

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5949937号
(P5949937)

(45) 発行日 平成28年7月13日(2016.7.13)

(24) 登録日 平成28年6月17日(2016.6.17)

(51) Int.Cl. F I
G 0 6 F 12/00 (2006.01) G O 6 F 12/00 5 3 5 F
 G O 6 F 12/00 5 1 8 A

請求項の数 10 (全 24 頁)

(21) 出願番号 特願2014-544217 (P2014-544217)
 (86) (22) 出願日 平成25年8月14日(2013.8.14)
 (86) 国際出願番号 PCT/JP2013/004860
 (87) 国際公開番号 W02014/068820
 (87) 国際公開日 平成26年5月8日(2014.5.8)
 審査請求日 平成27年4月10日(2015.4.10)
 (31) 優先権主張番号 特願2012-237597 (P2012-237597)
 (32) 優先日 平成24年10月29日(2012.10.29)
 (33) 優先権主張国 日本国(JP)

(73) 特許権者 000004237
 日本電気株式会社
 東京都港区芝五丁目7番1号
 (74) 代理人 100109313
 弁理士 机 昌彦
 (74) 代理人 100124154
 弁理士 下坂 直樹
 (72) 発明者 海老山 知生
 東京都港区芝五丁目7番1号
 日本電気株式会社内

審査官 篠塚 隆

最終頁に続く

(54) 【発明の名称】 トランザクションシステム

(57) 【特許請求の範囲】

【請求項1】

キーとバリューの組で表現されるレコードを記憶するレコードデータ記憶部と、
 複数の前記レコードを1つのキーバリューデータに集約したデータグループ単位で、その更新履歴を記憶するWAL記憶部と、

前記データグループを参照および更新するトランザクションのコミット時、前記WAL記憶部に記憶された更新履歴に基づいて、前記レコード単位で排他違反の有無を判定し、排他違反がなければ前記トランザクションを成功させ、その更新履歴を前記WAL記憶部に記録するトランザクション処理部と、

前記WAL記憶部に記憶されている前記更新履歴中のレコードの更新内容を前記レコードデータ記憶部に反映するWAL反映部と
 を有するトランザクションシステム。

10

【請求項2】

前記WAL記憶部は、キーとバリューの組で表現される更新ログを記憶し、前記キーは複数の前記レコードを集約したデータであるデータグループを一意に識別する値であり、前記バリューは前記データグループの更新履歴を表し、また、前記更新履歴は、1以上の部分更新履歴を有し、1つの前記部分更新履歴は、前記データグループを構成する前記レコードの更新内容と更新の順番を表すタイムスタンプとを有し、

前記トランザクション処理部は、

アプリケーションプログラムから入力されたトランザクションで使用する前記データグ

20

ループを構成する前記レコードの最新の値を前記W A L記憶部および前記レコードデータ記憶部から取得し、また、当該データグループに対応する前記更新ログ中の最新の前記タイムスタンプをスタート時タイムスタンプとして前記W A L記憶部から取得し、前記トランザクションを実行するW A L管理部と、

前記トランザクションのコミット時、前記トランザクションのスタート時に取得された前記スタート時タイムスタンプと、前記W A L記憶部に現に記憶されている前記データグループの最新の前記タイムスタンプであるコミット時タイムスタンプとを比較し、前記スタート時タイムスタンプと前記コミット時タイムスタンプとが相違するときは、前記スタート時タイムスタンプより後のタイムスタンプを有する前記部分更新履歴中に、前記トランザクションが使用した前記レコードと一致するレコードが存在するか否かを判定し、前記比較と前記判定の結果とに基づいて、前記トランザクションを成功させるか、失敗させるかを決定し、成功させると決定した場合、前記トランザクションで更新した前記レコードの更新内容と前記コミット時タイムスタンプより後のタイムスタンプとを有する新たな部分更新履歴を、前記データグループに対応して前記W A L記憶部に現に記憶されている前記更新ログに追加する排他チェック処理部と
を有する請求項1に記載のトランザクションシステム。

【請求項3】

前記W A L反映部は、前記レコードデータ記憶部に反映した前記更新ログの複数の部分更新履歴のうち、タイムスタンプの新しい上位所定個数の部分更新履歴は削除せずに残しておく

請求項2に記載のトランザクションシステム。

【請求項4】

前記W A L反映部は、削除せずに残しておいた反映済の前記部分更新履歴を識別する情報を前記更新履歴に記録する

請求項3に記載のトランザクションシステム。

【請求項5】

前記排他チェック処理部は、前記判定では、前記スタート時タイムスタンプより後のタイムスタンプを有する前記部分更新履歴中に、前記トランザクションが参照または更新した前記レコードと一致するレコードが存在するか否かを判定する

請求項2乃至4の何れかに記載のトランザクションシステム。

【請求項6】

前記排他チェック処理部は、前記判定では、前記スタート時タイムスタンプより後のタイムスタンプを有する前記部分更新履歴中に、前記トランザクションが更新した前記レコードと一致するレコードが存在するか否かを判定する

請求項2乃至4の何れかに記載のトランザクションシステム。

【請求項7】

前記W A L反映部は、トランザクションシステムを構成するC P Uのアイドル状態と前記W A L記憶部を構成する記憶装置の使用量とに基づいて、前記反映のタイミングを決定する

請求項1乃至6の何れかに記載のトランザクションシステム。

【請求項8】

キーとバリューの組で表現されるレコードを記憶するレコードデータ記憶部と、

複数の前記レコードを1つのキーバリューデータに集約したデータグループ単位で、その更新履歴を記憶するW A L記憶部と、

トランザクション処理部と、

W A L反映部と

を有するトランザクションシステムの制御方法であって、

前記トランザクション処理部が、前記データグループを参照および更新するトランザクションのコミット時、前記W A L記憶部に記憶された更新履歴に基づいて、前記レコード単位で排他違反の有無を判定し、排他違反がなければ前記トランザクションを成功させ、

