合は空集合になる。その場合、処理結果通知5 2の処理結果の項目は、空になるか、データがない旨の情報を含む。そして、データベース制御部1 1 2は、処理結果通知5 2をクライアント装置2 0 0に送信する。

【0 1 7 3】

図2 2は、限定範囲の算出処理の例を示すフローチャートである。図2 2の処理は、ステップS 4 9で実行される。以下、図2 2に示す処理をステップ番号に沿って説明する。

30

（ステップS 8 1）処理内容判定部1 1 1は、検索処理がインデックス検索か判定する。すなわち、処理要求通知5 1において検索範囲が指定されている列について、インデックス1 3 1が作成されているか判定する。検索処理がインデックス検索である場合、処理をステップS 8 3へ進める。検索処理が非インデックス検索である場合、処理をステップS 8 2へ進める。

【0 1 7 4】

（ステップS 8 2）処理内容判定部1 1 1は、ステップS 4 2で確認した検索範囲、検索サーバ装置数aおよび稼働優先度bに基づいて、列の値の限定範囲をインデックス1 3 1を用いずに算出する。すなわち、処理内容判定部1 1 1は、検索サーバ装置数aに応じて検索範囲を分割する。詳細は、図2 3で説明する。

40

【0 1 7 5】

（ステップS 8 3）処理内容判定部1 1 1は、稼働優先度bが“1”であるか判定する。稼働優先度bが“1”である場合（サーバ装置1 0 0が現用系である場合）、処理をステップS 8 4へ進める。稼働優先度bが“1”以外である場合（サーバ装置1 0 0が待機系である場合）、処理をステップS 8 5へ進める。

【0 1 7 6】

（ステップS 8 4）処理内容判定部1 1 1は、限定範囲を“ $p < R(1)$ ”と算出する。pは、処理要求通知5 1の検索条件に記載された列の値を示す変数であり、以下同様とする。“ $R(q)$ ”は、ルートノードに含まれるキーのうち、左から（前方から）q番目

50

のキーを示し、以下同様とする。例えば、1番目のルートノードのキーは、“R(1)”と示される。そして、限定範囲算出を終了する。

【0177】

(ステップS85) 処理内容判定部111は、稼働優先度bの値が最大か判定する。具体的には、処理内容判定部111は、稼働優先度bが検索サーバ装置数aと一致するか判定する。稼働優先度の値が最大である場合(サーバ装置100の稼働優先度が検索に参加するサーバ装置の中で最も低い場合)、処理をステップS86へ進める。稼働優先度の値が最大でない場合、処理をステップS87へ進める。

【0178】

(ステップS86) 処理内容判定部111は、限定範囲を“R(稼働優先度b-1)p”と算出する。そして、限定範囲算出を終了する。

10

(ステップS87) 処理内容判定部111は、限定範囲を“R(稼働優先度b-1)p < R(稼働優先度b)”と算出する。

【0179】

図23は、非インデックス検索の場合の限定範囲の算出処理の例を示すフローチャートである。図23の処理は、前述のステップS82で実行される。以下、図23に示す処理をステップ番号に沿って説明する。

【0180】

(ステップS91) 処理内容判定部111は、図20のステップS42で確認した検索範囲において、上限値maxと下限値minとの差xを算出する。

20

(ステップS92) 処理内容判定部111は、差xを検索サーバ装置数aで割った商yを算出する。

【0181】

(ステップS93) 処理内容判定部111は、稼働優先度bが“1”であるか判定する。稼働優先度bが“1”である場合(サーバ装置100が現用系である場合)、処理をステップS94へ進める。稼働優先度bが“1”以外である場合(サーバ装置100が待機系である場合)、処理をステップS95へ進める。

【0182】

(ステップS94) 処理内容判定部111は、限定範囲を“p < min + 稼働優先度b * y”と算出する。そして、限定範囲算出を終了する。

30

(ステップS95) 処理内容判定部111は、稼働優先度bの値が最大か判定する。稼働優先度bの値が最大である場合(サーバ装置100の稼働優先度が検索に参加するサーバ装置の中で最も低い場合)、処理をステップS96へ進める。稼働優先度bの値が最大でない場合、処理をステップS97へ進める。

【0183】

(ステップS96) 処理内容判定部111は、限定範囲を“min + (稼働優先度b-1) * y p”と算出する。そして、限定範囲算出を終了する。

(ステップS97) 処理内容判定部111は、限定範囲を“min + (稼働優先度b-1) * y p < min + 稼働優先度b * y”と算出する。

40

【0184】

次に、図24~25を用いて、データベースのデータの更新処理について説明する。

図24は、クライアント装置による更新制御の例を示すフローチャートである。以下、図24に示す処理をステップ番号に沿って説明する。

【0185】

(ステップS101) 処理要求部221は、アプリケーションソフトウェア210から更新を要求される。更新の要求には、更新対象となるテーブルを示す情報、更新対象となる列およびその値を示す情報、更新処理の種別を示す情報(追加、更新または削除)および更新対象のレコードの条件が含まれる。アプリケーションソフトウェア210が出力する上記の情報は、SQL文として記述されていてもよい。

【0186】

50

(ステップS102) 処理要求部221は、アプリケーションソフトウェア210からの更新の要求に基づいて、処理要求通知51を生成する。処理要求通知51は、処理要求部221により次のように生成される。

【0187】

処理要求部221は、処理要求通知51の制御情報の項目に、通知制御情報を設定する。また、処理要求部221は、種別の項目に、要求された更新処理の種別を設定する。また、処理要求部221は、テーブルの項目に、更新対象のテーブルを示す情報を設定する。また、処理要求部221は、列の項目に、更新対象の列およびその値を示す情報を設定する。また、処理要求部221は、条件の項目に、更新対象のレコードの条件を示す情報を設定する。

10

【0188】

そして、処理要求部221は、生成した処理要求通知51を複数のサーバ装置(例えば、サーバ装置100, 100a, 100b)に一斉送信する。

(ステップS103) 実行結果制御部222は、現在時刻を処理開始時刻として一時的に格納する。

【0189】

(ステップS104) 実行結果制御部222は、複数のサーバ装置の何れか1つから処理結果通知52を受信したか判定する。何れか1つのサーバ装置から処理結果通知52を受信した場合、処理をステップS108へ進める。何れのサーバ装置からも処理結果通知52を受信していない場合、処理をステップS105へ進める。

20

【0190】

(ステップS105) 実行結果制御部222は、ステップS103で設定した処理開始時刻と、現在時刻との差(処理要求通知51を送信してからの経過時間)が閾値未満か判定する。時刻差が閾値未満である場合、処理をステップS106へ進め、処理結果通知52を待つ。時刻差が閾値以上である場合、処理をステップS107へ進める。

【0191】

(ステップS106) 実行結果制御部222は、一定時間(例えば、10ミリ秒または100ミリ秒)経過するのを待つ。一定時間は、第2の実施の形態のシステムのユーザに設定されてもよいし、クライアント装置200の備えるHDD等の記憶領域に予め記憶されていてもよい。そして、処理をステップS104へ進め、処理結果通知52を受信するのを待つ。

30

【0192】

(ステップS107) 実行結果制御部222は、要求された更新がタイムアウトにより終了した旨をアプリケーションソフトウェア210に通知する。アプリケーションソフトウェア210は、タイムアウトに対応する処理を実行する。そして、クライアントの更新制御を終了する。

【0193】

(ステップS108) 実行結果制御部222は、処理結果通知52の処理結果の項目を参照し、ステップS102で送信した処理要求通知51に対応する更新処理が成功したか確認する。実行結果制御部222は、確認した更新処理の結果をアプリケーションソフトウェア210に通知する。

40

【0194】

図25は、サーバ装置による更新制御の例を示すフローチャートである。図25で説明する処理は、サーバ装置100が実行しているものとする。以下、図25に示す処理をステップ番号に沿って説明する。

【0195】

(ステップS111) 処理内容判定部111は、クライアント装置200から処理要求通知51を受信する。処理内容判定部111は、処理要求通知51の種別の項目が更新処理の種別(追加、更新または削除)を示す情報であることを確認する。

【0196】

50

(ステップS 1 1 2) 処理内容判定部 1 1 1 は、検索サーバ装置数 a を取得する。

(ステップS 1 1 3) 処理内容判定部 1 1 1 は、稼働優先度 b を取得する。

(ステップS 1 1 4) 処理内容判定部 1 1 1 は、検索サーバ装置数 a および稼働優先度 b に基づいて、処理結果通知 5 2 に設定する制御情報を生成する。ただし、制御情報の生成は、処理結果通知 5 2 を送信する直前に行ってもよい。

