

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6268116号
(P6268116)

(45) 発行日 平成30年1月24日(2018.1.24)

(24) 登録日 平成30年1月5日(2018.1.5)

(51) Int.Cl. F 1
G 0 6 F 1 2 / 0 0 (2 0 0 6 . 0 1) G 0 6 F 1 2 / 0 0 5 6 0 G

請求項の数 16 (全 21 頁)

(21) 出願番号	特願2015-58741 (P2015-58741)	(73) 特許権者	317006041 東芝メモリ株式会社 東京都港区芝浦一丁目1番1号
(22) 出願日	平成27年3月20日(2015.3.20)	(74) 代理人	100091982 弁理士 永井 浩之
(65) 公開番号	特開2016-177688 (P2016-177688A)	(74) 代理人	100091487 弁理士 中村 行孝
(43) 公開日	平成28年10月6日(2016.10.6)	(74) 代理人	100082991 弁理士 佐藤 泰和
審査請求日	平成29年3月7日(2017.3.7)	(74) 代理人	100105153 弁理士 朝倉 悟
		(74) 代理人	100107582 弁理士 関根 毅
		(74) 代理人	100118843 弁理士 赤岡 明

最終頁に続く

(54) 【発明の名称】 データ処理装置、データ処理方法およびコンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

キーを指定した読み出し要求を受け付け、前記読み出し要求を解釈する要求解釈部と、複数のキーを所定のルールに従った順序で保持する順序キー列における前記キーに対して前記読み出し要求で指示された位置関係にあるキーを特定する第1アクセス部と、

第1アドレスと第2アドレスとを対応づけて管理する内部または外部の第1記憶装置から、前記第1アクセス部で特定されたキーのハッシュ値に基づく第1アドレスに対応する第2アドレスを取得する第2アクセス部と、

内部または外部の第2記憶装置において前記第2アクセス部で取得された前記第2アドレスに格納されたデータを読み出す第3アクセス部と、
を備えたデータ処理装置。

【請求項2】

前記要求解釈部で解釈された前記読み出し要求の内容に応じて、前記第1アクセス部に前記キーと前記読み出し要求で指示された位置関係の情報とを通知し、前記第1アクセス部で特定されたキーを受け取り、受け取ったキーを前記第2アクセス部に通知する処理振分部

をさらに備えた請求項1に記載のデータ処理装置。

【請求項3】

前記第2記憶装置の第2アドレスを管理するアドレス割当部を備え、

前記要求解釈部は、キーとバリューとを指定した書き込み要求を受け付け、前記書き込

み要求を解釈し、

前記第 1 アクセス部は、前記書き込み要求で指定されたキーを前記順序キー列に追加することにより前記順序キー列を更新し、

前記アドレス割当部は、前記書き込み要求で指定されたキーに対して第 2 アドレスを割り当て、

前記第 2 アクセス部は、前記第 1 記憶装置において前記指定されたキーのハッシュ値に基づく第 1 アドレスに対応づけて前記第 2 アドレスを登録し、

前記第 3 アクセス部は、前記第 2 記憶装置における前記第 2 アドレスに前記バリューを含むデータを書き込む

請求項 1 または 2 に記載のデータ処理装置。

10

【請求項 4】

前記順序キー列は、前記第 2 記憶装置に格納されており、

前記第 1 記憶装置は、内部キーのハッシュ値に基づく第 1 アドレスと、前記順序キー列が格納された第 2 アドレスとの対応を管理し、

前記第 2 アクセス部は、前記第 1 記憶装置から前記内部キーに対応する前記第 2 アドレスを取得し、

前記第 3 アクセス部は、前記第 2 アクセス部で取得された前記第 2 アドレスに格納されたデータである前記順序キー列を前記第 2 記憶装置から読み出し、

前記第 1 アクセス部は、前記第 3 アクセス部により読み出された前記順序キー列を用いて前記キーを特定する

20

請求項 1 ないし 3 のいずれか一項に記載のデータ処理装置。

【請求項 5】

前記順序キー列は、前記第 2 記憶装置に格納されており、

前記第 1 記憶装置は、内部キーのハッシュ値に基づく第 1 アドレスと、前記順序キー列が格納された第 2 アドレスとの対応を管理し、

前記第 2 アクセス部は、前記第 1 記憶装置から前記内部キーに対応する前記第 2 アドレスを取得し、

前記第 3 アクセス部は、前記第 2 アクセス部で取得された前記第 2 アドレスに格納されたデータに含まれる前記順序キー列を前記第 2 記憶装置から読み出し、

前記第 1 アクセス部は、前記第 3 アクセス部により読み出された前記順序キー列に、前記書き込み要求で指定されたキーを追加することにより前記順序キー列を更新し、

30

前記アドレス割当部は、前記更新された順序キー列に対して第 2 アドレスを割り当て、

前記第 2 アクセス部は、前記第 1 記憶装置において前記内部キーのハッシュ値に基づく第 1 アドレスに対応づけて前記割り当てられた第 2 アドレスを登録し、

前記第 3 アクセス部は、前記第 2 記憶装置における前記割り当てられた第 2 アドレスに前記順序キー列を含むデータを書き込む

請求項 3 に記載のデータ処理装置。

【請求項 6】

前記内部キーを生成する内部キー生成部

をさらに備えた請求項 4 または 5 に記載のデータ処理装置。

40

【請求項 7】

前記順序キー列を複数に分割した各部分に、それぞれ異なる内部キーが割り当てられ、

前記内部キー生成部は、前記読み出し要求または前記書き込み要求で指定されたキーが、前記複数の各部分のいずれに属するかに応じて前記内部キーを生成する

請求項 6 に記載のデータ処理装置。

【請求項 8】

前記順序キー列は、前記複数のキーにそれぞれ対応づけて前記第 2 アドレスを管理し、

前記第 1 アクセス部は、前記特定したキーに対応する前記第 2 アドレスを特定し、

前記第 3 アクセス部は、前記第 2 記憶装置における前記特定した第 2 アドレスからデータを読み出す

50

請求項 1 ないし 7 のいずれか一項に記載のデータ処理装置。

【請求項 9】

前記順序キー列は、前記複数のキーを辞書順で管理する

請求項 1 ないし 8 のいずれか一項に記載のデータ処理装置。

【請求項 10】

前記読み出し要求は、1つのキーを指定し、前記キーの次の位置のキー、前記キーの前の位置のキー、またはこれらの両方を読み出すことを要求する、

もしくは、

前記読み出し要求は、2つのキーを含み、前記2つのキーの間に位置するキーの読み出しを要求する

10

請求項 1 ないし 9 のいずれか一項に記載のデータ処理装置。

【請求項 11】

キーを指定した読み出し要求を受け付け、前記読み出し要求を解釈する要求解釈ステップと、

複数のキーを所定のルールに従った順序で保持する順序キー列における前記キーに対して前記読み出し要求で指示された位置関係にあるキーを特定する第1アクセスステップと

、
第1アドレスと第2アドレスとを対応づけて管理する内部または外部の第1記憶装置から、前記第1アクセス部で特定されたキーのハッシュ値に基づく第1アドレスに対応する第2アドレスを取得する第2アクセスステップと、

20

内部または外部の第2記憶装置において前記第2アクセス部で取得された前記第2アドレスに格納されたデータを読み出す第3アクセスステップと、

をコンピュータが実行するデータ処理方法。

【請求項 12】

前記要求解釈ステップで解釈された前記読み出し要求の内容に応じて、前記第1アクセスステップに前記キーと前記読み出し要求で指示された位置関係の情報とを通知し、前記第1アクセスステップで特定されたキーを受け取り、受け取ったキーを前記第2アクセスステップに通知する処理振分ステップ

をさらに備えた請求項 11 に記載のデータ処理方法。

【請求項 13】

30

前記第2記憶装置の第2アドレスを管理するアドレス割当ステップを備え、

前記要求解釈ステップは、キーとバリュウとを指定した書き込み要求を受け付け、前記書き込み要求を解釈し、

前記第1アクセスステップは、前記書き込み要求で指定されたキーを前記順序キー列に追加することにより前記順序キー列を更新し、

前記アドレス割当ステップは、前記書き込み要求で指定されたキーに対して第2アドレスを割り当て、

前記第2アクセスステップは、前記第1記憶装置において前記指定されたキーのハッシュ値に基づく第1アドレスに対応づけて前記第2アドレスを登録し、

前記第3アクセスステップは、前記第2記憶装置における前記第2アドレスに前記バリュウを含むデータを書き込む

40

請求項 11 または 12 に記載のデータ処理方法。

【請求項 14】

キーを指定した読み出し要求を受け付け、前記読み出し要求を解釈する要求解釈ステップと、

複数のキーを所定のルールに従った順序で保持する順序キー列における前記キーに対して前記読み出し要求で指示された位置関係にあるキーを特定する第1アクセスステップと

、
第1アドレスと第2アドレスとを対応づけて管理する内部または外部の第1記憶装置から、前記第1アクセス部で特定されたキーのハッシュ値に基づく第1アドレスに対応する