その更新履歴を前記WAL記憶部に記録し、

前記WAL反映部が、前記WAL記憶部に記憶されている前記更新履歴中のレコードの更新内容を前記レコードデータ記憶部に反映する

トランザクションシステム制御方法。

【請求項9】

前記WAL記憶部は、キーとバリューの組で表現される更新ログを記憶し、前記キーは複数の前記レコードを集約したデータであるデータグループを一意に識別する値であり、前記バリューは前記データグループの更新履歴を表し、また、前記更新履歴は、1以上の部分更新履歴を有し、1つの前記部分更新履歴は、前記データグループを構成する前記レコードの更新内容と更新の順番を表すタイムスタンプとを有し、

10

前記トランザクション処理部は、

アプリケーションプログラムから入力されたトランザクションで使用する前記データグループを構成する前記レコードの最新の値を前記WAL記憶部および前記レコードデータ記憶部から取得し、また、当該データグループに対応する前記更新ログ中の最新の前記タイムスタンプをスタート時タイムスタンプとして前記WAL記憶部から取得し、前記トランザクションを実行し、

前記トランザクションのコミット時、前記トランザクションのスタート時に取得された前記スタート時タイムスタンプと、前記WAL記憶部に現に記憶されている前記データグループの最新の前記タイムスタンプであるコミット時タイムスタンプとを比較し、前記スタート時タイムスタンプと前記コミット時タイムスタンプとが相違するときは、前記スタート時タイムスタンプより後のタイムスタンプを有する前記部分更新履歴中に、前記トランザクションが使用した前記レコードと一致するレコードが存在するか否かを判定し、前記比較と前記判定の結果とに基づいて、前記トランザクションを成功させるか、失敗させるかを決定し、成功させると決定した場合、前記トランザクションで更新した前記レコードの更新内容と前記コミット時タイムスタンプより後のタイムスタンプとを有する新たな部分更新履歴を、前記データグループに対応して前記WAL記憶部に現に記憶されている前記更新ログに追加する

20

請求項8に記載のトランザクションシステム制御方法。

【請求項10】

キーとバリューの組で表現されるレコードを記憶するレコードデータ記憶部と、

30

複数の前記レコードを1つのキーバリューデータに集約したデータグループ単位で、その更新履歴を記憶するWAL記憶部と

を有するコンピュータを、

前記データグループを参照および更新するトランザクションのコミット時、前記WAL記憶部に記憶された更新履歴に基づいて、前記レコード単位で排他違反の有無を判定し、排他違反がなければ前記トランザクションを成功させ、その更新履歴を前記WAL記憶部に記録するトランザクション処理部と、

前記WAL記憶部に記憶されている前記更新履歴中のレコードの更新内容を前記レコードデータ記憶部に反映するWAL反映部と

して機能させるためのプログラム。

40

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、各々がキーバリューストアである複数のレコードに対して参照および更新するトランザクションを実行するトランザクションシステム、トランザクションシステム制御方法、およびプログラムに関する。

【背景技術】

【0002】

近年、スケーラブルなアプリケーションの実行基盤として、多数のコンピュータを用いるいわゆるクラウドが注目されている。クラウドのインフラストラクチャにおいて、スケ

50

ーラブルなデータ永続化手段として、キーバリューストア (Key Value Store (KVS)) というものが広く知られている。しかしながら、通常のKVSでは1つのレコードに対するトランザクションしか提供できない。

【0003】

そこで、各々がKVSである複数のレコードに対してトランザクションをサポートするシステムが本発明に関連する第1の関連技術として提案されている (例えば非特許文献1参照)。この第1の関連技術では、KVSで複数のレコードに対するトランザクションを実現するために、トランザクションで扱われる複数のレコードをグルーピングし、1つのKey Value (KV) データとしてKVSに書き込むことで、複数レコードに対するトランザクションを実現している。複数のレコードを1つのKVデータに集約したデータを、本明細書では、データグループと呼ぶ。

10

【0004】

また上記第1の関連技術では、KVSに対する排他制御の仕組みとして、楽観的排他を使用している。具体的には、レコード (KVデータ) ごとにタイムスタンプが採番されており、あるトランザクションがスタートしたときに当該タイムスタンプを取得し、レコードを更新してコミットする際にタイムスタンプがトランザクションのスタート時から変化していなければ、タイムスタンプを1つインクリメントしてKVSに書き込む。コミット時にタイムスタンプが変化していれば、当該トランザクションがスタートしてからコミットするまでの間に他のトランザクションによって先にレコードが更新されていることになるので、当該トランザクションを失敗させる。楽観的排他では、排他制御情報などの共有資源を使用しないため、スケールアウトが容易になる。

20

【0005】

他方、悲観的排他と呼ばれる排他制御の仕組みを採用したトランザクションシステムが本発明に関連する第2の関連技術として提案されている (例えば特許文献1参照)。この第2の関連技術では、トランザクションからのレコード更新要求に従うデータファイルの更新処理に際しては、更新するレコードへの同時アクセスを排他制御機構にて排他制御する。また、そのレコードの更新をバッファ上で行わせると共に、バッファ上でのレコード更新の履歴を生成する。トランザクションの確定時には、生成しておいたトランザクションに対応するレコード単位の更新履歴をジャーナルファイルに記録する。そして、当該ジャーナルファイルに記録した更新履歴に基づき書き込みバッファを通してブロック単位でデータファイルを更新する。

30

【0006】

また、マスタサイトのマスタデータとローカルサイトの複製データとの同期を管理する共有データ管理システムにおいて、マスタデータの更新データをローカルサイトの複製データに反映させるタイミングを、更新の緊急度等に基づいて決定することが、本発明に関連する第3の関連技術として提案されている (例えば特許文献2参照)。

【先行技術文献】

【特許文献】

【0007】

【特許文献1】特開平9-62550号公報

40

【特許文献2】特開2001-350661号公報

【非特許文献】

【0008】

【非特許文献1】“ビッグデータの活用に適したスケールアウト型 新データベース「InfoFrame Relational Store」、祐成光樹、田村稔、NEC技報、Vol.65 No.2/2012.