【0 1 9 7】

(ステップS 1 1 5) 処理内容判定部 1 1 1 は、稼働優先度 b が “ 1 ” であるか判定する。稼働優先度が “ 1 ” である場合 (サーバ装置 1 0 0 が現用系である場合)、処理をステップS 1 1 6 へ進める。稼働優先度が “ 1 ” 以外である場合 (サーバ装置 1 0 0 が待機系である場合)、処理をステップS 1 2 0 へ進める。

【0 1 9 8】

(ステップS 1 1 6) データベース制御部 1 1 2 は、ステップS 1 1 1 で受信した処理要求通知 5 1 に基づいて、更新処理を実行する。更新するテーブルは、処理要求通知 5 1 のテーブルの項目を確認する。更新するレコードは、処理要求通知 5 1 の条件の項目に基づいて検索する。

【0 1 9 9】

(ステップS 1 1 7) データベース制御部 1 1 2 は、更新後のテーブルについて、同期を要求するための同期要求通知を他のサーバ装置それぞれへ送信する。例えば、同期要求通知には、データベース 1 2 0 の更新履歴が含まれる。更新履歴は、データベース制御部 1 1 2 がデータベース 1 2 0 に対して行った操作を示すコマンド (SQL 文でもよい) を含んでもよいし、更新後のデータを含んでもよい。

【0 2 0 0】

(ステップS 1 1 8) データベース制御部 1 1 2 は、更新したテーブルの特定の列に対応するインデックスがインデックス情報記憶部 1 3 0 に記憶されている場合、そのインデックスを更新する。ただし、テーブルの更新内容によっては、インデックスを更新しなくてもよい場合がある。

【0 2 0 1】

(ステップS 1 1 9) データベース制御部 1 1 2 は、処理結果通知 5 2 の制御情報の項目に、ステップS 1 1 4 で生成した制御情報を設定する。また、データベース制御部 1 1 2 は、処理結果通知 5 2 の処理結果の項目に、ステップS 1 1 6 の更新処理の結果を示す情報 (例えば、更新成功を示す情報) を設定する。

【0 2 0 2】

データベース制御部 1 1 2 は、処理結果通知 5 2 をクライアント装置 2 0 0 に送信する。そして、サーバの更新制御を終了する。

(ステップS 1 2 0) データベース制御部 1 1 2 は、受信した処理要求通知 5 1 に応じた更新処理は行わず、現用系のサーバ装置に更新処理を任せる。そして、データベース制御部 1 1 2 は、同期要求通知を現用系のサーバ装置から受信する。

【0 2 0 3】

(ステップS 1 2 1) データベース制御部 1 1 2 は、同期要求通知が示すデータについて、データベース 1 2 0 を現用系のサーバ装置が有するデータベースと同期させる。例えば、データベース制御部 1 1 2 は、受信した同期要求通知が示す更新履歴に基づいて、現用系のサーバ装置が行ったものと同様の操作をデータベース 1 2 0 に対して行い、現用系のデータベースの状態を再現する。

【0 2 0 4】

(ステップS 1 2 2) データベース制御部 1 1 2 は、更新したテーブルの特定の列に対応するインデックスがインデックス情報記憶部 1 3 0 に記憶されている場合、そのインデックスを更新する。ただし、テーブルの更新内容によっては、インデックスを更新しなくてもよい場合がある。

【0 2 0 5】

次に、図 2 6 ~ 2 8 を用いて、第 2 の実施の形態のシステムにより検索されるデータ例

10

20

30

40

50

について説明する。

図26は、非インデックス検索により検索されるデータの例を示す図である。クライアント装置200は、サーバ装置100およびサーバ装置100aと接続している。サーバ装置100はデータベース120を有し、サーバ装置100aはデータベース120aを有する。データベース120、120aそれぞれは、テーブル(T01)を有する。テーブル(T01)は、列(C01)および列(C02)を有する。

【0206】

以下、データベース120、120aに格納されたレコードを示す情報を“列(C01)の値、列(C02)の値”のように記載する。例えば、データベース120、120aにおいて、列(C01)に“20”が設定され、列(C02)に“bb”が設定されているレコードは、“20, bb”と記載される。

10

【0207】

データベース120、120aは、“3, aa”、“20, bb”、“25, cc”、“75, dd”および“200, ee”を含むテーブル(T01)を有する。ただし、データベース120、120aは、列(C01)に対応するインデックスを有していない。

【0208】

この状態で、クライアント装置200は、テーブル(T01)から“10 < C01 < 100”を満たすレコードを検索するよう、サーバ装置100、100aに要求するものとする。すると、サーバ装置100、100aは、テーブル(T01)に対し、非インデックス検索により次のようにデータを検索する。

20

【0209】

サーバ制御部110は、クライアント装置200から指定された検索範囲、検索に参加するサーバ装置の数(=2台)およびサーバ装置100の稼働優先度(=“1”)に基づいて、“10 < C01 < 55”を部分検索範囲#1として算出する。そして、サーバ制御部110は、部分検索範囲#1に限定してデータベース120からデータを検索する。その結果、“20, bb”および“25, cc”が抽出され、サーバ装置100の検索結果114としてクライアント装置200に送信される。

【0210】

サーバ制御部110aは、クライアント装置200から指定された検索範囲、検索に参加するサーバ装置の数(=2台)およびサーバ装置100aの稼働優先度(=“2”)に基づいて、“55 < C01 < 100”を部分検索範囲#2として算出する。そして、サーバ制御部110aは、部分検索範囲#2に限定してデータベース120aからデータを検索する。その結果、“75, dd”が抽出され、サーバ装置100aの検索結果114aとしてクライアント装置200に送信される。

30

【0211】

クライアント装置200は、サーバ装置100の検索結果114(“20, bb”および“25, cc”)と、サーバ装置100aの検索結果114a(“75, dd”)とをマージする。ここでは、検索結果114のレコードと検索結果114aのレコードとを連結して1つのリストにする。マージされた検索結果211は、検索要求に対する応答としてアプリケーションソフトウェア210に提供される。

40

【0212】

図27は、インデックス検索により検索されるデータの例を示す図である。図27~28では、図26で説明した内容と同様の点について、説明を省略する。データベース120は、列(C01)に対応するインデックス131を有し、データベース120aは、列(C01)に対応するインデックス131aを有する。インデックス131、131aそれぞれのルートノードのキーには“24”が設定されている。

【0213】

この状態で、クライアント装置200は、図26の場合と同様に、“10 < C01 < 100”を満たすレコードを検索するよう、サーバ装置100、100aに要求するものとする。すると、サーバ装置100、100aは、テーブル(T01)に対し、インデック

50

ス検索により次のようにデータを検索する。

【0214】

サーバ制御部110は、インデックス131、検索に参加するサーバ装置の数(=2台)およびサーバ装置100の稼働優先度(="1")に基づいて、限定範囲#1として“C01<24”を算出する。サーバ装置100からは、データベース120に、“3, aa”および“20, bb”を含み、“25, cc”、“75, dd”および“200, ee”を含まない仮想パーティションが形成されているように見える。

【0215】

サーバ制御部110aは、インデックス131a、検索に参加するサーバ装置の数(=2台)およびサーバ装置100aの稼働優先度(="2")に基づいて、限定範囲#2として“24 C01”を算出する。サーバ装置100aからは、データベース120aに、“25, cc”、“75, dd”および“200, ee”を含み、“3, aa”および“20, bb”を含まない仮想パーティションが形成されているように見える。

10

【0216】

図28は、インデックス検索により検索されるデータの例を示す図(続き)である。

サーバ制御部110は、クライアント装置200から指定された検索範囲と、限定範囲#1とが重複する範囲“10<C01<24”を、部分検索範囲#1として算出する。そして、サーバ制御部110は、部分検索範囲#1に限定してデータベース120からデータを検索する。その結果、“20, bb”が抽出され、サーバ装置100の検索結果115としてクライアント装置200に送信される。

20

【0217】

サーバ制御部110aは、クライアント装置200から指定された検索範囲と、限定範囲#2とが重複する範囲“24 C01<100”を、部分検索範囲#2として算出する。そして、サーバ制御部110aは、部分検索範囲#2に限定してデータベース120aからデータを検索する。その結果、“25, cc”および“75, dd”が抽出され、サーバ装置100aの検索結果115aとしてクライアント装置200に送信される。