50

第2アドレスを取得する第2アクセスステップと、

内部または外部の第2記憶装置において前記第2アクセス部で取得された前記第2アドレスに格納されたデータを読み出す第3アクセスステップと、

をコンピュータに実行させるためのコンピュータプログラム。

【請求項15】

前記要求解釈ステップで解釈された前記読み出し要求の内容に応じて、前記第1アクセスステップに前記キーと前記読み出し要求で指示された位置関係の情報とを通知し、前記第1アクセスステップで特定されたキーを受け取り、受け取ったキーを前記第2アクセスステップに通知する処理振分ステップ

をさらに備えた請求項14に記載のコンピュータプログラム。

10

【請求項16】

前記第2記憶装置の第2アドレスを管理するアドレス割当ステップをさらに前記コンピュータに実行させ、

前記要求解釈ステップは、キーとバリューとを指定した書き込み要求を受け付け、前記書き込み要求を解釈し、

前記第1アクセスステップは、前記書き込み要求で指定されたキーを前記順序キー列に追加することにより前記順序キー列を更新し、

前記アドレス割当ステップは、前記書き込み要求で指定されたキーに対して第2アドレスを割り当て、

前記第2アクセスステップは、前記第1記憶装置において前記指定されたキーのハッシュ値に基づく第1アドレスに対応づけて前記第2アドレスを登録し、

20

前記第3アクセスステップは、前記第2記憶装置における前記第2アドレスに前記バリューを含むデータを書き込む

請求項14または15に記載のコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、データ処理装置、データ処理方法およびコンピュータプログラムに関する。

【背景技術】

30

【0002】

近年、Webメールやソーシャルネットワークなどに代表される高度なWebサービスが急速に普及している。このようなWebサービスでは、WebサーバがユーザからのWebページ表示要求を受けると、Webサーバは必要な個々のテキスト、画像や動画などのデータを、逐一バックエンドのストレージに問い合わせ取得し、ユーザに回答する。このようなストレージとしては、従来からのファイルシステムベースのデータベースも使われているが、近年では、一意のURLを指定してコンテンツデータにアクセスするという、インターネット上のコンテンツアクセス方式と親和性の高いオブジェクトストレージが、活用され始めている。このようなオブジェクトストレージでは、Webサーバといったクライアント装置から、イーサネット（登録商標）越しに可変長の一意のKey（キー）で、可変長のValue（バリュー）にアクセスするというキーバリューストア（KV S）型のI/Fがとられている。Webサーバの場合、キーはURL、バリューはコンテンツデータであることが一般的である。キーに関して、オブジェクトIDが可変長でなく、固定長になるものもあるが、基本的な原理等は同様である。ストレージデバイスとしてはHDDやSSDが用いられる。

40

【0003】

最もシンプルなアクセスの方法は、例えばKey = "001"に対する読み出しアクセスは、GET (Key = "001") など、読み出したいバリューに対応するキーを指定した要求を送り、応答としてバリューを受けるものである。また、書き込みアクセスは、PUT (Key = "002", Value = "GHI") など、書き込みを行いたいキー

50

とバリューを指定して、要求を送るという形態である。なお、応答として、指定したキーに対応するバリューのみを受け取る形態、指定したキーと、対応するバリューとの両方を受け取る形態など種々の態様がある。

【0004】

より高度なアクセスとして、指定したキーの辞書順で、次の/前のキーに対応するバリューを読み出す、GETNEXT要求/GETPREV要求などもある。例えば (Key, Value) = (“000”, “ABC”)、 (“002”, “DEF”)、 (“004”, “GHI”)、 (“005”, “JKL”) という、4つのキーおよびバリューが格納されているとする。このとき、GETNEXT (Key = “002”) の要求を送ると、応答として次のキーである Key = “004” に対応するバリューである “GHI” が返される。GETPREV (Key = “002”) の場合は、前のキーである Key = “000” に対応するバリューである “ABC” が返される。

10

【0005】

このような要求は、大きなデータを複数のバリューに分割して、各バリューにキーを対応付けておき、各キーに対応するバリューを順次読み出す場合などに用いられる。なお、要求で指定されたキーが存在しない場合もある。例えば Key = “001” などのような場合もある。その場合でも、GETPREV (Key = “001”) に対する応答は、Key = “000” に対応するバリューである “ABC” であり、GETNEXT (Key = “001”) に対する応答は、Key = “002” に対応するバリューである “DEF” となる。

20

【0006】

同様に2つのキーの間に存在する全てのキーおよびバリューを応答するGETRANGE要求もある。例えばGETRANGE (Key 1 = “001”, Key 2 = “004”) の要求を送ると、(Key, Value) = (“002”, “DEF”)、 (“004”, “GHI”) が応答される。この例では、応答としてバリューと、キーとの両方が返されている。

【0007】

キーバリューストア (KVS) の典型的な形態として、ハッシュテーブルベースのものが挙げられる。ハッシュテーブルベースのKVSは、キーに対応するバリューのストレージ上のアドレス (以下、KVアドレス) を、ハッシュテーブルで管理する。典型的には、キーを入力とする所定のハッシュ関数を計算し、ハッシュテーブル上のそのハッシュ値 (通常、キーよりも十分短いバイト列である) に基づくアドレス (以下、ハッシュアドレス) のエントリに、バリューが記憶されているストレージ上のKVアドレスを格納しておく。なお、ハッシュ値から、ハッシュテーブル上のハッシュアドレスが一意に決まる限り、ハッシュ値とハッシュアドレスとが必ずしも一致する必要はないが、以下では説明のためハッシュ値とハッシュアドレスが一致する場合を想定する。このようにすることによって、GET要求 (バリューの読み出し) などを受けた場合、要求で指定されたキーに対して、ハッシュ関数を用いて瞬時にバリューが格納されているストレージ上のKVアドレスを特定することができる。よって、高速にバリューにアクセスして、応答を返すことができる。

30

40

【0008】

しかしながら、ハッシュテーブルベースのKVSは、上記のGETNEXT要求/GETPREV要求/GETRANGE要求に対する応答を返すことが出来ない。なぜなら、ハッシュテーブルベースのKVSは、指定されたキーのバリューを読み出すことはできるものの、指定されたキーの前後にどのようなキーが格納されているかを知る事が出来ないためである。各キーはハッシュ関数を介してハッシュ値に変換されて、ハッシュテーブル上のハッシュアドレスとして扱われる。ハッシュテーブル上のハッシュアドレスの前後関係は、実際のキーの前後関係とは異なるものであり、ハッシュテーブルからは実際のキーの前後関係を辿ることができない。よって、所定のキーの辞書順で、指定されたキーの次のキーを得るには、全検索を行うしかなく、それは現実的ではない。

50

【0009】

上記のようなGETNEXT要求/GETPREV要求/GETRANGE要求などに
 応答するため、キーを辞書順でアクセス可能なKVS(以下、順序型KVS)がある。順序型KVSは、ハッシュテーブルではなく、辞書順でアクセス可能なデータ構造でキーを
 管理し、GETNEXT要求/GETPREV要求/GETRANGE要求などに対応す
 るものである。オープンソース・ソフトウェアとしては、LevelDB、RocksDB
 などが有名である。このような順序型KVSは、B-treeやLSM-treeなど
 のデータ構造でキーを管理し、通常のGET要求などに対して、上記データ構造上でキー
 の探索を行う。探索の結果、キーが見つかった場合は、キーに対応するストレージ上のバ
 リューのKVアドレスを取得し、取得したKVアドレスに基づきストレージ上のバリュー
 にアクセスして、そのバリューを応答する。上記データ構造では、キーを辞書順で並べて
 管理しているため、GETNEXT要求を受けた場合も、通常のGET要求と同様に、ま
 ず指定されたキーを探索し、その後その次のキーを辿って、当該次のキーに対応するバ
 リューのKVアドレスにアクセスすることができる。よって、GETNEXT要求にも応
 答可能である。

10

【0010】

前述の順序型KVSは、GETNEXT要求等に応答できる機能を持つものの、ハッ
 シュテーブルベースのものとは異なり、キーを「探索」する必要があるという問題がある。応
 答速度を高めるために数々のデータベース技術によって改良が試みられているものの、順
 序型KVSでは、基本的に探索が不要なハッシュテーブルベースのKVSに比べて、キー
 を特定するまでの手順が複雑になるという問題がある。このため目的のバリューを検索す
 るまでの時間が長くなり、結果的に応答速度の悪化につながる。