【発明の概要】

【発明が解決しようとする課題】

【0009】

上述したようにスケラブルなデータ永続化手段としてKVSが注目されており、また

50

K V Sで複数のレコードに対してトランザクションをサポートするシステムが提案されている。しかしながら、上述したK V Sに対する楽観的な排他制御の仕組みでは、1つのレコード(K Vデータ)単位の排他しか実現できない。複数のレコードを1つのK Vデータにまとめてトランザクションを実行するシステムにおいては、まとめられた大きな1つのK Vデータが排他の単位となるため、本来は排他違反にならないようなトランザクション同士も排他違反になってしまう可能性があり、トランザクションの同時実行性が下がり、スループットが下がってしまう問題がある。本発明に関連する第2の関連技術に示される悲観的排他を適用し、大きな1つのK Vデータにまとめられる前のレコード単位で排他することが考えられるが、悲観的排他ではスケールアウトが困難になり、K Vデータを使用する意味がなくなる。

10

【0010】

本発明の目的は、上述した課題、すなわち、複数のレコードを1つのK Vデータにまとめてトランザクションを実行するシステムではトランザクションの同時実行性が低下する、という課題を解決するトランザクションシステムを提供することにある。

【課題を解決するための手段】**【0011】**

本発明の第1の観点に係るトランザクションシステムは、
キーとバリューの組で表現されるレコードを記憶するレコードデータ記憶部と、
複数の上記レコードを1つのキーバリューデータに集約したデータグループ単位で、その更新履歴を記憶するW A L記憶部と、

20

上記データグループを参照および更新するトランザクションのコミット時、上記W A L記憶部に記憶された更新履歴に基づいて、上記レコード単位で排他違反の有無を判定し、排他違反がなければ上記トランザクションを成功させ、その更新履歴を上記W A L記憶部に記録するトランザクション処理部と、

上記W A L記憶部に記憶されている上記更新履歴中のレコードの更新内容を上記レコードデータ記憶部に反映するW A L反映部と
を有する。

【0012】

また本発明の第2の観点に係るトランザクションシステム制御方法は、
キーとバリューの組で表現されるレコードを記憶するレコードデータ記憶部と、
複数の上記レコードを1つのキーバリューデータに集約したデータグループ単位で、その更新履歴を記憶するW A L記憶部と、

30

トランザクション処理部と、
W A L反映部と
を有するトランザクションシステムの制御方法であって、

上記トランザクション処理部が、上記データグループを参照および更新するトランザクションのコミット時、上記W A L記憶部に記憶された更新履歴に基づいて、上記レコード単位で排他違反の有無を判定し、排他違反がなければ上記トランザクションを成功させ、その更新履歴を上記W A L記憶部に記録し、

上記W A L反映部が、上記W A L記憶部に記憶されている上記更新履歴中のレコードの更新内容を上記レコードデータ記憶部に反映する。

40

また本発明の第3の観点に係るプログラムは、

キーとバリューの組で表現されるレコードを記憶するレコードデータ記憶部と、
複数の上記レコードを1つのキーバリューデータに集約したデータグループ単位で、その更新履歴を記憶するW A L記憶部と
を有するコンピュータを、

上記データグループを参照および更新するトランザクションのコミット時、上記W A L記憶部に記憶された更新履歴に基づいて、上記レコード単位で排他違反の有無を判定し、排他違反がなければ上記トランザクションを成功させ、その更新履歴を上記W A L記憶部に記録するトランザクション処理部と、

50

上記WAL記憶部に記憶されている上記更新履歴中のレコードの更新内容を上記レコードデータ記憶部に反映するWAL反映部として機能させる。

【発明の効果】

【0013】

本発明は上述した構成を有するため、複数のレコードを1つのKVデータにまとめてトランザクションを実行するシステムにおけるトランザクションの同時実行性を高めることが可能である。

【図面の簡単な説明】

【0014】

【図1】本発明の第1の実施形態のブロック図である。

【図2】本発明の第1の実施形態におけるトランザクション処理部の動作の一例を示すフローチャートである。

【図3】本発明の第1の実施形態におけるレコードデータおよびデータグループの説明に使用するデータの例を示す図である。

【図4】本発明の第1の実施形態におけるデータグループの例を示す図である。

【図5】本発明の第1の実施形態におけるWAL記憶手段に記憶されるWALの説明図である。

【図6】本発明の第1の実施形態における排他チェック処理部で排他違反無しと判断される一例を示す図である。

【図7】本発明の第1の実施形態における排他チェック処理部で排他違反有りとして判断される一例を示す図である。

【図8】本発明の第1の実施形態における排他チェック処理部で排他違反無しと判断される他の例を示す図である。

【図9】本発明の第1の実施形態におけるWAL反映手段の動作の一例を示すフローチャートである。

【図10】本発明の第1の実施形態におけるWAL反映手段で反映されるレコードデータの一例を示す図である。

【図11】本発明の第2の実施形態のブロック図である。

【図12】本発明の第2の実施形態においてタイムスタンプの新しい上位所定個数のWALを削除せずに残しておく説明図である。

【図13】本発明の第3の実施形態のブロック図である。

【図14】本発明の第3の実施形態におけるWAL反映手段の動作の一例を示すフローチャートである。

【図15】本発明の第4の実施形態のブロック図である。

【発明を実施するための形態】

【0015】

次に本発明の実施の形態について図面を参照して詳細に説明する。

[第1の実施形態]

[構成]

図1を参照すると、本発明の第1の実施形態にかかるトランザクションシステムは、トランザクション処理部100、WAL記憶手段200、WAL反映手段300、レコードデータ記憶手段400、システムカタログ500を備え、トランザクション処理部100はさらにWAL管理手段101、リード/ライトセット管理手段102、排他チェック処理部103を備え、排他チェック処理部103はさらにWALデータ分解手段113とレコードデータ照合手段123を備える。

【0016】

トランザクション処理部100は、1以上備えられている。各々のトランザクション処理部100は、アプリケーション(アプリケーションプログラム)600から入力されたトランザクションを受け付け、当該トランザクションを処理する手段である。各々のトラ

10

20

30

40

50

ンザクション処理部100は、他のトランザクション処理部100と並行して動作する。

【0017】

WAL管理手段101は、アプリケーション600から入力されたトランザクションに対して、当該トランザクションが使用するレコードを論理的にまとめたデータグループを作成し、その更新内容をCAS(Compare And Swap)コマンド等を使用して、WALとしてWAL記憶手段200に対して書き込んだり、読み出したりする手段である。

【0018】

ここで、WAL(Write Ahead Log)について説明する。一般にWrite Ahead Loggingとは、データベースに対する操作前の変更内容をログに書き出すことを指し、WALとはログそのものを指す。このように、データベースへの変更は、先ずログに記録され、次いで、ログからデータベースへ反映される。