【0218】

クライアント装置200は、サーバ装置100の検索結果115(“20, bb”)と、サーバ装置100aの検索結果115a(“25, cc”および“75, dd”)とをマージする。ここでは、検索結果115のレコードと検索結果115aのレコードとを連結して1つのリストにする。マージされた検索結果212は、検索要求に対する応答としてアプリケーションソフトウェア210に提供される。

30

【0219】

図29は、システムによる検索時間の例を示す図である。図29では、同期がとられたデータを格納している3台のサーバ装置(1台の現用系と2台の待機系)を有するシステムにおいて、クライアント装置が検索要求を送信してから検索結果を取得するまでの間の時間について説明する。

【0220】

図29の上側は、第2の実施の形態のシステム以外のシステムにより、データを検索する場合の処理時間の例を示す図である。図29の上側で説明するシステムでは、現用系のサーバ装置のみが検索処理を実行する。時間T1は、現用系のサーバ装置による検索処理の時間である。時間T2は、現用系のサーバ装置からクライアント装置に検索結果を送信するための通信時間である。このシステムによるクライアント装置上のアプリケーションソフトウェアから見た検索のレスポンス時間は、“T1+T2”となる。

40

【0221】

図29の下側は、第2の実施の形態のシステムにより、データを並列検索する場合の処理時間の例を示す図である。時間T3は、3台のサーバ装置それぞれによる検索処理の時間である。時間T4は、3台のサーバ装置それぞれからクライアント装置に検索結果を送信するための通信時間である。また、時間T5は、クライアント装置が3台のサーバ装置から受信した検索結果をマージする時間である。この場合、第2の実施の形態のシステム

50

によるクライアント装置上のアプリケーションソフトウェアから見た検索のレスポンス時間は、“ $T3 + T4 + T5$ ”となる。

【0222】

ここで、第2の実施の形態のシステムでは、同じ検索処理を3台のサーバ装置が並列して実行しているため、時間 $T3$ は、時間 $T1$ の約 $1/3$ となる。また、検索されるデータの総量は第2の実施の形態のシステム以外のシステムと変わらないため、時間 $T4$ は時間 $T2$ とほぼ同じと考えられ、また、時間 $T3$ と比べて十分に短いと考えられる。また、データをマージする処理の時間は、時間 $T3$ と比べると、非常に短いと考えられる。そのため、第2の実施の形態のシステムにおけるレスポンス時間は、図29の上側で説明したシステムの場合と比べ、時間 $T1$ の $2/3$ だけ短縮される。また、各サーバ装置の検索の負荷が大きい場合は、レスポンス時間が全体として $1/3$ 程度になると期待できる。

10

【0223】

第2の実施の形態のシステムによれば、クライアント装置200は、サーバ装置100、100a、100bに対して、同じ検索条件を指定してデータの検索を要求する。次に、データの検索を要求された各サーバ装置は、自己の稼働優先度や、検索に参加するサーバ装置の数等に基づいて、検索条件の示す検索範囲のうち自サーバ装置が担当する部分検索範囲を算出し、算出した部分検索範囲に限定して検索処理を行い、その検索結果をクライアント装置200に送信する。そして、クライアント装置200は、サーバ装置100、100a、100bから受信した検索結果をマージし、マージした検索結果をアプリケーションソフトウェアに提供する。これにより、クライアント装置200から見たレスポンス時間を短縮できる。すなわち、データの検索を高速化できる。

20

【0224】

また、データを検索するときサーバ装置100、100a、100bが異なるデータベースにアクセスするため、共通のデータベースにアクセスする場合と比べて、アクセス競合を抑制できスループットが向上しやすくなる。また、データベース120、120a、120bのデータが同期されているため、あるサーバが故障しても、残りのサーバを用いてデータ全体へのアクセスを確保でき、耐故障性が向上する。

【0225】

また、各サーバ装置が担当する部分検索範囲は、重複しないように算出される。これにより、各サーバ装置による無駄な検索処理を削減できる。

30

また、各サーバ装置は、部分検索範囲を算出する際、B-Tre eインデックス等の木構造のインデックスを用いる。これにより、部分検索範囲の間で検索の負荷（例えば、部分検索範囲に属するデータの量）がほぼ均等になることが期待できるため、各サーバ装置による検索処理の時間が均等化される。よって、検索処理の時間の差により生じる遅延が抑制され、本システムによるデータの検索時間がより短縮される。

【0226】

〔第3の実施の形態〕

次に、第3の実施の形態を説明する。前述の第2の実施の形態との違いを中心に説明し、第2の実施の形態と同様の事項については説明を省略する。第3の実施の形態のシステムでは、各サーバ装置の限定範囲（検索処理を担当する範囲）を算出するときに用いるインデックスの構造が、第2の実施の形態のものと異なる。

40

【0227】

例えば、第2の実施の形態のシステムでは、インデックスのルートノードのキーの数を“検索に参加するサーバ装置の数 - 1”としている。そのため、サーバ装置の故障や復旧により検索に参加するサーバ装置の数が増減する度に、インデックスを更新することになる。

【0228】

そこで、第3の実施の形態では、検索に参加するサーバ装置の数が増減しても更新が不要になるようなインデックスを生成することにする。ただし、第3の実施の形態では、検索に参加するサーバ装置の数は、“1～システムの運用当初のサーバ装置の数”の範囲内

50

で変化するものとする。すなわち、サーバ装置の数が運用当初よりも増えない限り、サーバ装置の故障や復旧が生じて、インデックスを再構成しなくてよいようにする。

【0229】

図30は、インデックスの変形例を示す図である。インデックス132は、図6等で説明した第2の実施の形態のインデックス131の変形例である。インデックス132は、システムの運用当初の検索に参加するサーバ装置の数を“x”とすると（以下、同様とする）、“x”が2以上である場合、ルートノードのキーの数が「（“1”から“x”までの自然数の最小公倍数）-1」となるように生成される。“x”が1の場合は、上記の方法でインデックスを生成するとルートノードのキーの数が0になるため、ルートノードのキーの数が1以上の任意の数となるようにインデックスを生成する。

10

【0230】

例えば、システムの運用当初の検索に参加するサーバ装置の数が“3”とすると、インデックスのルートノードのキーの数は、「（“1”、“2”、“3”の最小公倍数）-1」により“6-1=5”となる。

【0231】

このように、ルートノードのキーの数を「（“1”から“x”までの自然数の最小公倍数）-1」とすることで、運用当初のサーバ装置のうちの任意の台数が故障しても、インデックスを生成し直さなくてよくなる。

【0232】

例えば、インデックス132は、ルートノードに“4”、“11”、“18”、“25”および“32”の5つのキーが設定されている。この状態で、3台のサーバ装置で検索処理を並列して実行する場合、図30の上側のように、“ $p < 11$ ”、“ $11 \leq p < 25$ ”および“ $25 \leq p$ ”の限定範囲が算出される。また、この状態で、3台のサーバ装置の何れか1台が故障した場合、検索に参加するサーバ装置の数は2台となり、図30の下側のように、“ $p < 18$ ”および“ $18 \leq p$ ”限定範囲が算出される。

20

【0233】

このように、“ルートノードのキーの数+1”は、“1”からシステムの運用当初のサーバ装置の数までの全ての自然数で割りきれられる。このため、検索に参加するサーバ装置が任意の数だけ故障しても、ルートノードのキーによって分割される区間を、稼働中のサーバ装置に均等に割り振ることができる。よって、インデックス132を再構成しなくても、稼働中の複数のサーバ装置の間で、検索の負荷をほぼ均等にすることが可能となる。

30

【0234】

図31は、限定範囲の算出処理の変形例を示すフローチャートである。図31の処理は、第2の実施の形態の限定範囲算出（図22の処理）に代えて、前述のステップS49において実行される。第2の実施の形態との違いは、ステップS81とステップS83の間にステップS81aが追加される点と、ステップS84、S86、S87に代えてステップS84a、S86a、S87aが実行される点である。以下、ステップS81a、S84a、S86a、S87aについて説明する。