20

【先行技術文献】

【特許文献】

【0011】

【特許文献1】米国出願公開第2011/0276744号明細書

【発明の概要】

【発明が解決しようとする課題】

【0012】

本発明の実施形態は、指定されたキーとある位置関係にあるキーに対するデータを高速
 に検索可能することを目的とする。

30

【課題を解決するための手段】

【0013】

本発明の実施形態としてのデータ処理装置は、要求解釈部と、第1アクセス部と、第2
 アクセス部と、第3アクセス部とを、備える。前記要求解釈部は、キーを指定した読み出
 し要求を受け付け、前記読み出し要求を解釈する。前記第1アクセス部は、複数のキーを
 所定のルールに従った順序で保持する順序キー列における前記キーに対して前記読み出
 し要求で指示された位置関係にあるキーを特定する。前記第2アクセス部は、第1アドレ
 スと第2アドレスとを対応づけて管理する内部または外部の第1記憶装置から、前記第1ア
 クセス部で特定されたキーのハッシュ値に基づく第1アドレスに対応する第2アドレスを
 取得する。前記第3アクセス部は、内部または外部の第2記憶装置において前記第2アク
 セス部で取得された前記第2アドレスに格納されたデータを読み出す。

40

【図面の簡単な説明】

【0014】

【図1】本発明の実施形態に係るデータ処理システムを示す図。

【図2】第1の実施形態に係るデータ処理装置のブロック図。

【図3】ハッシュテーブル、順序キー列、およびバリューを含むデータ群を示す図。

【図4】ハッシュテーブルの他の例を示す図。

【図5】第1の実施形態の動作例のフローチャート。

【図6】第1の実施形態の他の動作例のフローチャート。

50

【図7】第2の実施形態に係るデータ処理装置のブロック図。

【図8】順序キー列をバリューとして保存する様子を示す図。

【図9】第2の実施形態の動作例のフローチャート。

【図10】第2の実施形態の他の動作例のフローチャート。

【発明を実施するための形態】

【0015】

以下、図面を参照しながら、本発明の実施形態について説明する。

【0016】

図1は、本発明の実施形態に係るデータ処理システムを示す。データ処理システムは、データ処理装置11と、複数のクライアント装置21~23を備えている。

10

【0017】

データ処理装置11は、イーサネットなどのネットワーク越しに、クライアント装置21、22、23からの各種読み出しまたは書き込みの要求を受け付けるキーバリュースタ装置である。クライアント装置21~23は、例えばWebサーバ等のサーバであり、Webサーバの場合、一例として、URLをキー、Webページ等のコンテンツデータをバリューとして、キーとバリュー（キーバリュー）の書き込み要求をデータ処理装置11に行う。クライアント装置21~23は、インターネット等のネットワークを介してユーザ端末と接続されており、例えば、ユーザ端末からURLを指定した、コンテンツデータの書き込み指示を受ける。クライアント装置21~23は、この書き込み指示に基づき、URLをキー、コンテンツデータをバリューとして、キーバリューの書き込み要求を生成して、データ処理装置11に送る。データ処理装置11は、キーバリューの書き込み要求を受けると、後述する各実施形態に従ってキーバリューの書き込み処理を行う。書き込みが完了すると、書き込み完了の応答をクライアント装置21~23に返す。

20

【0018】

また、クライアント装置21~23は、ユーザ端末から、一例として、キーを指定した、バリューの読み出し指示を受けると、キーを指定した、キーバリューの読み出し要求を生成してデータ処理装置11に送る。データ処理装置11は、クライアント装置21~23からの読み出し要求を受けて、指定されたキーに基づき、後述する各実施形態に従って読み出し処理を行ってバリューを取得し、取得したバリューを含む応答をクライアント装置21~23に返す。以下、このようなデータ処理装置の各種の実施形態について説明する。

30

【0019】

(第1の実施形態)

図2は、本発明の第1の実施形態に係るデータ処理装置を示す。

【0020】

図2のデータ処理装置は、要求解釈部101、処理振分部102、ハッシュテーブルアクセス部103、順序キー列アクセス部104、KVアドレス割当部105、KVアクセス部106、応答生成部107、ハッシュテーブル記憶部108、順序キー列記憶部109、およびKV記憶部110を備える。順序キー列アクセス部104は一例として第1アクセス部に対応し、ハッシュテーブルアクセス部103は一例として第2アクセス部に対応し、KVアクセス部106は一例として第3アクセス部に対応する。このデータ処理装置は、典型的な形態としては、CPU、メモリ、ストレージ、ネットワークインタフェースなどで構成される。一部の機能がFPGAやASICなどの専用ハードウェアによって実現されていてもよい。

40

【0021】

要求解釈部101は、GET要求、GETNEXT要求、GETPREV要求、または、GETRANGE要求またはPUT要求などの要求を受け付け、受信した要求を解釈する。GET要求は、指定したキーに対応するバリューを読み出すことの要求である。GETNEXT要求は、指定したキーの辞書順で、次のキーに対応するバリューを読み出すことの要求である。GETPREV要求は、指定したキーの辞書順で、前のキーに対応する

50

バリューを読み出すことの要求である。GETRANGE要求は、2つのキーを指定し、当該指定した2つのキーの間に辞書順で存在する全てのキーに対応するバリューを読み出すことの要求である。PUT要求は、キーとバリューとを指定し、キーに対応づけてバリューを書き込むことの要求である。

【0022】

KV記憶部110は、キーとバリューとを含むデータを記憶している。キー毎に対応して、キーバリューが格納されている。図3の右下において複数のキーに対応する複数のデータが格納されている様子が示される。例えば、KV(Bruce)は、キーであるBruceと、当該キーに対応するバリューとを含むデータを表している。図3では、複数のデータ(キーバリュー)が、連続したアドレスに書き込まれている様子が示されるが、これは一例であり、どのように格納するかは、書き込むKVアドレスを割り当てる方式等に依存する。本例では、データは、キーとバリュー(すなわちキーバリュー)を含むが、データが、バリューを含み、キーは含まない構成もあり得る。またデータが、キーとバリュー以外の情報(例えばバリューのサイズ情報など)を含んでもかまわない。以下では、データが、少なくともバリューを含むことを前提として、説明を続ける。KV記憶部110は、ハードディスクまたはSSD、あるいは任意の不揮発性メモリ、あるいはその他の任意の永続的にデータを記憶する装置で構成されることができる。

10

【0023】

ハッシュテーブル記憶部108は、ハッシュテーブルを保持する。ハッシュテーブルとは、キーを所定のハッシュ関数にかけて算出されたハッシュ値に基づくアドレス(以下、ハッシュアドレス)のエントリに、そのキーに関わる情報、特にバリューを含むデータのKV記憶部110上のKVアドレス(以下、KVアドレス)等を格納するテーブルである。ハッシュアドレスは、ハッシュ値そのものでもよいし、ハッシュ値から一意に特定される別の値でもかまわない。本例では、ハッシュアドレスが、ハッシュ値そのものである場合を想定する。図3の左にハッシュテーブルの例を示す。H1、H5等はハッシュ値を表している。ADR(Bruce)は、キーであるBruceに対応するデータが格納された記憶部上のKVアドレスを表す。ハッシュテーブル記憶部108は、SSDまたはDRAMなど、任意の記憶媒体により構成でき、メモリで構成する場合、不揮発性メモリでも揮発性メモリでもよい。

20

【0024】

順序キー列記憶部109は、順序キー列を管理する。順序キー列は、複数のキーを所定のルールに従った順序で保持したものである。例えば、順序キー列は、B-treeまたはLSM-treeなどのデータ構造であり、複数のキーを辞書順の並びで保持する。図3の右上に順序キー列の例を示す。図3の例では、Bob、Bruce、Jones、Kenがこの順序で並んでいる。このデータ構造は、ハッシュテーブルと異なり、指定されたキーの前後のキーに容易にアクセスできるという特長を持つ。順序キー列は、単純にキーの並びだけ表してもよいし、キーのみならず、そのキーに関わる情報、特にバリューのKVアドレスをキーに関連付けて含んでも良い。以降では簡単のため、KVアドレスをキーに関連付けて含まない例を用いてKVアドレスを得るためにハッシュテーブルにアクセスする実施の形態を説明する。しかしながら、順序キー列にアクセスした後にハッシュテーブルにアクセスしてKVアドレスを得る代わりに、順序キー列内にKVアドレスをキーに関連付けて含むことで、その後のハッシュテーブルアクセスを省略する形態でもよい。順序キー列の構成は、具体的には様々な形態があり得るが、そのアルゴリズム自体は本実施形態とは無関係であり、本実施形態としては何でも良い。なお、順序キー列の管理には、一例としてファイルシステムが利用され、個々の部分データは個々のファイルとして実装される。順序キー列記憶部109は、HDD、SSDまたはDRAMなど、任意の記憶媒体により構成でき、メモリで構成する場合、不揮発性メモリでも揮発性メモリでもよい。

30

40

【0025】

なお、図3では、KV記憶部110に格納されているデータに対し矢印付きの線が引かれている。当該線は、ハッシュテーブルのエントリに格納されているKVアドレスと、順

50

序キー列のキー（エントリ）とからそれぞれ引かれている。これは、ハッシュテーブルのエントリに格納されたK Vアドレスが指す領域に、当該エントリに対応するキーに対するデータが格納されていることを表現している。