10

【0019】

WAL記憶手段200は、WAL管理手段101によって生成、または更新されたWALを保存する手段である。WAL記憶手段200は、複数のトランザクション処理部100に共通に備えられる。本実施形態において、WAL記憶手段200はKVSによって構成されており、WAL自体もKeyとValueからなるKVデータとして保存される。KVSは、例えばCASコマンドを使用することにより、1つのKVデータに対するアトミック性は保証することができるが、複数のレコードを一度のトランザクションでアトミックに処理することはできない。このため、複数のレコードを論理的にまとめたデータグループを作成し、このデータグループに対する更新内容をWALとしてKVSに書き込むことで複数レコードをトランザクショナルに処理することを可能としている。

20

【0020】

リード/ライトセット管理手段102は、アプリケーション600から入力されたトランザクションで参照や更新したレコードの情報などを抽出する手段である。リード/ライトセット管理手段102で抽出された参照レコードや更新レコードの情報などは、排他チェック処理部103に渡される。

【0021】

排他チェック処理部103は、トランザクションをコミットする際に、当該トランザクションがスタートしてからコミットするまでにWAL記憶手段200に追加されたWALのデータと、リード/ライトセット管理手段102から渡される、当該トランザクションで参照や更新したレコードとを比較して、同じレコードが存在する場合はトランザクションを失敗させ、同じレコードが存在しない場合はトランザクションを成功させる。トランザクションが成功した場合は、WAL管理手段101が、当該トランザクションによる更新内容をWALとしてWAL記憶手段200に書き込む。

30

【0022】

WALデータ分解手段113は、WAL記憶手段200から取得したWALデータをレコード単位に分解してWALデータに含まれるレコードデータを取得する手段である。

【0023】

レコードデータ照合手段123は、当該トランザクションで参照や更新したレコードと、WALに含まれるレコードデータとを比較して、同じレコードがあるかどうかを調べ、同じレコードが存在する場合はトランザクションを失敗させ、同じレコードが存在しない場合はトランザクションを成功させる手段である。

40

【0024】

WAL反映手段300は、WAL記憶手段200に記憶されているWALに含まれるレコードデータをレコードデータ記憶手段400に反映し、反映したWALをWAL記憶手段200から削除する。WAL反映手段300が動作するタイミングはタイマーで指定してもよいし、システム管理者等が任意のタイミングで行っても良い。

【0025】

レコードデータ記憶手段400は、レコードデータを保存しておく手段である。本実施

50

形態において、レコードデータ記憶手段400はKVSによって構成されており、レコードデータはKeyとValueからなるKVデータとして保存される。なお、WAL記憶手段200とレコードデータ記憶手段400は、図1では別々のKVSとして表現されているが、1つのKVSにまとめても良い。

【0026】

システムカタログ500は、レコードデータ記憶手段400に保存されているデータのメタ情報や、トランザクションを実行する際に作成するデータグループの情報などのデータを保存する手段である。アプリケーション600からトランザクションが入力された時、WAL管理手段101はシステムカタログ500からデータグループに関連する情報を抽出し、抽出された情報を元に当該トランザクションで使用するデータグループを特定する。

10

【0027】

[動作]

次に、本実施形態の動作について説明する。本実施形態は図2に示すフローチャートに従って動作を行う。

【0028】

まず、トランザクション入力ステップS001で、アプリケーション600が本システムに対してトランザクションを入力する。

【0029】

次に、WAL特定のステップS002で、WAL管理手段101が、入力されたトランザクションを解析し、システムカタログ500の情報を用いて当該トランザクションで処理するデータが所属するデータグループを特定し、当該データグループのこれまでの更新内容が保存されているWALを特定する。本システムにおいて、WALとはトランザクションで扱われる複数のレコードを論理的にまとめたデータグループの更新内容の履歴を1つのKVデータにまとめたものである。

20

【0030】

ここで、本システムにおけるデータグループの概念を説明する。図3に示すような二種類のデータ、すなわち、商品テーブルT001と入札テーブルT002があるとする。このようなデータに対して、以下のようなトランザクションを実行するとする。

「ユーザーが新たに入札を行った際に、商品テーブルと入札テーブルのデータを更新する」

30

【0031】

このようなトランザクションを処理する場合、商品テーブルのデータと入札テーブルのデータをアトミックに更新しなければならない。しかし、KVSは1つのKVデータに対するアトミック性しか保証していないため、例えば、商品テーブルのデータを更新したあと、入札テーブルのデータの更新に失敗した場合、商品テーブルのデータだけが更新された状態になってしまい、アトミック性が保証されないという問題がある。

【0032】

そこで、トランザクションに関連する複数のデータを論理的にまとめて1つのKVデータとし、複数のデータ更新を1回のKVSに対する書き込みで実現することでトランザクションのアトミック性を保証する。

40

【0033】

データグループを作成するには、まず各データで共通するKey(図3の例では商品ID)を抽出する。図3に示すデータからは、図4に示すような3つのグループ(G001、G002、G003)が論理的にまとめられる。

【0034】

次に、本システムにおけるWALの概念を説明する。図3に示すデータに対して以下のような2つのトランザクションが順番に実行されたとする。

最初のトランザクション：「入札テーブルT002のIDがb1のレコードの入札額を12000に更新し、商品テーブルT001の商品IDがi1の最高金額を12000に更

50

新する。」

二番目のトランザクション：「入札テーブルT002のIDがb2のレコードの入札額を13000に更新し、商品テーブルT001の商品IDがi1の最高金額を13000に更新する。」

【0035】

まず、最初のトランザクション1によって商品テーブルT001の1行目のデータと入札テーブルT002の1行目のデータが更新され、二番目のトランザクション2によって商品テーブルT001の1行目のデータと入札テーブルT002の2行目のデータが更新される。これらの2つのトランザクションで更新されるデータは全て同じデータグループに所属しているため、これらのトランザクションによる更新結果は1つのWALに記録される。WALは以下のような情報で構成される。