【0235】

（ステップS81a）処理内容判定部111は、“（ルートノードのキーの数+1）/ 検索サーバ装置数a”により変数“n”の値を算出する。

40

（ステップS84a）処理内容判定部111は、限定範囲を“ $p < R(n)$ ”と算出する。そして、限定範囲算出を終了する。

【0236】

（ステップS86a）処理内容判定部111は、限定範囲を“ $R((稼働優先度b-1) * n) \leq p$ ”と算出する。そして、限定範囲算出を終了する。

（ステップS87a）処理内容判定部111は、限定範囲を“ $R((稼働優先度b-1) * n) < p \leq R(稼働優先度b * n)$ ”と算出する。

【0237】

図30～図31で説明したように、各サーバ装置は、ルートノードのキーの数が「（“

50

「 1 」から「 x 」までの自然数の最小公倍数) - 1」となるように生成されたインデックスを用いて限定範囲を算出し、この限定範囲から部分検索範囲を算出する。これにより、検索に参加するサーバ装置の数が増加したとき、インデックスを更新する処理が不要となるため、各サーバ装置の処理の負荷を軽減できる。

【 0 2 3 8 】

なお、第 3 の実施の形態のシステムでは、ルートノードのキーの数を「 (「 1 」から「 x 」までの自然数の最小公倍数) - 1」となるように生成するが、「 (「 1 」から「 x 」までの自然数の公倍数) - 1」としても同様の効果を得ることが可能である。例えば、「 1 」から「 x 」までの自然数の公倍数」として、最小公倍数を N 倍 (N は 2 以上の整数) したものをを用いることができる。また、例えば、「 1 」から「 x 」までの自然数の公倍数」として、「 1 」から「 x 」までの自然数の積を用いることができる。

10

【 0 2 3 9 】

[第 4 の実施の形態]

次に、第 4 の実施の形態を説明する。前述の第 2 の実施の形態との違いを中心に説明し、第 2 の実施の形態と同様の事項については説明を省略する。第 4 の実施の形態のシステムは、複数のサーバ装置の間の更新処理と検索処理の分担が第 2 の実施の形態と異なる。

【 0 2 4 0 】

図 3 2 は、サーバ装置の機能構成の変形例を示す図である。第 4 の実施の形態のシステムは、サーバ装置 1 0 0 - 1 , 1 0 0 a - 1 , 1 0 0 b - 1 を有する。サーバ装置 1 0 0 - 1 は、現用系のサーバ装置である。第 4 の実施の形態では、現用系のサーバ装置は、更新処理のみ実行し、検索処理を実行しない。そのため、サーバ装置 1 0 0 - 1 は、インデックス (例えば、インデックス 1 3 1 a , 1 3 1 b に相当するもの) を有さなくてもよい。

20

【 0 2 4 1 】

サーバ装置 1 0 0 a - 1 , 1 0 0 b - 1 は、待機系のサーバ装置である。第 4 の実施の形態では、待機系のサーバ装置は、検索処理のみ実行し、更新処理を実行しない。サーバ装置 1 0 0 a - 1 , 1 0 0 b - 1 は、並列に検索処理を実行することができる。サーバ装置 1 0 0 a - 1 はインデックス 1 3 1 a - 1 を有し、サーバ装置 1 0 0 b - 1 はインデックス 1 3 1 b - 1 を有する。インデックス 1 3 1 a - 1 , 1 3 1 b - 1 それぞれは、「検索に参加するサーバ装置の数 - 1」の数のキーをルートノードに有する。例えば、図 3 2 では、検索に参加するサーバ装置の数は、稼働中の待機系のサーバ装置の数である「 2 」であるため、ルートノードのキーの数は、「 1 」となる。

30

【 0 2 4 2 】

図 3 3 は、システムによる検索の変形例を示す図である。サーバ装置 1 0 0 - 1 はサーバ制御部 1 1 0 - 1 を有し、サーバ装置 1 0 0 a - 1 はサーバ制御部 1 1 0 a - 1 を有し、サーバ装置 1 0 0 b - 1 はサーバ制御部 1 1 0 b - 1 を有する。サーバ制御部 1 1 0 - 1 , 1 1 0 a - 1 , 1 1 0 b - 1 は、第 2 の実施の形態のサーバ制御部 1 1 0 , 1 1 0 a , 1 1 0 b に対応する。

【 0 2 4 3 】

クライアント装置 2 0 0 は、サーバ装置 1 0 0 - 1 , 1 0 0 a - 1 , 1 0 0 b - 1 に、同じ検索条件を指定した検索要求を送信する (S 2 , S 2 a , S 2 b)。サーバ制御部 1 1 0 - 1 は、クライアント装置 2 0 0 からデータの検索を要求されたとき、検索処理を実行せず、クライアント装置 2 0 0 へ応答しない。

40

【 0 2 4 4 】

サーバ制御部 1 1 0 a - 1 , 1 1 0 b - 1 それぞれは、クライアント装置 2 0 0 からデータの検索を要求されたとき、自己のサーバ装置の種別と検索に参加するサーバ装置の数とに基づき部分検索範囲を算出する。例えば、サーバ制御部 1 1 0 a - 1 は部分検索範囲 # 1 を算出し、サーバ制御部 1 1 0 b - 1 は部分検索範囲 # 2 を算出する (S 3 a - 1 , S 3 b - 1)。そして、サーバ制御部 1 1 0 a - 1 は、部分検索範囲 # 1 に限定してデータベース 1 2 0 a からデータを検索し、クライアント装置 2 0 0 に検索結果を送信する (

50

S 4 a)。サーバ制御部 1 1 0 b - 1 は、部分検索範囲 # 2 に限定してデータベース 1 2 0 b からデータを検索し、クライアント装置 2 0 0 に検索結果を送信する (S 4 b)。

【 0 2 4 5 】

第 4 の実施の形態のシステムでは、現用系のサーバ装置 1 0 0 - 1 は更新処理のみを担当し、待機系のサーバ装置 1 0 0 a - 1 , 1 0 0 b - 1 は検索処理のみを担当する。これにより、現用系のサーバ装置 1 0 0 - 1 の負荷を軽減し、データの更新および並列検索をサーバ装置 1 0 0 - 1 , 1 0 0 a - 1 , 1 0 0 b - 1 の間で適切に分担することができる。

【 0 2 4 6 】

なお、前述のように、第 1 の実施の形態の情報処理は、サーバ 1 0 , 1 0 a や検索要求装置 2 0 にプログラムを実行させることで実現できる。第 2 および第 3 の実施の形態の情報処理は、サーバ装置 1 0 0 , 1 0 0 a , 1 0 0 b やクライアント装置 2 0 0 にプログラムを実行させることで実現できる。第 4 の実施の形態の情報処理は、サーバ装置 1 0 0 - 1 , 1 0 0 a - 1 , 1 0 0 b - 1 にプログラムを実行させることで実現できる。このようなプログラムは、コンピュータ読み取り可能な記録媒体 (例えば、記録媒体 4 3) に記録しておくことができる。記録媒体としては、例えば、磁気ディスク、光ディスク、光磁気ディスク、半導体メモリ等を使用できる。磁気ディスクには、F D および H D D が含まれる。光ディスクには、C D、C D - R (Recordable) / R W (Rewritable)、D V D および D V D - R / R W が含まれる。

【 0 2 4 7 】

プログラムを流通させる場合、例えば、当該プログラムを記録した可搬記録媒体が提供される。コンピュータは、例えば、可搬記録媒体に記録されたプログラムを、記憶装置 (例えば、H D D 1 0 3) に格納し、当該記憶装置からプログラムを読み込んで実行する。ただし、可搬記録媒体から読み込んだプログラムを直接実行してもよい。また、上記の情報処理の少なくとも一部を、D S P、A S I C、P L D (Programmable Logic Device) 等の電子回路で実現することも可能である。

【符号の説明】

【 0 2 4 8 】

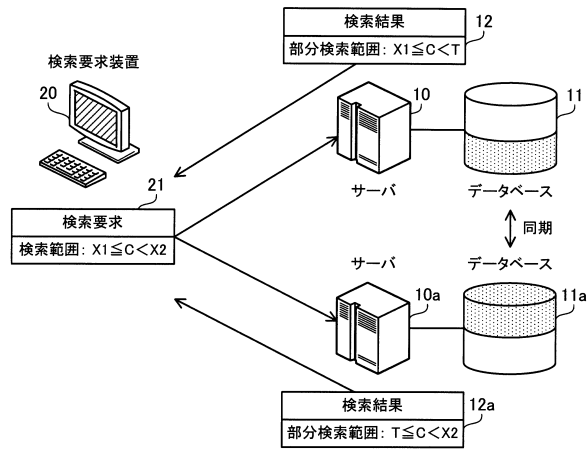
- 1 0 , 1 0 a サーバ
- 1 1 , 1 1 a データベース
- 1 2 , 1 2 a 検索結果
- 2 0 検索要求装置
- 2 1 検索要求