【0026】

K V記憶部110、ハッシュテーブル記憶部108および順序キー列記憶部109はそれぞれ別々のハードウェアでもよいし、共通のハードウェアでもよい。またはこれらの内の一部が同じハードウェアで、残りが別のハードウェアでもよい。

【0027】

処理振分部102は、要求解釈部101で解釈された要求の内容に応じて、ハッシュテーブルアクセス部103または順序キー列アクセス部104またはこれらの両方に、処理を依頼する。

10

【0028】

ハッシュテーブルアクセス部103は、指定されたキーのハッシュ値を計算し、当該ハッシュ値に対応するハッシュアドレスに基づき、ハッシュテーブルのエントリへアクセスし、エントリに格納されているK V記憶部110上のK Vアドレスを取得する。

【0029】

K Vアドレス割当部105は、K V記憶部110に新規にデータ（バリューまたはキーバリューなど）を書き込むときのK Vアドレスを割り当てる。K Vアドレスの割り当て方法は、ログ形式でPUT要求が来た順番にシーケンシャルにK Vアドレスを割り当てて、その順番で書き込みを行う方式が一般的であるが、これに限定されるものではない。

20

【0030】

K Vアクセス部106は、K V記憶部110にアクセスして、ハッシュテーブルアクセス部103（または順序キー列アクセス部104）によって取得されたK Vアドレスに格納されているデータを取得する。

【0031】

応答生成部107は、K Vアクセス部106により取得されたデータ（バリューまたはキーバリューなど）に基づき応答を生成して、生成した応答を、要求解釈部101で受けた要求の送信元に送信する。

【0032】

ここでハッシュテーブルについて補足説明を行う。前述したハッシュテーブルの構成方法には、様々な形態があり得る。特に異なるキーが同じハッシュ値をとってしまうという、いわゆる衝突が発生した場合に、どのようにその同じハッシュ値を持つ複数のキーのエントリを格納するかについては、様々なアルゴリズムがある。

30

【0033】

例えば、衝突した場合には、所定のルールに基づいて他の空いているエントリが見つかるまで探索する開番地法、一つのエントリに、例えばリスト構造などで複数のエントリを格納する連鎖法などがある。別の例として、指定されたキーと、対応するエントリが一致するかどうかを確認する方法がある。この方法の場合、ハッシュテーブルの各エントリに、K Vアドレスに加えて、キーそのものを格納して、指定されたキーと、エントリに格納されたキーの一致性を確認してもよい。

40

【0034】

ハッシュテーブルの容量を減らすために、ハッシュテーブルの各エントリに、キーを所定の関数で縮退させたビット列をK Vアドレスに加えて格納してもよい。例えばキーであるBruceを関数Sig()で縮退させてビット列Sig(Bruce)を、ADR(Bruce)に加えてエントリに格納する。縮退させたビット列は、縮退キーとも呼ぶ。このハッシュテーブルの例を図4に示す。この場合は、最終的に、指定されたキーとの完全マッチングを確認するために、K V記憶部110上にキーバリューを格納し、そのキーバリューにアクセスして、そこに含まれるキーとの最終的なマッチングを図る必要がある。すなわち、まずは、指定されたキーを縮退させてビット列を取得し、そのビット列と、ハッシュテーブルの該当するエントリに格納されている、縮退させたビット列とのマッチ

50

ングを確認する。この段階でマッチングしていなければ、マッチしていない（つまり当該エントリに含まれるK Vアドレスに格納されているバリューは、指定されたキーに対応するバリューではない）と判断し、エントリの探索を再開する。マッチしていれば、マッチしている可能性は高いが、完全にはマッチすることはまだ確認されていないため、そのエントリに格納されているK V記憶部110上のK Vアドレスを元にK V記憶部110にアクセスしてキーバリューを読み出す。そして、そこに含まれているキーと、指定されたキーとのマッチングを確認する。マッチしたら、そのときに読み出したバリューをそのまま採用し、マッチしなかった場合は再びハッシュテーブルに戻り、エントリの探索を再開し、同様にマッチするキーが見つかるか、あるいはマッチするキーが無いことが確認できるまで繰り返す。

10

【0035】

このようなハッシュテーブルの構造、および探索アルゴリズム自体は、本実施形態の本質とは無関係であり、本実施形態では何でもよい。以下の説明では、K V記憶部110上に少なくともバリューを含むデータを格納することを仮定するが、本実施形態は、これに限定するものではない。

【0036】

本実施形態は、要求解釈部101が受けた要求の種類によって、ハッシュテーブルと順序キー列とを使い分けるところを特徴の1つとする。以下、GET要求を受けたとき、GETNEXT要求/GETPREV要求/GETRANGE要求を受けたとき、および、PUT要求を受けたときの動作について、それぞれ説明する。

20

【0037】

まず、キーを指定したGET要求を受けたときの動作を説明する。要求解釈部101が、その要求がGETであることを解釈し、処理振分部102は、キーをハッシュテーブルアクセス部103に渡す。ハッシュテーブルアクセス部103は、所定のハッシュ関数を用いて、指定されたキーのテーブルエントリにアクセスし、K Vアドレスを取得する。例えば指定されたキーのハッシュ値を計算し、ハッシュ値に対応するハッシュテーブル上のエントリを特定する。そして、特定したエントリに格納されているK Vアドレスに基づきK Vアクセス部106が、K V記憶部110にアクセスして、K V記憶部110に記憶されているデータ（キーバリュー）を読み出す。読み出したキーバリューに含まれるキーが、指定されたキーにマッチすることを確認した上で、応答生成部107が、読み出したバリューを含む応答を生成して、GET要求の送信元に送信する。読み出したキーバリューに含まれるキーが、指定されたキーにマッチするかの判断は、K Vアクセス部106が行っても良いし、処理振分部102が行っても良いし、当該判断を行う処理部を独立して設けてもよい。指定されたキーにマッチしない場合は、マッチしなかった場合は再びハッシュテーブルに戻り、エントリの探索を再開し、同様にマッチするキーが見つかるか、あるいはマッチするキーが無いことが確認できるまで繰り返す。

30

【0038】

次に、GETNEXT要求、GETPREV要求、または、GETRANGE要求を受けた場合の動作を説明する。GETNEXT要求は、指定したキーの辞書順で、次のキーに対応するバリューを読み出すことの要求である。GETPREV要求は、指定したキーの辞書順で、前のキーに対応するバリューを読み出すことの要求である。GETRANGE要求は、2つのキーを指定し、当該指定した2つのキーの間に存在する全てのキー（当該範囲を規定する2つのキーを含むように定義しても、含まないように定義してもよい）に対応するバリューを読み出すことの要求である。

40

【0039】

処理振分部102が、GETNEXT要求、GETPREV要求、または、GETRANGE要求を受けたことを解釈し、順序キー列アクセス部104に、指定されたキーと、要求の内容を通知する。

【0040】

順序キー列アクセス部104は、GETNEXT要求であれば、指定されたキーの次の

50

キー、GETPREV要求であれば、指定されたキーの前のキー、GETRANGE要求であれば、指定された2つのキーの間にある全てのキーを、順序キー列記憶部109内の順序キー列に基づき、特定する。なお、前述したように、順序キー列において、キーと関連づけて、バリューが格納されているKV記憶部110上のKVアドレスを保持してもよく、この場合は、特定したキーに対応するバリューのKV記憶部110上のKVアドレスも特定してもよい。

【0041】

順序キー列アクセス部104は、特定したキーを、処理振分部102に渡し、処理振分部102は、当該特定したキーをハッシュテーブルアクセス部103に渡す。ハッシュテーブルアクセス部103は、ハッシュテーブルを用いて、当該特定したキーに対応するバリューのKV記憶部110上のKVアドレスを特定する。なお、順序キー列に関連づけられたKVアドレスを特定した場合は、順序キー列アクセス部104は、特定したKVアドレスをKVアクセス部106に渡せばよい。

10

【0042】

この後の処理は、GET要求の場合と同様である。すなわち、KVアクセス部106が、KV記憶部110にアクセスしてバリュー（もしくはキーバリューなど）を読み出し、応答生成部107が、読み出したバリュー（もしくはキーバリューなど）を含めた応答を生成して、要求の送信元に送信する。

【0043】

次に、PUT要求を受けた場合の動作を説明する。PUT要求は、キーとバリューとを指定し、キーに対応づけてバリューを書き込むことの要求である。

20

【0044】

処理振分部102が、PUT要求を受けたことを解釈し、ハッシュテーブルアクセス部103にキーとバリューとを渡し、順序キー列アクセス部104にキーを渡す。ハッシュテーブルアクセス部103は、指定されたキーのハッシュ値を計算し、ハッシュテーブル上の当該ハッシュアドレスに対応するエントリを生成する。KVアドレス割当部105は、指定されたキーとバリュー（少なくともバリュー）とを含むデータを格納するKV記憶部110上のKVアドレスを、当該キーに対して割り当てる。ハッシュテーブルアクセス部103は、生成したエントリに、KVアドレス割当部105が決定したKV記憶部110上のKVアドレスを格納する。