- ・当該WALを識別するためのKey値
- ・トランザクションによる更新結果

【0036】

WALはKVデータであるため、(Key, Value)の組で表される。Key値としてデータグループで共通のKey値にWALであることを示す接頭辞を付加したものをを用いる。上記の例では、WAL_i1が当該WALのKey値となる。Value値にはトランザクションによる更新結果が含まれる。具体的には、以下のようにWALが構成される。

WAL=(Key, Value)=(WAL_i1, <[{i1, AAA, 10, 5, 12000}], {b1, i1, u3, 12000}, タイムスタンプ1], [{i1, AAA, 10, 5, 13000}], {b2, i1, u1, 13000}, タイムスタンプ2]>)

【0037】

タイムスタンプは当該WALに対するトランザクションが実行された順番を表している。タイムスタンプは単純な数値で表されていれば良い。このように、WALはあるデータグループに対する更新履歴を表している。

【0038】

WAL特定のステップS002でWAL管理手段101は、具体的には以下のようにWALを特定する。

【0039】

まず、WAL管理手段101は、入力されたトランザクションから、当該トランザクションで利用されるテーブルを抽出し、システムカタログ500から当該トランザクションで利用されるテーブル間のデータグループに関係する情報を抽出する。

【0040】

システムカタログ500には、テーブルのスキーマ情報のほかに、どのテーブルがどの列でデータグループを構成するのかという情報が格納されている。例えば、図3に示すようなデータの場合には、「商品テーブルT001と入札テーブルT002でデータグループを構成し、データグループを構成する際のKeyは商品IDである。」という情報が格納されている。

【0041】

WAL管理手段101は、入力されたトランザクションとシステムカタログ500の情報から、当該トランザクションが利用するデータグループを特定する。例えば、「入札テーブルT002のIDがb1のレコードの入札額を12000に更新し、商品テーブルT001の商品IDがi1の最高金額を12000に更新する」というトランザクションが入力された場合は、当該トランザクションが利用するテーブルが商品テーブルT001と入札テーブルT002であることが分かり、システムカタログ500から商品テーブルT001と入札テーブルT002が商品IDでグループ化されていることが分かる。これにより、取得すべきWALがWAL_i1であることが分かり、当該トランザクションが利用するWALが特定できる。

【0042】

10

20

30

40

50

次に、WAL特定のステップS002で特定されたWALの情報を用いて、WAL管理手段101がWAL記憶手段200からWALを取得する(ステップS003)。ここで、WAL記憶手段200にWALが存在する場合は、WAL取得のステップS004でWAL管理手段101がWAL記憶手段200からWALを取得する。このとき、後の排他チェックで使用するために、当該WALの最新のタイムスタンプをWAL記憶手段200から取得しておく。

【0043】

ここで、図5を用いてWAL記憶手段200にWALがどのように保存されているのかを説明する。

【0044】

WALはWAL記憶手段200の中でそれぞれのWALごとに履歴を持って保存されている。例えば、あるトランザクションがWAL_{i1}を用いてトランザクションを実行すると、WAL_{i1}のタイムスタンプ1としてWALが蓄積され、次のトランザクションでWAL_{i1}を用いてトランザクションを実行すると、WAL_{i1}のタイムスタンプ2としてWALが蓄積される。次のトランザクションでWAL_{i2}を用いてトランザクションを実行するとWAL_{i2}のタイムスタンプ1としてWALが蓄積される。

【0045】

それぞれのWALに含まれるレコードについて、タイムスタンプの値が古いものから順番にWALのデータを確認すれば、各レコードの更新履歴や最新のレコード値を知ることができる。

【0046】

WAL記憶手段200に該当するWALが存在しない場合は、レコード取得のステップS006で、WAL管理手段101がレコードデータ記憶手段400から当該トランザクションで利用するレコードデータを取得する。

【0047】

ここでレコードデータ記憶手段400にレコードがどのように保存されているのかを説明する。レコードデータ記憶手段400はKVSであるため、レコードデータは(Key, Value)の組で表されるデータとして保存されている。

【0048】

図3に示すようなレコードデータは、主キーがKey値となり、それ以外の列がValue値として保存されている。具体的には、以下のように保存されている。

(Key, Value) = (i1, {AAA, 10, 5, 10000})

(Key, Value) = (b1, {i1, u3, 9000})

【0049】

WAL記憶手段200に該当するWALが存在する場合でも、WAL記憶手段200に保存されているWALの中に必要なレコードデータが全て含まれていない場合がある。以下にその場合の例を示す。

【0050】

例えば、トランザクション1で以下のようなトランザクションを実行したとする。

「入札テーブルT002のIDがb1のレコードの入札額を12000に更新し、商品テーブルT001の商品IDがi1の最高金額を12000に更新する」

このときWALは存在せず、新たにWAL_{i1}のWALが作成され、WAL記憶手段200に保存されたとする。このとき、WAL_{i1}のWALとして保存されるWALは以下のようなものであり、レコードデータとしては、商品テーブルT001の1行目のレコードと入札テーブルT002の1行目のレコードが含まれる。

WAL = (WAL_{i1}, [{i1, AAA, 10, 5, 12000}, {b1, i1, u3, 12000}, 1])

次に、トランザクション2で以下のようなトランザクションを実行したとする。

「入札テーブルT002のIDがb2のレコードの入札額を13000に更新し、商品テーブルT001の商品IDがi1の最高金額を13000に更新する。」

10

20

30

40

50

W A L 管理手段 1 0 1 は上記のトランザクション 2 が入力されると、W A L _ i 1 を特定し、W A L 記憶手段 2 0 0 から W A L _ i 1 の W A L を取得する。このとき、W A L _ i 1 にはタイムスタンプが 1 の W A L のみが含まれており、当該 W A L の中にはレコードデータとして、商品テーブル T 0 0 1 の 1 行目のレコードと入札テーブル T 0 0 2 の 2 行目のレコードが含まれている。しかし、トランザクション 2 では、商品テーブル T 0 0 1 の 1 行目のレコードと、入札テーブル T 0 0 2 の 2 行目のレコードを更新するため、W A L を取得しただけでは、トランザクションに必要なデータが全て揃わない。