10

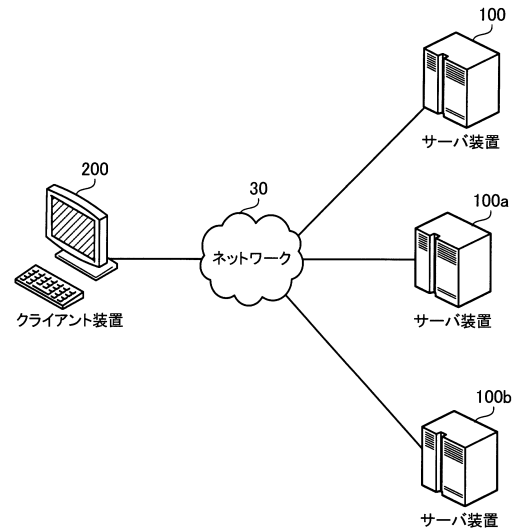
20

30

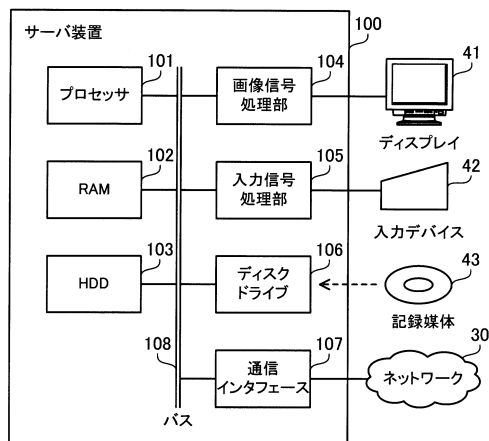
【図 1】



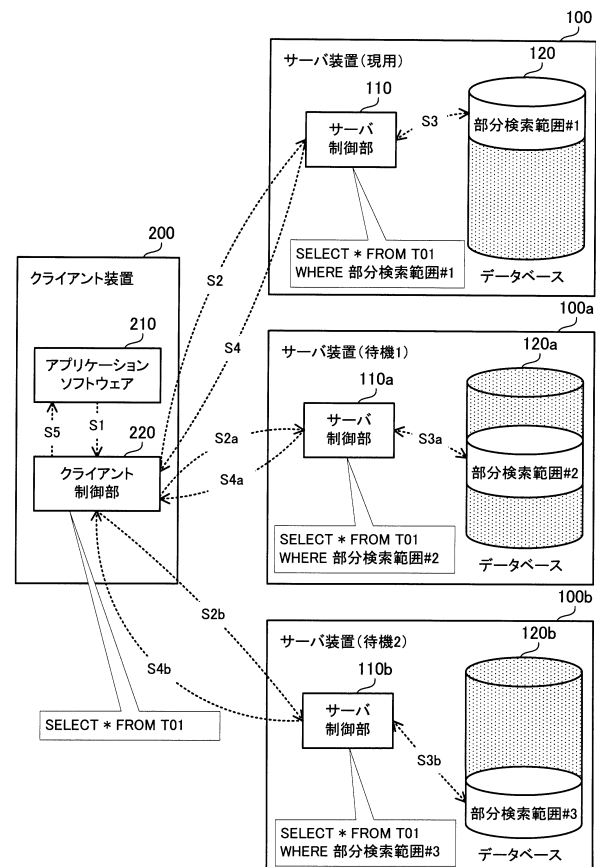
【図 2】



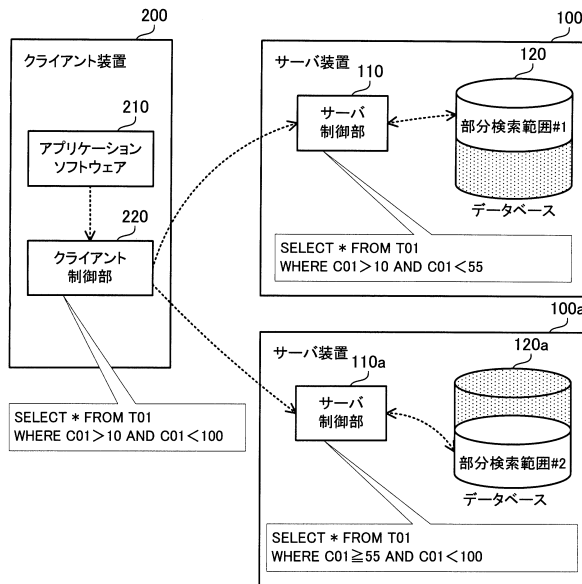
【図 3】



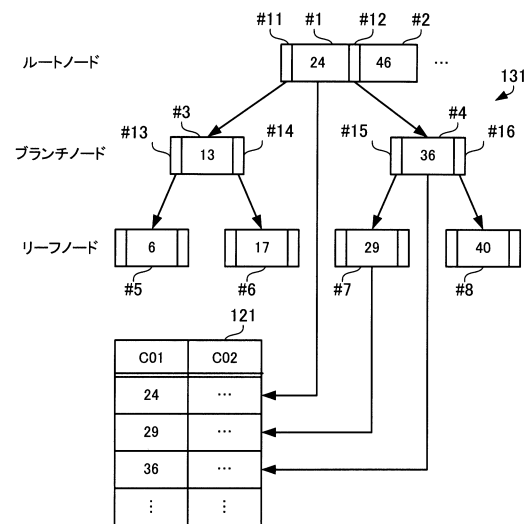
【図 4】



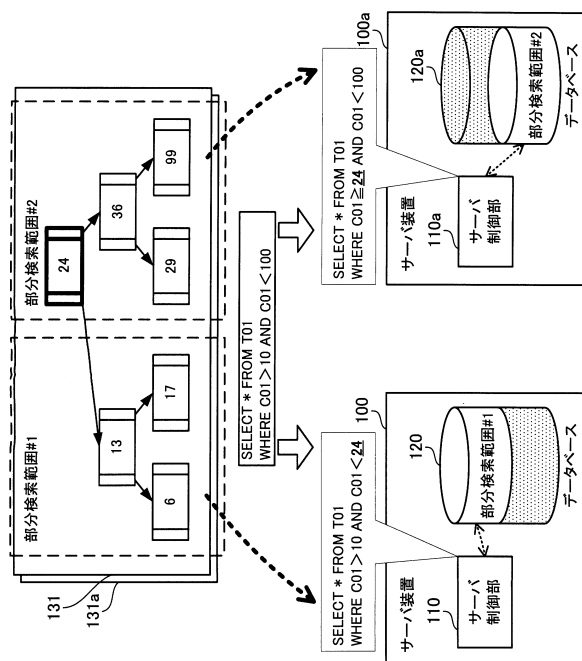
【図 5】



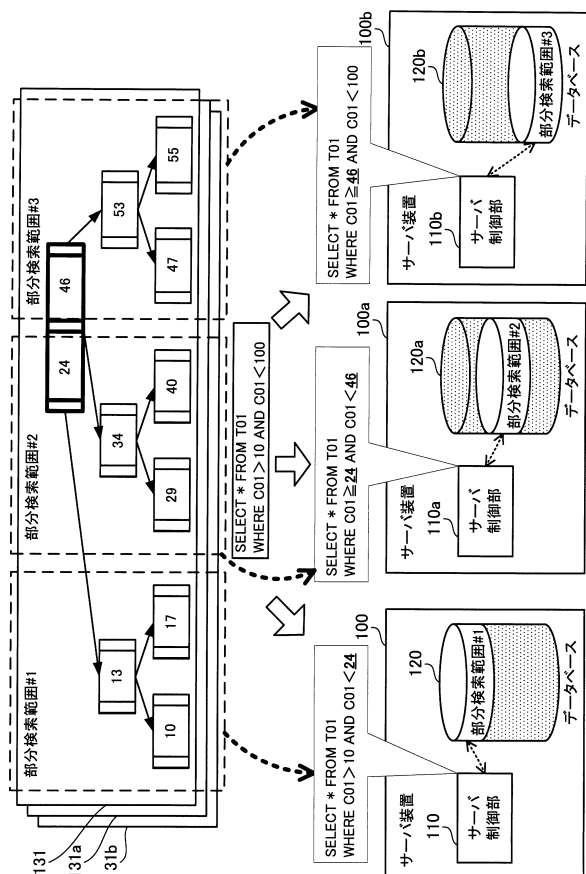
【図 6】



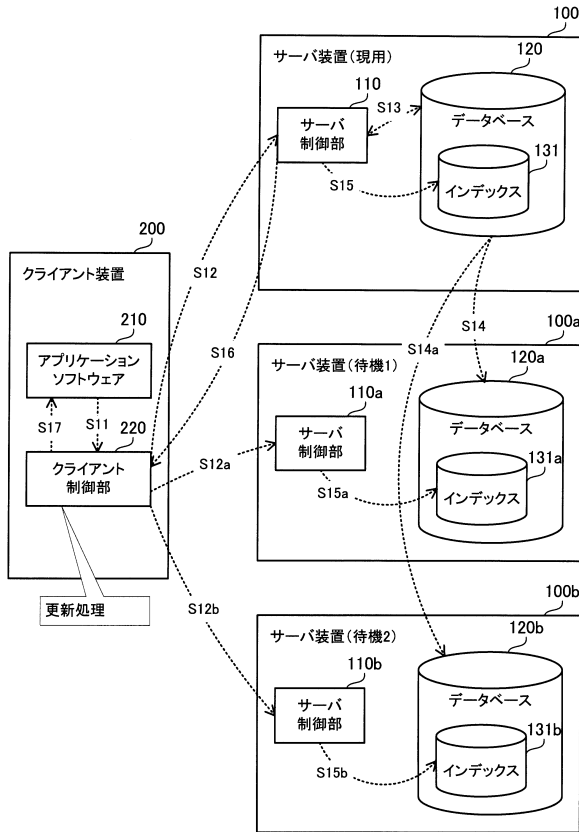
【図 7】



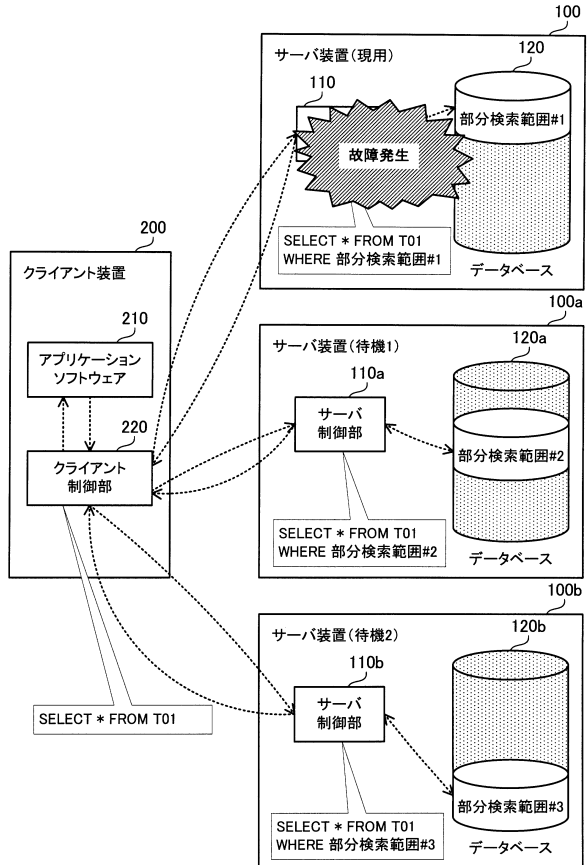
【図 8】