30

【0045】

また順序キー列アクセス部104は、順序キー列の構造に従って、指定されたキーのエントリを順序キー列に追加する。キーに関連づけて、当該KVアドレス割当部105が割り当てたKVアドレスを格納してもよい。このKVアドレスは、ハッシュテーブルに生成したエントリに設定したKVアドレスと同じである。

【0046】

なお、KVアドレス割当部105は、ハッシュテーブルアクセス部103または順序キー列アクセス部104の少なくとも一方の要求に応じて、KVアドレスを割り当ててもよい。あるいは、処理振分部102からKVアドレスの割り当て要求を受けて、指定されたバリューを含むデータを格納するKV記憶部110上のKVアドレスを割り当ててもよい。

40

【0047】

KVアクセス部106は、KVアドレス割当部105で割り当てられたKV記憶部110上のKVアドレスに、指定されたバリューを含むデータ（ここではキーバリュー）の書き込みを行う。

【0048】

応答生成部107は、KVアクセス部106による書き込みが完了すると、PUT要求を実行できた旨の応答を生成して、PUT要求の送信元に送信する。

【0049】

ここで、指定されたキーと同じキーに対応するキーバリューが既にKV記憶部110さ

50

れているかどうかを事前に確認する必要がある場合は、当該キーにマッチするエントリがハッシュテーブルに無いことを確認できた場合、または、K Vアクセス部 1 0 6 でキーを読み出した時点で、マッチするエントリが無いことを確認できた場合に、K Vアドレス割当部 1 0 5 は、K V記憶部 1 1 0 上に新規に格納するデータ（キーまたはキーバリューなど）のK Vアドレスの割り当てを行えばよい。指定されたキーに対応するデータが既に保持されていた場合は、既にハッシュテーブルと順序キー列に該当エントリがあることになるので、それぞれに既に存在しているエントリのK Vアドレスを、新しく割り当てられたK Vアドレスで上書きすればよい。

【 0 0 5 0 】

なお、指定されたキーと同じキーに対応するデータ（バリューまたはキーバリューなど）を既に登録しているかどうかの確認は、上述ではハッシュテーブルアクセス部 1 0 3 経由で行う例を示したが、順序キー列アクセス部 1 0 4 を用いてもよい。すなわち、順序キー列に、指定されたキーと同じキーが存在するかを確認すればよい。なお、ハッシュテーブルアクセス部 1 0 3 経由で行った方が、処理が高速になるという利点がある。

10

【 0 0 5 1 】

図 5 に、本実施形態の動作例のフローチャートを示す。ここでは読み出しに係る要求として、G E T N E X T 要求を受け付けた場合の動作例を示す。

【 0 0 5 2 】

まず、要求解釈部 1 0 1 がクライアント装置（図 1 参照）から、ネットワークを介して G E T N E X T 要求を受け付ける（S 1 0 1）。G E T N E X T 要求は、インターネット

20

【 0 0 5 3 】

処理振分部 1 0 2 は、G E T N E X T 要求で指定されたキーと、要求の処理内容とを順序キー列アクセス部 1 0 4 に通知する。順序キー列アクセス部 1 0 4 は、順序キー列に基づき、指定されたキーの次のキー（指定されたキーに関連するキー）を特定し、処理振分部 1 0 2 に渡す（S 1 0 2）。

【 0 0 5 4 】

処理振分部 1 0 2 は、当該キーをハッシュテーブルアクセス部 1 0 3 に渡し、ハッシュテーブルアクセス部 1 0 3 は、当該キーのハッシュ値を計算する。そして、ハッシュテーブルにおいて当該ハッシュ値に基づくハッシュアドレスに対応するエントリに格納されているK Vアドレスを取得する（S 1 0 3）。

30

【 0 0 5 5 】

K Vアクセス部 1 0 6 は、K V記憶部 1 1 0 において当該K Vアドレスに格納されているデータ（キーバリューなど）を読み出し、当該データからバリューを取り出す（S 1 0 4）。

【 0 0 5 6 】

応答生成部 1 0 7 は、K Vアクセス部 1 0 6 により取得されたバリューを含む応答を生成し、当該応答をG E T N E X T 要求の送信元に送信する（S 1 0 5）。

【 0 0 5 7 】

ここでは、G E T N E X T 要求を受け付けた場合の動作を説明したが、G E T P R E V 要求、または、G E T R A N G E 要求などを受けた場合も同様に行われる。

40

【 0 0 5 8 】

図 6 に、本実施形態の他の動作例のフローチャートを示す。ここでは、書き込みに係る要求として、P U T 要求を受け付けた場合の動作例を示す。

【 0 0 5 9 】

まず、要求解釈部 1 0 1 がクライアント装置から、ネットワークを介してP U T 要求を受け付ける（S 2 0 1）。

【 0 0 6 0 】

処理振分部 1 0 2 は、P U T 要求で指定されたキーとバリューとをハッシュテーブルアクセス部 1 0 3 に渡し、ハッシュテーブルアクセス部 1 0 3 は、当該キーのハッシュ値を

50

計算する。またK Vアドレス割当部1 0 5にK Vアドレスの割り当てを要求してK Vアドレスの割り当てを受け(S 2 0 2)、ハッシュテーブルにおいて当該ハッシュ値に基づくハッシュアドレスに対応するエントリにK Vアドレスを格納する(S 2 0 3)。

【0 0 6 1】

処理振分部1 0 2は、P U T要求で指定されたキーを順序キー列アクセス部1 0 4に渡す。順序キー列アクセス部1 0 4は、順序キー列に、指定されたキーを追加する(S 2 0 4)。

【0 0 6 2】

K Vアクセス部1 0 6は、K Vアドレス割当部1 0 5で割り当てられたK V記憶部1 1 0上のK Vアドレスに、指定されたバリューを含むデータ(キーバリューなど)の書き込みを行う(S 2 0 5)。

【0 0 6 3】

応答生成部1 0 7は、K Vアクセス部1 0 6による書き込みが完了すると、P U T要求を実行できた旨の応答を生成して、P U T要求の送信元に送信する(S 2 0 6)。

【0 0 6 4】

以上、本実施形態によれば、G E T要求に対してはハッシュテーブル方式並みの高速な検索を行って応答を返しつつ、G E T N E X T要求/G E T P R E V要求/G E T R A N G E要求などのキー順序を考慮する必要がある要求にも、高速に検索して応答することが可能になる。

【0 0 6 5】

(第2の実施形態)

図7に、本発明の第2の実施形態に係るデータ処理装置を示す。図7のデータ処理装置は、第1の実施形態のデータ処理装置を拡張したものである。第1の実施形態と異なり、内部キー生成部1 1 1が追加され、順序キー列記憶部1 0 9が除去されている。第2の実施形態では、順序キー列をバリューとして、K V記憶部1 1 0で管理するところが異なる。以下、第1の実施形態との差分を中心に説明し、重複する構成および動作の説明は適宜省略する。

【0 0 6 6】

第1の実施形態では、順序キー列は順序キー列記憶部1 0 9(図2参照)で管理され、一例としてファイルシステムで管理されることを述べた。順序キー列は、全てのキーを辞書順に並べて記憶しているため、それだけで容量が大きくなるため、一般に複数の部分に分割されている。個々の部分は全体の部分集合となり、それ自体でキーを辞書順に並べて管理している。当然ながら既存のキーの並び間に、新たなキーを挿入したり、あるいはキーの並びの中の途中のキーを削除したりする必要がある。このことから、キーの位置や長さが柔軟に変更可能である必要があり、また、部分自体を新規に生成したり削除したりすることも容易にできる必要がある。これには上述したファイルシステムを利用することで対応でき、個々の部分は個々のファイルとして実装される。これにより実装が容易となるが、一方でファイルシステムを用いると、汎用で複雑な機構を設ける必要があり、構成全体として複雑になり、性能的にもあまり効率が良くなる可能性がある。

【0 0 6 7】

そこで本実施形態では、ファイルシステムを使わずに、順序キー列を、内部的なキーバリューとしてK V記憶部1 1 0で管理する構成をとる。

【0 0 6 8】

順序キー列アクセス部1 0 4は、P U T要求を受けた場合、P U T要求で指定されたキーを追加した新たな順序キー列を生成(詳細は後述)する。内部キー生成部1 1 1は、予め定めた方法で内部キーを生成する。一例として内部キーとして、常に同じ値を生成することがある。以下の説明では、内部キー生成部1 1 1は、常に同じ値の内部キーを生成するとする。処理振分部1 0 2は、内部キーと新たな順序キー列とをハッシュテーブルアクセス部1 0 3に渡す。ハッシュテーブルアクセス部1 0 3は、内部キーからハッシュ値を計算し、ハッシュテーブルに当該ハッシュ値に基づくハッシュアドレスに対応するエント