【 0 0 5 1 】

このように、W A L 管理手段 1 0 1 は W A L 記憶手段 2 0 0 から W A L を取得した後に、当該トランザクションの処理に必要なデータが当該 W A L の中に全て含まれているかを
10
チェックする（ステップ S 0 0 5）。もし、必要なデータが W A L の中に含まれていない場合は、レコードデータ記憶手段 4 0 0 から必要な分のレコードを取得する（ステップ S 0 0 6）。

【 0 0 5 2 】

W A L やレコードデータからトランザクションに必要なデータを全て取得したら、それらのデータを使用してトランザクションの処理を進める。そして、トランザクション実行のステップ S 0 0 7 で、W A L 管理手段 1 0 1 が、当該トランザクションによるデータグループの更新イメージを新たに W A L として作成する。この時点ではまだ W A L 管理手段 1 0 1 は W A L 記憶手段 2 0 0 に新たに作成した W A L を書き込まない。

【 0 0 5 3 】

次に、参照・更新レコード抽出のステップ S 0 0 8 で、リード/ライトセット管理手段 1 0 2 は当該トランザクションで参照や更新したレコード情報を抽出する。具体的には、W A L 管理手段 1 0 1 から当該トランザクションで更新や参照したレコードの情報と、当該トランザクションのスタート時点のタイムスタンプ情報を取得する。なお、リード/ライトセット管理手段で保持しておくレコード情報はレコードデータ全体である必要はなく、各レコードの K e y 値だけで十分である。
20

【 0 0 5 4 】

リード/ライトセット管理手段 1 0 2 は、排他チェック処理部 1 0 3 に対して、保持しているレコード情報を出力するが、このとき、参照したレコード情報と更新したレコード情報とを含む全てのレコード情報を出力してもよいし、更新したレコードの情報だけを出力してもよい。参照と更新のすべてのレコード情報を出力するか、更新したレコード情報だけを出力するかで、トランザクション間の I s o l a t i o n レベルを調整することができる。どの種別のレコード情報を出力するかは、システムが起動する前にシステム管理者などが予め設定しておけばよい。
30

【 0 0 5 5 】

次に、排他チェックのステップ S 0 0 9 で、排他チェック処理部 1 0 3 は当該トランザクションがスタートしてから現時点までに W A L 記憶手段 2 0 0 に追加された W A L データと、リード/ライトセット管理手段 1 0 2 から渡される当該トランザクションで参照や更新したレコードとを比較して（ステップ S 0 1 0）、同じレコードが存在する場合はトランザクションを失敗させ（ステップ S 0 1 2）、同じレコードが存在しない場合はトランザクションを成功させる（ステップ S 0 1 1）。
40

【 0 0 5 6 】

具体的には、排他チェック処理部 1 0 3 はまず、リード/ライトセット管理手段 1 0 2 から当該トランザクションで参照や更新されたレコードデータとスタート時点のタイムスタンプの値を受け取る。次に、W A L 記憶手段 2 0 0 から当該トランザクションで使用された W A L を取得する。

【 0 0 5 7 】

W A L 記憶手段 2 0 0 から取得した W A L の最新のタイムスタンプがリード/ライトセット管理手段 1 0 2 から受け取ったスタート時点のタイムスタンプと同じ値の場合は、当該トランザクションがスタートしてから現時点までの間に当該 W A L に関連する他のトラ
50

ンザクションが実行されていないことになるため、トランザクションを成功させる。

【0058】

WAL記憶手段200から取得したWALの最新のタイムスタンプが、リード/ライトセット管理手段102から受け取ったスタート時点のタイムスタンプの値より大きな値の場合は、当該トランザクションがスタートしてから現時点までの間に当該WALに関連する他のトランザクションが実行されていることを意味するため、排他違反がないかチェックする。

【0059】

排他違反チェックは以下のようにして実行される。まず、WALデータ分解手段113が、WAL記憶手段200から取得したWALデータのうち、リード/ライトセット管理手段102から受け取ったタイムスタンプの値より大きなタイムスタンプの値を持つWALデータを全て抽出する。そして、抽出したWALデータをレコード単位に分解して当該WALデータに含まれるレコードデータを抽出する。

【0060】

次に、レコードデータ照合手段123が、抽出したWALデータに含まれるレコードデータと、リード/ライトセット管理手段102から受け取ったレコードデータに一致するものが含まれるかどうかをチェックする。もし、一致するものが1つでも存在する場合は排他違反が発生しているため、トランザクションを失敗させる。一致するものが存在しない場合は排他違反が発生していないため、トランザクションを成功させる。

【0061】

排他チェック処理部103がトランザクションを成功させる場合、失敗させる場合を説明する概念図を図6、図7、図8に示す。

【0062】

図6は、トランザクションがスタートしてから排他チェック処理部103が排他チェックを実行するまでの間に当該WALに更新が無かった場合の例である。この場合、当該トランザクションが実行されている間に当該WALに関連する他のトランザクションが実行されていないため、当該トランザクションは成功する。

【0063】

図7は、トランザクションがスタートしてから排他チェック処理部103が排他チェックを実行するまでの間に当該WALに更新があった場合の例である。この場合、排他チェック処理部103で持っているスタート時点のタイムスタンプ値より大きなタイムスタンプ値を持つWALデータを抽出する。図7の例では、タイムスタンプが3のWALデータを取得する。図7の例では、タイムスタンプが3のWALデータにはレコードAとレコードDが含まれている。つまり、当該トランザクションがスタートしてから排他チェック処理部103が排他チェックを実行するまでの間に、他のトランザクションがレコードAとレコードDを更新するトランザクションを実行したことを意味している。ここで、当該トランザクションを実行するとレコードAに関して排他違反が発生するため、当該トランザクションを失敗させる。

【0064】

図8は、トランザクションがスタートしてから排他チェック処理部103が排他チェックを実行するまでの間に当該WALに更新があった場合の例である。この場合、排他チェック処理部103で持っているスタート時点のタイムスタンプ値以上のタイムスタンプ値を持つWALデータを抽出する。図8の例ではタイムスタンプ値が3のWALデータを取得する。図8の例では、タイムスタンプが3のWALデータにはレコードCとレコードDが含まれている。つまり、当該トランザクションがスタートしてから排他チェック処理部103が排他チェックを実行するまでの間に、他のトランザクションがレコードCとレコードDを更新するトランザクションを実行したことを意味している。ここで、当該トランザクションを実行してもレコードレベルでは排他違反は発生していないため、当該トランザクションを成功させる。