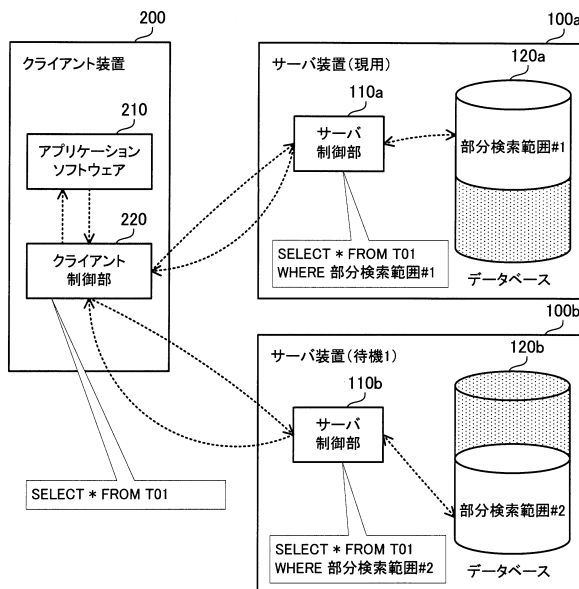
【図 9】



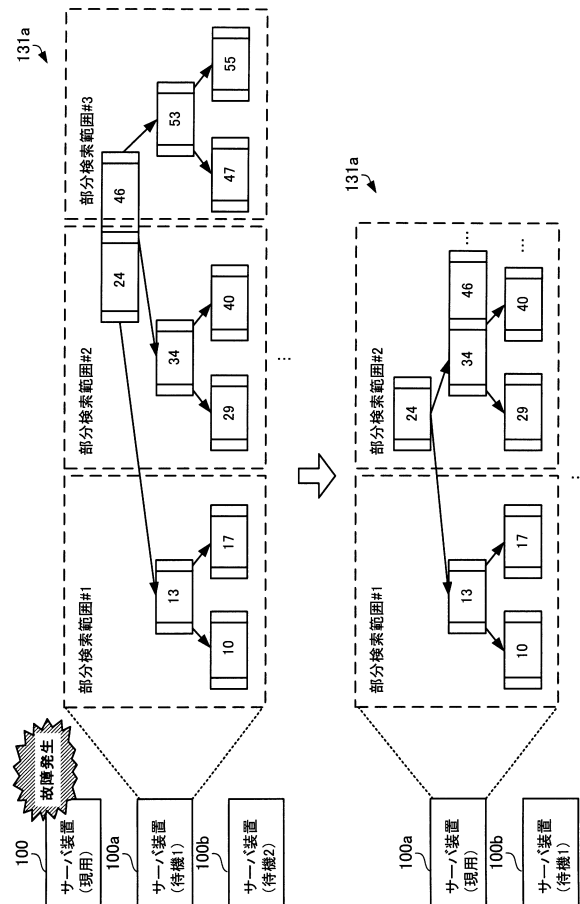
【図 10】



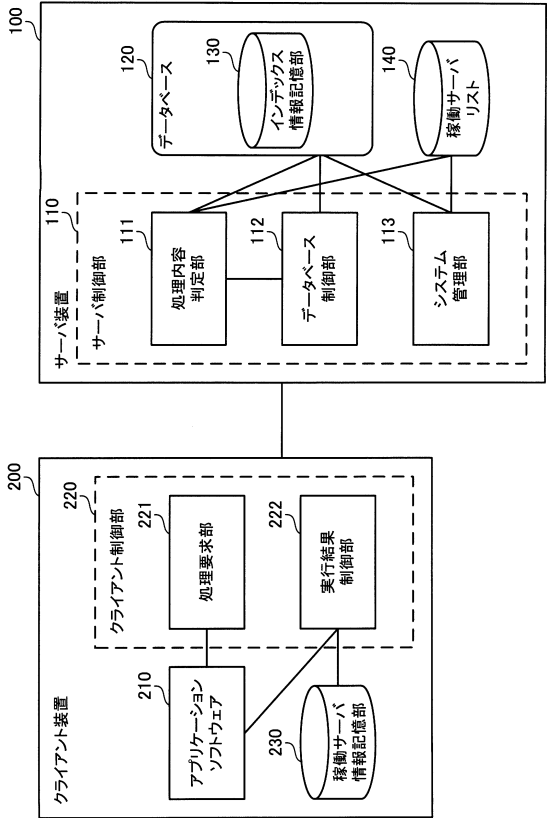
【図 11】



【図 12】



【図 13】



【図 14】

処理要求通知					
制御情報	並列フラグ	種別	テーブル	列	条件
情報A	TRUE	検索	T01	C01,C02	C01>10 AND C01<100

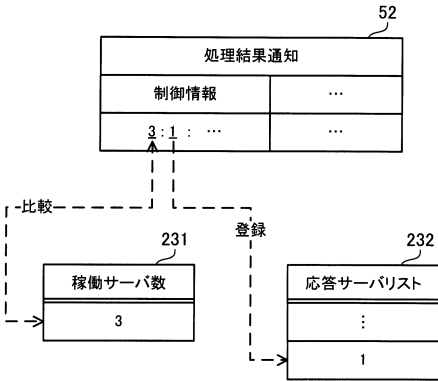
【図 15】

稼働サーバリスト	
優先度	サーバ
1	SV#A → 現用
2	SV#B → 待機1
3	SV#C → 待機2
⋮	

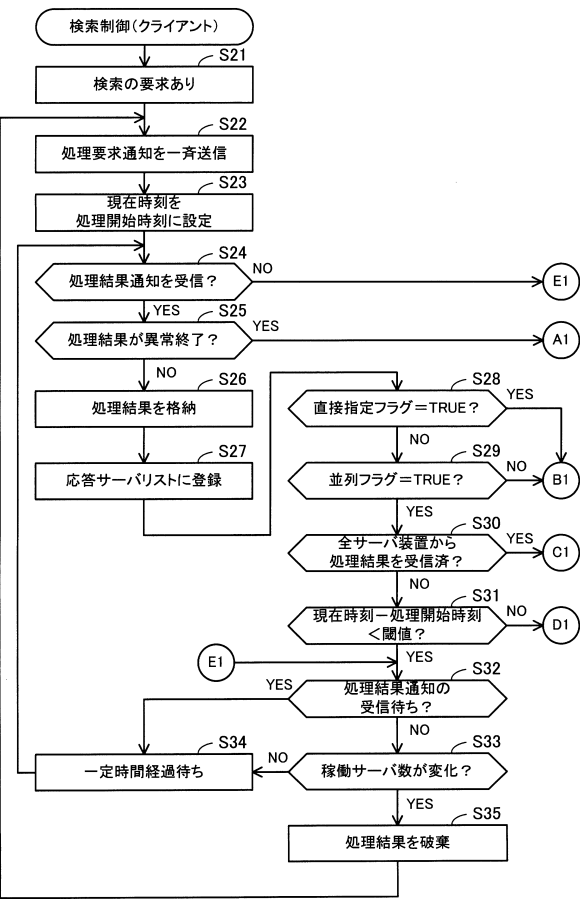
【図 16】

処理結果通知		
制御情報	直接指定フラグ	処理結果
3 : 1 : 情報A	TRUE	(20, aa), (25, bb)