10

20

30

40

50

りを生成し、生成したエントリに、新たな順序キー列を格納するためのKV記憶部110上のKVアドレスを格納する。当該KVアドレスは、KVアドレス割当部105に要求するなどして割り当てを受ければよい。KVアクセス部106は、新たな順序キー列を、通常のバリュー（データ）として、KV記憶部110上のKVアドレスに格納する。

【0069】

具体例として、例えばデータ処理装置を立ち上げて、最初のPUT (Key = "John", Value = "GHI") 要求を受けたとする。このとき、指定されたキーとバリューが初めて、この装置に格納されることになるので、順序キー列としては、この指定されたキー (John) のみが存在することになる。よって、このキーを最初の順序キー列として、KV記憶部110に、内部キーに対応するデータ (バリュー) として格納するべく、例えば、内部キー生成部111は、Key = "key_seq0" を生成し、Value = "John" とする。これらを処理振分部102に渡す。処理振分部102は、内部キーと "John" とをハッシュテーブルアクセス部103に渡す。ハッシュテーブルアクセス部103は、この内部キー "key_seq0" のハッシュ値を計算し、ハッシュテーブルにおける当該ハッシュ値のエントリに、"John" が格納されるKVアドレスを格納する。KVアドレスは、KVアドレス割当部105が割り当てる。KVアクセス部106は、"John" を、通常のバリューとしてKV記憶部110に格納する。なお、PUT要求そのものに対する処理も、第1の実施形態と同様にして行う。すなわち、"John" のハッシュ値に対応するエントリに、"GHI" を格納するKV記憶部110上のKVアドレスを格納し、KV記憶部110上の当該KVアドレスには、"GHI" を含むデータを格納する。どちらの処理を先にするかは任意である。

【0070】

さらに、PUT (Key = "Bob", Value = "JKL") 要求を受けた場合、まず順序キー列アクセス部104は、現在の順序キー列の読み出しを行う。すなわち、内部キー生成部111に内部キーの生成を依頼する。内部キー生成部111は、Key = "key_seq0" を生成し、処理振分部102に渡す。処理振分部102は、これをハッシュテーブルアクセス部103に送る。ハッシュテーブルアクセス部103は、"key_seq0" のハッシュ値を計算し、ハッシュ値に対応するエントリに格納されたKVアドレスを取得する。KVアクセス部106は、KV記憶部110上の当該KVアドレスからデータ、すなわち、内部キー "key_seq0" のバリューである "John" を読み出す。KVアクセス部106は、これを順序キー列アクセス部104に渡す。これにより、順序キー列アクセス部104は、現在の順序キー列 "John" を得る。そして、今回のPUT要求で指定されたKey = "Bob" と、現在の順序キー列 "John" を比較し、"Bob" を "John" の前に配置することにより、新たな順序キー列として "Bob John" を生成する。これを新たなバリューとし、また内部キーは固定の値である "key_seq0" とする。そして、前回同様に、ハッシュテーブルに、当該内部キーのハッシュ値に対応するエントリに、"Bob John" を格納するためのKV記憶部110上のKVアドレスを上書きし、"Bob John" をKV記憶部110上のKVアドレスに格納する。当該アドレスはKVアドレス割当部105により割り当てを受ける。なお、PUT要求そのものに対する処理も行う。すなわち、"Bob" のハッシュ値に対応するエントリに、"JKL" を格納するKV記憶部110上のKVアドレスを格納し、KV記憶部110上の当該KVアドレスには、"JKL" を含むデータを格納する。どちらの処理を先にするかは任意である。

【0071】

この後、さらにPUT (Key = "Ken", Value = "MNO") 要求が来た場合、同様に、内部キーを "key_seq0" として、KV記憶部110からバリュー（データ） "Bob John Ken" を現在の順序キー列として読み出す。これに "Ken" を追加して、新たな順序キー列 "Bob John Ken" を得る。ハッシュテーブルおよびKV記憶部110を前回と同様にして更新する。なお、PUT要求そのものに対する処理も前回と同様にして行う。

10

20

30

40

50

【 0 0 7 2 】

さらに、PUT (Key = “ Bruce ”、 Value = “ P Q R ”) 要求が来た場合、同様に、内部キーを “ key __ seq 0 ” として生成し、KV記憶部 1 1 0 からバリュウー(データ) “ Bob John Ken ” を現在の順序キー列として読み出す。これに “ Bruce ” を追加して、新たな順序キー列 “ Bob Bruce John Ken ” を得る。ハッシュテーブルおよびKV記憶部 1 1 0 を前回と同様にして更新する。なお、PUT要求そのものに対する処理も前回と同様にして行う。

【 0 0 7 3 】

本実施形態におけるKV記憶部 1 1 0 のデータ構造を図 8 に示す。順序キー列を含むデータ(バリュウー)が、通常のキーに対応するバリュウーと同じ位置づけで、KV記憶部 1 1 0 に記憶されている。

【 0 0 7 4 】

次に、GETNEXT (Key = “ John ”) 要求を受けた場合、順序キー列アクセス部 1 0 4 は、前述同様に、まず現在の順序キー列を読み出す。すなわち、内部キー生成部 1 1 1 に内部キーの生成を要求し、内部キー生成部 1 1 1 は、内部キーとして固定の値である “ key __ seq 0 ” を生成して、処理振分部 1 0 2 に渡す。処理振分部 1 0 2 は、これをハッシュテーブルアクセス部 1 0 3 に渡して、内部キーに対応するバリュウー(データ)が格納されているKVアドレスをハッシュテーブルから取得する。KVアクセス部 1 0 6 は、KV記憶部 1 1 0 の当該取得したKVアドレスから、内部キー “ key __ seq 0 ” のバリュウーである “ Bob Bruce John Ken ” を読み出す。これを順序キー列アクセス部 1 0 4 に渡し、これにより順序キー列アクセス部 1 0 4 は現在の順序キー列を得る。この順序キー列の情報から、GETNEXTで指定されたキーである Johnの次のキーはKenであることを確認する。キー “ Ken ” を処理振分部 1 0 2 に渡し、処理振分部 1 0 2 は、これをハッシュテーブルアクセス部 1 0 3 に渡す。ハッシュテーブルアクセス部 1 0 3 は、 “ Ken ” のハッシュ値に対応するエントリからKVアドレスを取得し、KVアクセス部 1 0 6 が、KV記憶部 1 1 0 における当該KVアドレスからデータを読み出し、読み出したデータから、キー “ Ken ” に対応するバリュウーを取り出す。応答生成部 1 0 7 が、取り出したバリュウーを含む応答を生成して、GETNEXT要求の送信元に送信する。

【 0 0 7 5 】

なお、順序キー列はPUT要求が来る毎にサイズが増大していくので、適宜分割して、それぞれに異なる内部キーを割り当ててもよい。例えば内部キー key __ seq 0 に対応する順序キー列を 2 つに分割し、それぞれ、内部キーとして key __ seq 1、key __ seq 2 を割り当ててもよい。分割して管理するアルゴリズムは、LSM - tree など従来の様々な方式に従ってよく、本実施形態はその内容を問わない。この場合、内部キー生成部 1 1 1 は、GET要求等の読み出しの要求を受けた場合は、当該アルゴリズムに基づき内部キーを特定すればよい。単純には、分割の際、アルファベットの範囲で分けてもよい。例えば A ~ L、M ~ Z の 2 つのグループに分けて、内部キーとして key __ seq 1、key __ seq 2 を割り当てる。そして、読み出し要求で指定されたキーの先頭のアルファベットが A ~ L のときは key __ seq 1、M ~ Z のときは key __ seq 2 と決定してもよい。ここで述べた例は実現の可能性を示すための例に過ぎず、他にも様々な方法で実施可能である。

【 0 0 7 6 】

また、内部キー生成部 1 1 1 が生成する内部キーは、GET要求やPUT要求など、要求解析部で受け付ける要求で指定されるキー(外部キーと呼ぶ)と同じ値になる可能性もある。この問題を解消するため、外部キーと内部キーを明確に区別する情報を別途、外部キーおよび内部キーに付加してもよい。例えば、外部キーには例えば常に先頭に 0 を付け、内部キーには常に先頭に 1 を付ける。すなわち、外部キーが “ John ” であれば、先頭に 0 を付けて “ 0 John ” とし、内部キーが “ key __ seq 0 ” であれば、先頭に 1 を付けて、 “ 1 key __ seq 0 ” とする。このように、内部キーおよび外部キー