【0065】

10

20

30

40

50

排他チェック処理部 103 で排他違反が検出されなかった場合は、排他チェック処理部 103 は WAL 管理手段 101 にトランザクションのコミットを要求する。WAL 管理手段 101 はトランザクションのコミットが要求されると、当該トランザクションの WAL データを WAL 記憶手段 200 に書き込む (ステップ S011)。このとき、書き込む WAL のタイムスタンプの値は、WAL 記憶手段 200 に記憶されている最新のタイムスタンプの次の値とする。例えば図 6 の場合、更新後のレコード A、B に値 3 のタイムスタンプを付加して書き込む。また例えば図 8 の場合、更新後のレコード A、B に値 4 のタイムスタンプを付加して書き込む。

【0066】

排他チェック処理部 103 で排他違反が検出された場合は、排他チェック処理部 103 は WAL 管理手段 101 にトランザクションのアボートを要求する。WAL 管理手段 101 はトランザクションのアボートが要求されると、当該トランザクションの WAL データを破棄する (ステップ S012)。

10

【0067】

次に、本実施形態において WAL データをレコードデータ記憶手段 400 に反映する動作について説明する。WAL データの反映は図 9 に示すフローチャートに従って動作を行う。

【0068】

まず、WAL データ取得のステップ S101 で WAL 反映手段 300 が WAL 記憶手段 200 から、記憶されている WAL データを全て取得する。

20

【0069】

WAL 記憶手段 200 に WAL データが存在しない場合は WAL 反映の処理を終了する (ステップ S102 の no)。

【0070】

WAL 記憶手段 200 に WAL データが存在する場合 (ステップ S102 の yes) は、WAL データ反映のステップ S103 で WAL 反映手段 300 が WAL 記憶手段 200 から取得した WAL データに含まれるレコードデータを抽出し、抽出したレコードデータをレコードデータ記憶手段 400 に反映させる。

【0071】

WAL 反映手段 300 が WAL 記憶手段 200 から取得した WAL データには同じレコードデータが複数含まれている場合がある。同じレコードデータが複数含まれている場合は、タイムスタンプがもっとも大きな値の WAL に含まれるレコードデータが当該レコードデータについて最新のデータであるため、これを反映する。

30

【0072】

図 10 を用いて具体的に説明する。WAL 反映手段 300 が WAL 記憶手段 200 から図 10 に示すような WAL を取得した場合、レコード A については、タイムスタンプが 1 の WAL データに含まれるものが最新のため、このデータをレコードデータ記憶手段 400 に反映する。レコード B とレコード C については、タイムスタンプが 3 の WAL データに含まれるものが最新のため、このデータをレコードデータ記憶手段 400 に反映する。

40

【0073】

最後に、WAL データ削除のステップ S104 で WAL 反映手段 300 が WAL 記憶手段 200 から抽出し、反映した WAL データを削除する。WAL 反映手段 300 が WAL 記憶手段 200 からデータを取得し、反映するまでの間にトランザクションが実行され、新たな WAL が追加されている可能性があるが、新たに追加された WAL については削除しない。

【0074】

あるレコードについて、WAL にデータが存在する場合は、WAL に保存されているレコードデータが当該レコードについて最新の情報であり、WAL 反映手段 300 がレコードデータ記憶手段 400 にレコードデータを反映した時点で当該レコードに対する最新の

50

データがレコードデータ記憶手段400に保存される。あるレコードについてWALにデータが存在しない場合は、レコードデータ記憶手段400に保存されているレコードデータが当該レコードについて最新の情報となる。

【0075】

なお、アプリケーションから入力されたトランザクションを処理するプロセスと、WALを反映するプロセスはそれぞれ並列に動作している。WALを反映するプロセスはタイマーによって一定時間ごとに行われてもよいし、システム管理者等が任意のタイミングで行っても良い。

【0076】

[効果]

本実施形態は、KVSに対して複数レコードのトランザクションをサポートするシステムに対して、トランザクションの排他制御を行う際に、トランザクションを実行するデータグループ単位で排他制御を行うのではなく、当該トランザクションで参照や更新したレコード単位で排他制御を行うことで、トランザクションの排他単位を小さくすることができ、トランザクションの同時実効性を高めスループットを上げることができる。

【0077】

[第2の実施形態]

次に、本発明の第2の実施形態について説明する。

【0078】

図11に示すように、本実施形態は図1に示した第1の実施形態の構成要素に加えて、WAL反映手段300の中にWAL削除制御手段301を備える点で異なる。

【0079】

本実施形態では、WAL反映手段300はWALデータをレコードデータ記憶手段400に反映させた後に、反映させた全てのWALデータを削除するのではなく、一定量のWALデータを削除せず残しておく。

【0080】

図12に本実施形態におけるWAL反映手段300によるWALデータ管理の概念図を示す。WAL反映手段300が処理実行時にタイムスタンプ1からタイムスタンプ6までのWALデータが存在したとする。WAL削除制御手段301には、WALデータ反映後のWALデータ削除時に、反映済みのWALデータのうち、どれぐらいのWALデータを削除せずに残しておくかの情報を持っている。図12に示す例ではWAL削除制御手段301は3個のWALデータを削除せずに残しておくという情報を持っている。なお、タイムスタンプ7からタイムスタンプ9までのWALデータは、WAL反映手段300によってWAL反映処理が開始してからWAL反映処理が終了するまでの間に実行されたトランザクションによって新たに生成されたWALデータである。このように新たなWALデータはトランザクションの実行によって逐次蓄積されている。

【0081】

本実施形態におけるWAL反映処理は基本的には図9に示すフローチャートと同様の手順で動作を行う。

【0082】

まず、WALデータ取得のステップS101でWAL反映手段300がWAL記憶手段200からそこに記憶されているWALデータを全て取得する。

【0083】

次に、WALデータ反映のステップS102でWAL反映手段300がWAL記憶手段200から取得したWALデータに含まれるレコードデータを抽出し、抽出したレコードデータをレコードデータ記憶手段400に反映させる。

【0084】

最後に、WALデータ削除のステップS103で、WAL削除制御手段301が、反映したWALデータのうちタイムスタンプが大きな値のものから順番にN個(Nは自然数)だけ残して、WALデータを削除する。すなわち、タイムスタンプの新しい上位所定個数