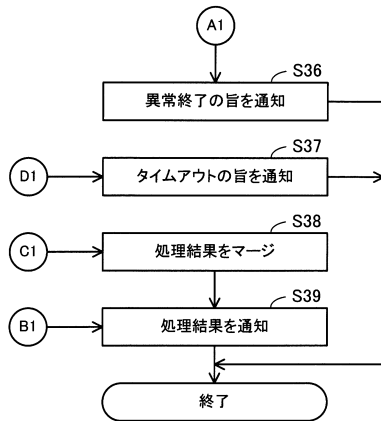
【図 17】



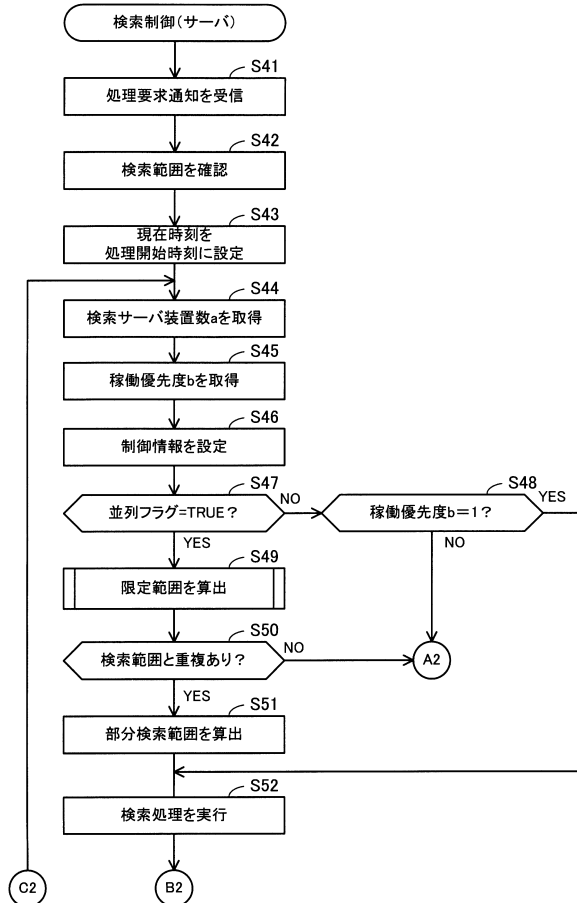
【図 18】



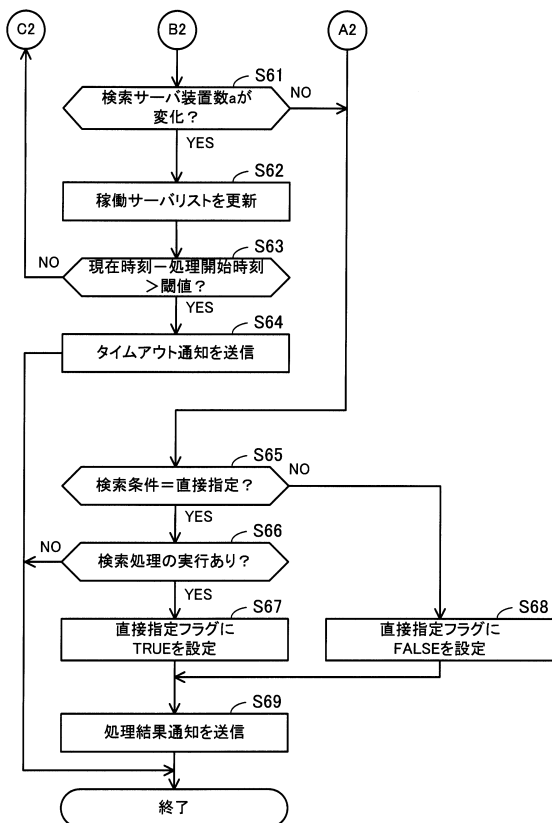
【図 19】



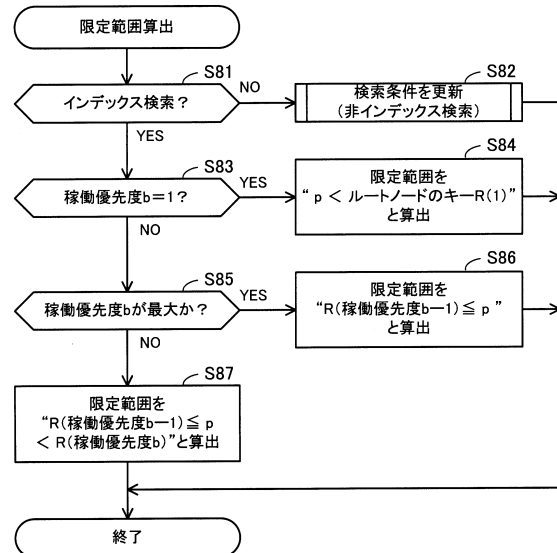
【図 20】



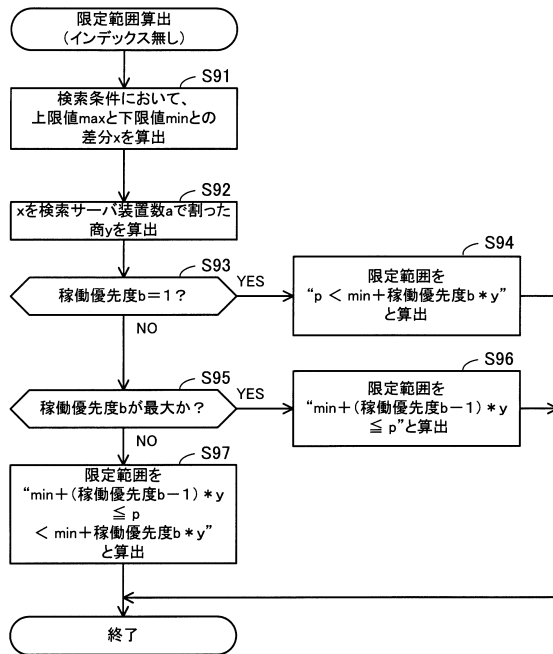
【図 21】



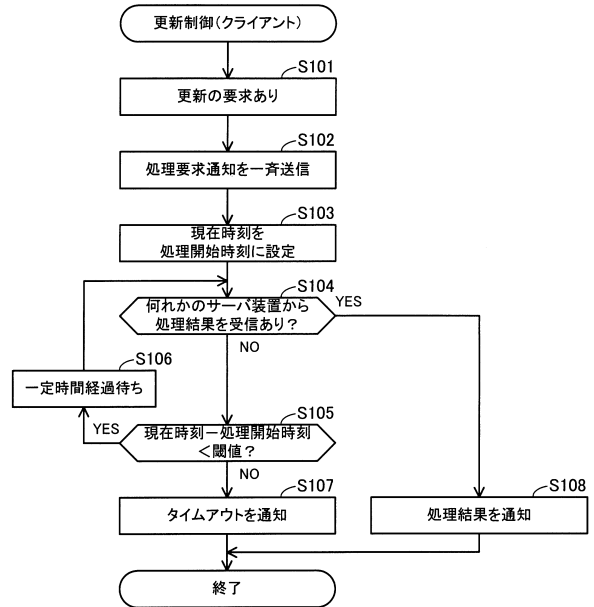
【図 22】



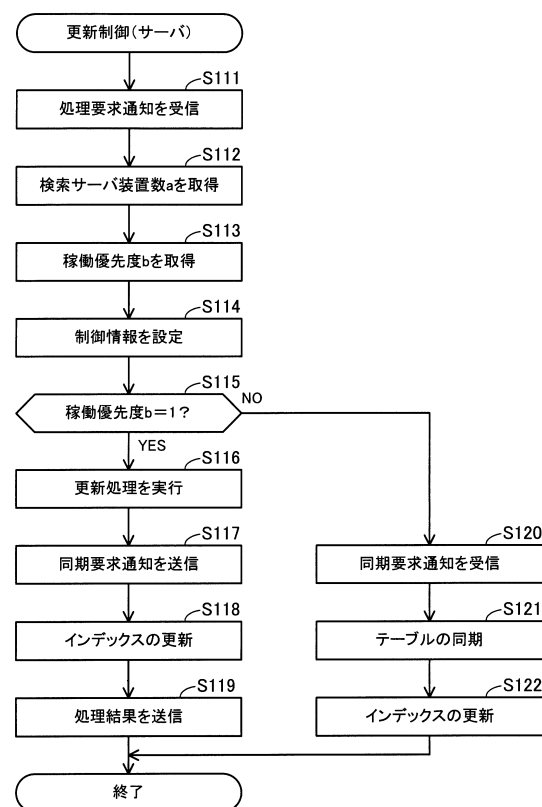
【図 23】



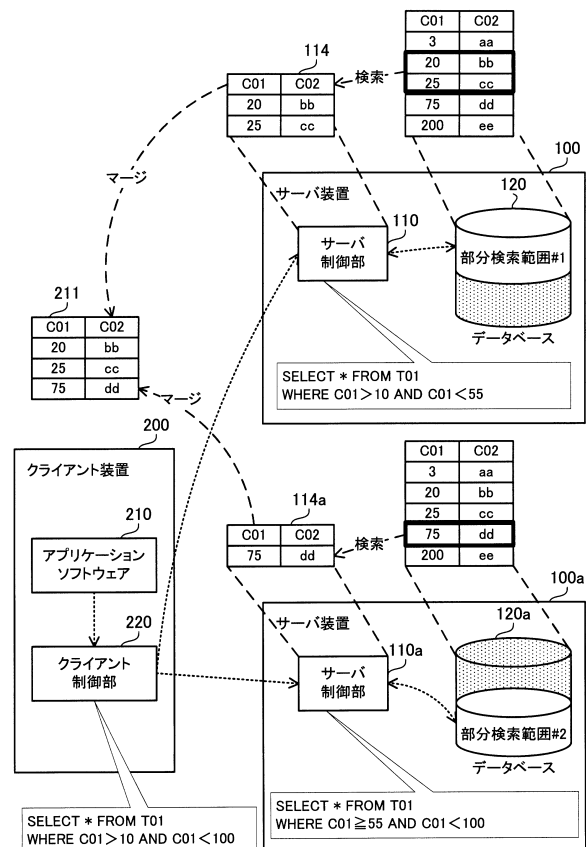