10

20

30

40

50

に、特定の識別子を付けることで、内部キーと外部キーとを明確に区別できる。識別子の形式や付け方はどのような形態でもよい。

【 0 0 7 7 】

図 9 に、本実施形態の動作例のフローチャートを示す。ここでは、読み出しに係る要求として、GETNEXT要求を受け付けた場合の動作例を示す。

【 0 0 7 8 】

まず、要求解釈部 1 0 1 がクライアント装置から、ネットワークを介してGETNEXT要求を受け付ける (S 3 0 1) 。

【 0 0 7 9 】

処理振分部 1 0 2 は、GETNEXT要求で指定されたキーを順序キー列アクセス部 1 0 4 に渡す。順序キー列アクセス部 1 0 4 は、内部キー生成部 1 1 1 に内部キーの生成を依頼し、内部キー生成部 1 1 1 は、内部キーを生成して、処理振分部 1 0 2 に送る (S 3 0 2) 。

【 0 0 8 0 】

処理振分部 1 0 2 は、内部キーをハッシュテーブルアクセス部 1 0 3 に渡し、ハッシュテーブルアクセス部 1 0 3 は、当該内部キーのハッシュ値を計算する。そして、ハッシュテーブルにおいて当該ハッシュ値に基づくハッシュアドレスに対応するエントリに格納されているKVアドレスを取得する (S 3 0 3) 。

【 0 0 8 1 】

KVアクセス部 1 0 6 は、KV記憶部 1 1 0 において当該KVアドレスに格納されているデータを読み出し、当該データからバリューである順序キー列を取り出す (S 3 0 4) 。

【 0 0 8 2 】

順序キー列アクセス部 1 0 4 は、当該順序キー列において、指定されたキーの次のキーを特定し、特定したキーを処理振分部 1 0 2 に送る (S 3 0 5) 。

【 0 0 8 3 】

処理振分部 1 0 2 は、当該キーをハッシュテーブルアクセス部 1 0 3 に渡し、ハッシュテーブルアクセス部 1 0 3 は、当該キーのハッシュ値を計算する。そして、ハッシュテーブルにおいて当該ハッシュ値に基づくハッシュアドレスに対応するエントリに格納されているKVアドレスを取得する (S 3 0 6) 。

【 0 0 8 4 】

KVアクセス部 1 0 6 は、KV記憶部 1 1 0 において当該KVアドレスに格納されているデータを読み出し、当該データからバリューを取り出す (S 3 0 7) 。

【 0 0 8 5 】

応答生成部 1 0 7 は、KVアクセス部 1 0 6 により取得されたバリューを含む応答を生成し、当該応答をGETNEXT要求の送信元に送信する (S 3 0 8) 。

【 0 0 8 6 】

ここでは、GETNEXT要求を受け受けた場合の動作を説明したが、GETPREV要求、または、GETRANGE要求などを受けた場合も同様に行われる。

【 0 0 8 7 】

図 1 0 に、本実施形態の他の動作例のフローチャートを示す。ここでは、書き込みに係る要求として、PUT要求を受け付けた場合の動作例を示す。

【 0 0 8 8 】

まず、要求解釈部 1 0 1 がクライアント装置から、ネットワークを介してPUT要求を受け付ける (S 4 0 1) 。

【 0 0 8 9 】

処理振分部 1 0 2 は、PUT要求で指定されたキーを順序キー列アクセス部 1 0 4 に渡す。順序キー列アクセス部 1 0 4 は、内部キー生成部 1 1 1 に内部キーの生成を依頼し、内部キー生成部 1 1 1 は、内部キーを生成して、処理振分部 1 0 2 に送る (S 4 0 2) 。

【 0 0 9 0 】

10

20

30

40

50

処理振分部 102 は、当該内部キーをハッシュテーブルアクセス部 103 に渡し、ハッシュテーブルアクセス部 103 は、当該内部キーのハッシュ値を計算する。ハッシュ値に対応するエントリから K V アドレスを取得し (S 4 0 3)、K V アクセス部 106 に渡す。K V アクセス部 106 は、K V 記憶部 110 における当該 K V アドレスからデータを読み出し、データから順序キー列を取得する (S 4 0 4)。

【 0 0 9 1 】

順序キー列アクセス部 104 は、取得された順序キー列に、P U T 要求で指定されたキーを追加して更新し、更新後の順序キー列をハッシュテーブルアクセス部 103 に渡す (S 4 0 5)。ハッシュテーブルアクセス部 103 は、K V アドレス割当部 105 に K V アドレスの割り当てを要求して、K V アドレスの割り当てを受ける (S 4 0 6)。ハッシュテーブルにおいて当該更新後の順序キー列のハッシュ値に基づくハッシュアドレスに対応するエントリに、K V アドレスを格納 (上書き) する (S 4 0 7)。K V アクセス部 106 は、K V 記憶部 110 における当該 K V アドレスに、更新後の順序キーを含むデータを書き込む (S 4 0 8)。

10

【 0 0 9 2 】

処理振分部 102 は、P U T 要求で指定されたキーとバリューとをハッシュテーブルアクセス部 103 に渡し、ハッシュテーブルアクセス部 103 は、当該指定されたキーのハッシュ値を計算する。ハッシュテーブルアクセス部 103 は、K V アドレス割当部 105 に K V アドレスの割り当てを要求して K V アドレスの割り当てを受け (S 4 0 9)、ハッシュテーブルにおいて当該ハッシュ値に基づくハッシュアドレスに対応するエントリに K V アドレスを格納する (S 4 1 0)。

20

【 0 0 9 3 】

K V アクセス部 106 は、K V 記憶部 110 上の当該 K V アドレスに、指定されたバリューを含むデータの書き込みを行う (S 4 1 1)。

【 0 0 9 4 】

応答生成部 107 は、K V アクセス部 106 による書き込みが完了すると、P U T 要求を実行できた旨の応答を生成して、P U T 要求の送信元に送信する (S 4 1 2)。

【 0 0 9 5 】

以上、本実施形態によれば、順序キー列をハッシュテーブルベースの K V S 上にバリューとして管理することができるので、順序キー列を管理するためにファイルシステムを使用する必要はなくなる。これにより、全体のデータ構造の簡略化を図ることができる。

30

【 0 0 9 6 】

尚、各実施形態のデータ処理装置は、例えば、汎用のコンピュータ装置を基本ハードウェアとして用いることでも実現可能である。すなわち、上記のコンピュータ装置に搭載されたプロセッサにプログラムを実行させることにより実現出来る。このとき、データ処理装置は、上記のプログラムをコンピュータ装置にあらかじめインストールすることで実現することや、各種の記憶媒体に記憶、あるいはネットワークを介して上記のプログラムを配布、このプログラムをコンピュータ装置に適宜インストールすることで実現出来る。また、データ処理装置内の各記憶部は、上記のコンピュータ装置に内蔵あるいは外付けされたメモリ、ハードディスクもしくは C D - R、C D - R W、D V D - R A M、D V D - R などの記憶媒体などを適宜利用して実現することができる。

40

【 0 0 9 7 】

本発明は上記実施形態そのままに限定されず、実施段階では要旨を逸脱しない範囲で構成要素を変形し具体化出来る。上記実施形態に開示されている複数の構成要素の適宜な組み合わせで、種々の発明を形成できる。例えば、実施形態に示される全構成要素から幾つかの構成要素を削除出来、異なる実施形態に渡る要素を適宜組み合わせることも出来る。

【 符号の説明 】

【 0 0 9 8 】

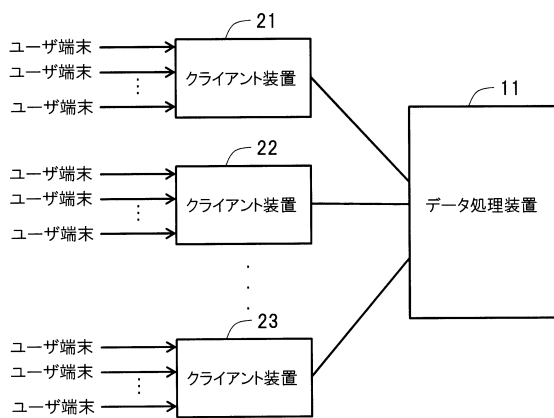
21 ~ 23 : クライアント装置

11 : データ処理装置

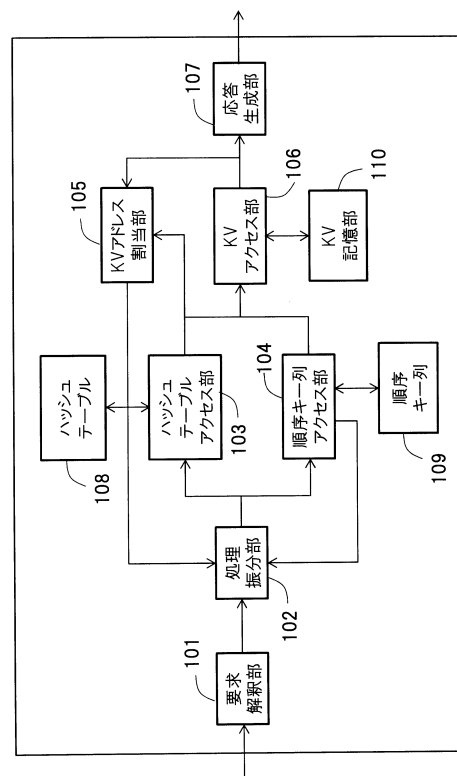
50

- 101 : 要求解釈部
- 102 : 処理振分部
- 103 ; ハッシュテーブルアクセス部
- 104 : 順序キー列アクセス部
- 105 : K V アドレス割当部
- 106 : K V アクセス部
- 107 : 応答生成部
- 108 : ハッシュテーブル記憶部
- 109 : 順序キー列記憶部
- 110 : K V 記憶部
- 111 : 内部キー生成部