【図 24】



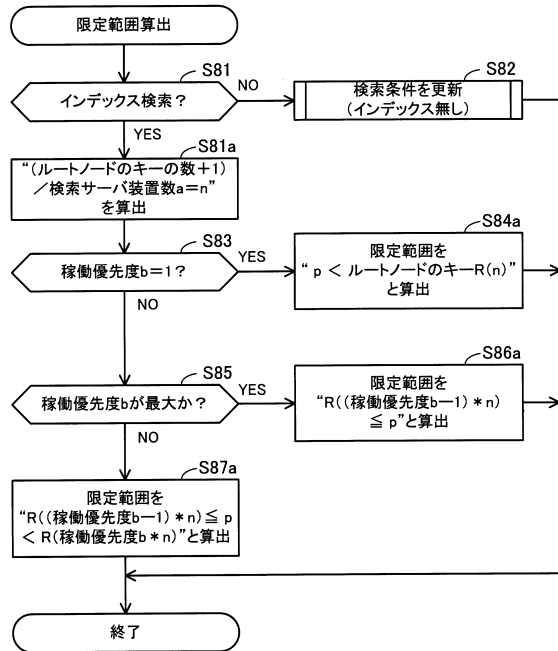
【図 25】



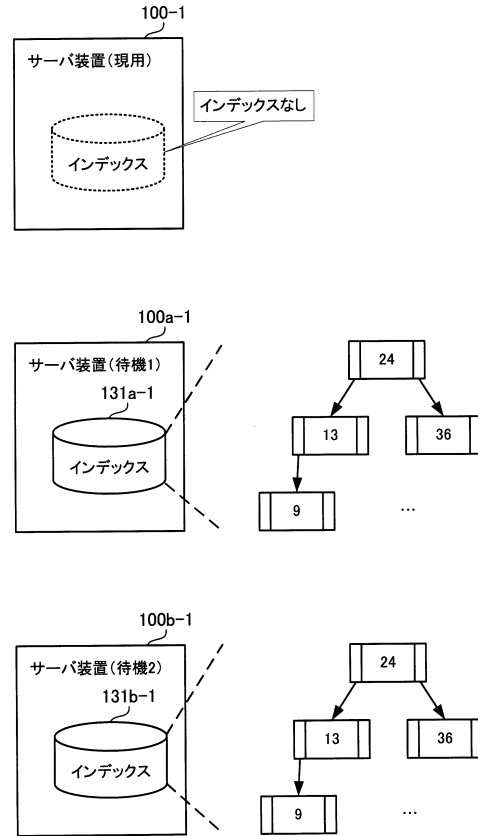
【図 26】



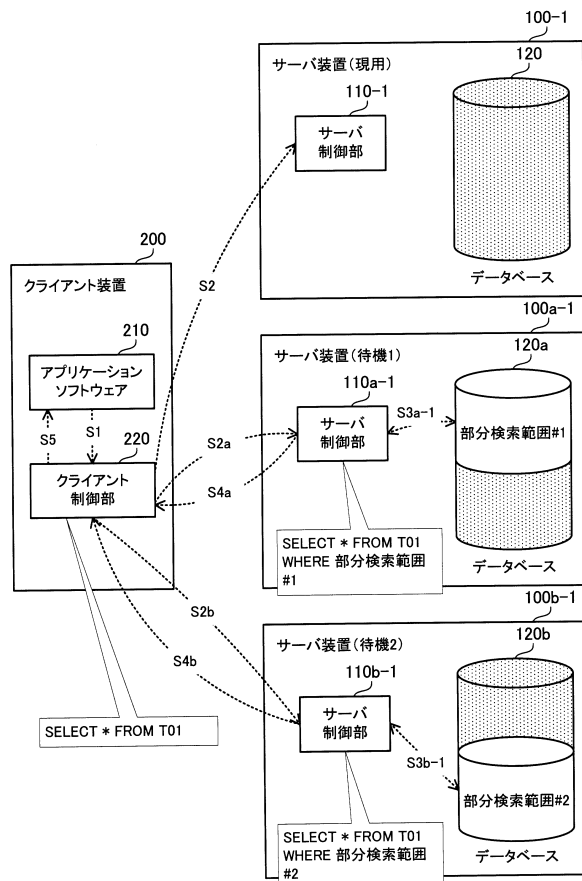
【 図 3 1 】



【 図 3 2 】



【 図 3 3 】



フロントページの続き

審査官 小太刀 慶明

(56)参考文献 特開2008-250722(JP,A)
特開平07-160557(JP,A)
米国特許出願公開第2010/0082655(US,A1)
米国特許第06439783(US,B1)
米国特許出願公開第2007/0239759(US,A1)

(58)調査した分野(Int.Cl., DB名)
G06F 17/30
G06F 12/00
G06F 13/00