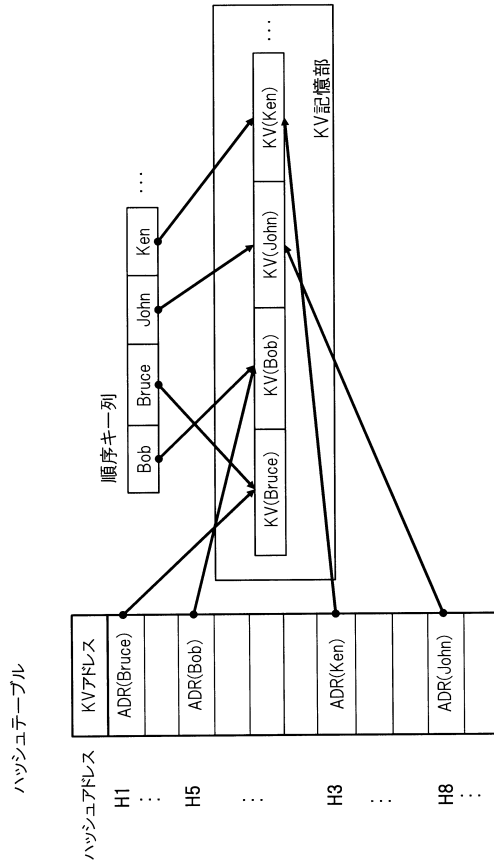
【図1】



【図2】



【図3】

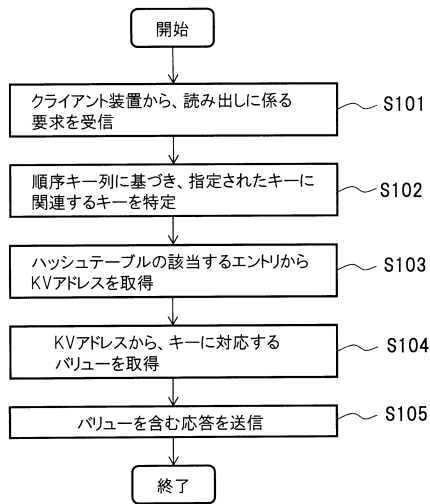


【図4】

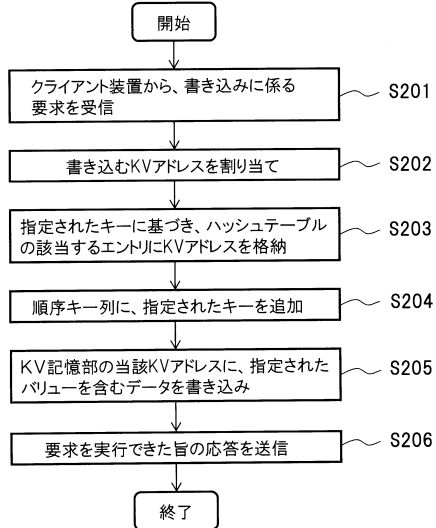
ハッシュテーブル

ハッシュアドレス	KVアドレス、縮退キー
H1	ADR(Bruce),Sig(Bruce)
...	
H5	ADR(Bob),Sig(Bob)
...	
H3	ADR(Ken),Sig(Ken)
...	
H8	ADR(John),Sig(John)
...	

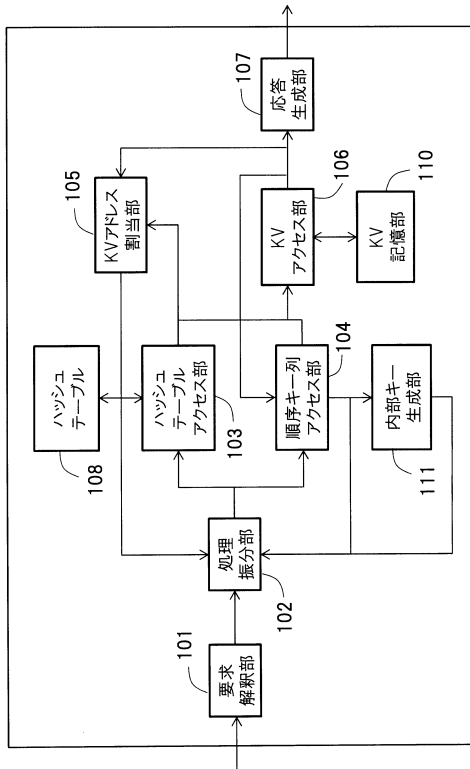
【図5】



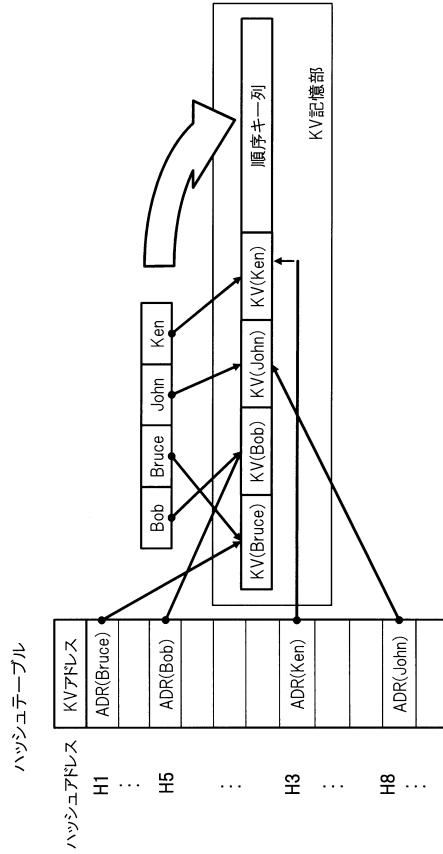
【図6】



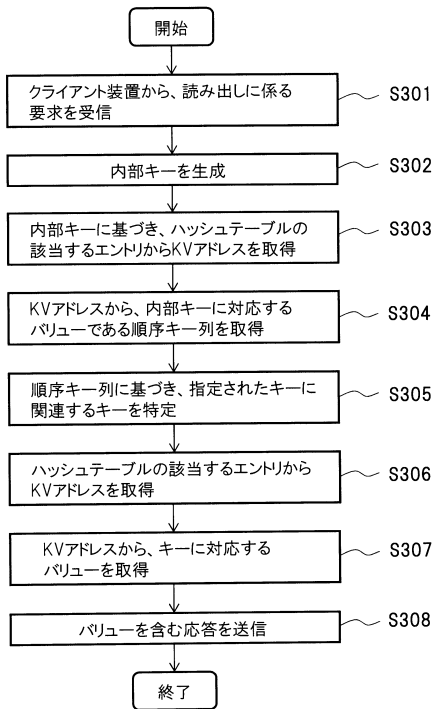
【図7】



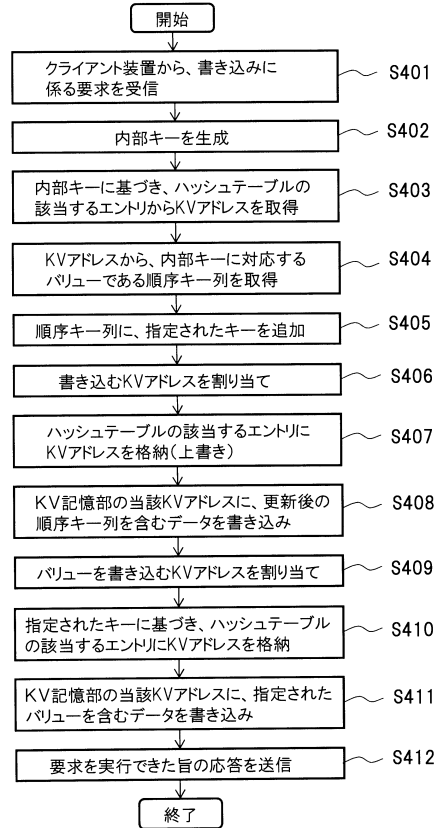
【図8】



【図9】



【図10】



フロントページの続き

(74)代理人 100118876

弁理士 鈴木 順生

(72)発明者 田中 信吾

東京都港区芝浦一丁目1番1号 株式会社東芝内

審査官 塚田 肇

(56)参考文献 特開2015-176407(JP,A)

国際公開第2014/141594(WO,A1)

特開2009-003541(JP,A)

特開2006-053711(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00