10

20

30

40

50

のWALデータは削除せずに残しておく。

【0085】

次に、反映したWALデータを全て削除せずに一部残しておくことの効果について説明する。トランザクションの実行プロセスとWAL反映のプロセスはそれぞれ並列に動作しているため、トランザクションがスタートしてからコミットするまでの間にWAL反映のプロセスが動作する可能性がある。

【0086】

もし、あるトランザクションがスタートしてからコミットするまでの間にWAL反映処理が動作した場合、当該トランザクションに関連するWALデータはWAL記憶手段200から全て削除されてしまう。WALデータが削除された後に当該トランザクションの排他チェックを実行しようとしても、WALデータが存在しないため排他違反があるのかわからないのが判断できない。このため、本来成功するはずのトランザクションも失敗することになる。これに対して、WAL反映処理によって全てのデータが削除されず、トランザクションの排他チェックのために十分なWALデータが残っていれば当該トランザクションの排他違反を正確に検出できる。

10

【0087】

例えば、本発明の第1の実施形態で参照した図6乃至図8においては、右側の破線内のWALの状態に基づいて、図6および図8の場合は排他違反無し、図7の場合は排他違反有りが検出された。若し、図6乃至図8において、右側の破線内のWALが全て削除されてしまっていると、そのような検出は困難になる。従って、安全を考えて図6乃至図8の全てのケースにおいて排他違反有りと判断せざるを得ず、トランザクションの同時実行性が低下する。

20

【0088】

このように、本実施形態では、WAL反映手段300によって反映したWALデータを全て削除するのではなく、一部のWALデータを残しておくことによって、WALデータが存在しないためにトランザクションの実行が失敗してしまうという事態を避け、本来失敗しないはずのトランザクションを成功させることができ、トランザクションの同時実行性を向上させスループットを上げることができる。

【0089】

なお、どれぐらいの量のWALを削除せずに残しておくかは、求めるトランザクションの同時実行性、使用する記憶容量とのトレードオフである。使用する記憶容量を抑えるためには、残しておくWALは少なくすべきであり、トランザクションの同時実行性を上げるためには、残しておくWALは多くすべきである。

30

【0090】

WAL削除制御手段301が管理する、残しておくWALの量は固定値であってもよいし、WALが存在しないために排他チェックが行えず失敗したトランザクションの量からフィードバックを受けて自動で調整されてもよい。つまり、排他チェックが行えず失敗したトランザクションの量が多いならば残しておくWALの量を大きくし、排他チェックが行えずに失敗したトランザクションの量が少ないなら残しておくWALの量を少なくすればよい。

40

【0091】

また、本実施形態におけるWALの構造は、以下のようになっている。

- ・当該WALを識別するためのKey値
- ・データグループに対する更新結果
- ・WAL反映位置

【0092】

即ち、第1の実施形態と比較して、新たにWAL反映位置が追加されている。これは、WALデータのうちレコードデータ記憶手段400に反映済みのデータ位置を示している。図12に示すWALデータの例では、WAL反映位置はタイムスタンプ6となる。WAL管理手段101などがWAL記憶手段200からWALデータを取得した際にWAL反

50

映位置を参照することによって、WALデータのうち、どのデータがレコードデータ記憶手段400に反映済みであるのかが分かる。同様に、WAL反映手段300は、WAL記憶手段200からWALデータを取得する際にWAL反映位置を参照することによって、WALデータのうち、どのデータがレコードデータ記憶手段400に既に反映済みであるのかが分かる。

【0093】

[第3の実施形態]

次に、本発明の第3の実施形態について説明する。

【0094】

図13に示すように、本実施形態は図1に示した第1の実施形態の構成要素に加えて、WAL反映タイミング制御手段310を備える点で異なる。

10

【0095】

本実施形態では、WAL反映手段300は一定時間ごとやシステム管理者等による動作指示ではなく、WAL反映タイミング制御手段310からの指示によってWAL反映処理を実行する。

【0096】

WAL反映処理はすべてのWALに含まれるレコードデータをレコードデータ記憶手段400に反映するため、負荷が高い処理である。このため、アプリケーションから頻りにトランザクションが入力されているようなシステムに高い負荷がかかっているときにWAL反映処理が動作するとシステムにさらに負荷がかかってしまい、トランザクションの処理性能に悪影響を及ぼす。他方、WALはレコードデータに対する更新履歴が全て含まれているため、WALを反映せずに残しておく、WALを保存するための記憶装置の使用容量が大きくなってしまふ。

20

【0097】

WAL反映タイミング制御手段310は、現在のシステムの負荷状況や記憶装置の使用量を監視して、システムの負荷が低い場合や、記憶装置の使用量が大きくなってきた場合にWAL反映手段300に動作を指示する。

【0098】

例えば、以下のような評価式を用いて、評価値がある閾値を超えた場合にWAL反映タイミング制御手段310は、WAL反映手段300に動作を指示する。

30

評価値 = $W1 * CPUのアイドル状況 + W2 * (記憶装置の使用量 / 記憶装置の全容量)$

【0099】

上記評価式において、W1やW2は重み係数であり、システムの負荷の状況と記憶装置の使用量のどちらを優先するかを調整することができる。例えば、W1を大きくした場合は、よりシステムの負荷状況を優先した動作になる。W1やW2などの重み係数や、評価式に対する閾値などはシステム管理者がシステム起動前に予め設定しておけばよい。

【0100】

本実施形態におけるWAL反映処理のフローチャートを図14に示す。

【0101】

まず、評価値計算のステップS301で、WAL反映タイミング制御手段310が評価値の計算を行う。

40

【0102】

次に、WAL反映タイミング制御手段310は、ステップS301で計算した評価値が閾値を超えているかどうかを判定する(ステップS302)。評価値が閾値を超えていない場合は、評価値計算のステップS301に戻る。評価値が閾値を超えている場合は、WAL反映タイミング制御手段310はWAL反映手段300に動作指示を行う。

【0103】

WAL反映タイミング制御手段310から動作指示を受けたWAL反映手段300は、WAL反映処理を実行する。ステップS304以降の動作は図9に示したWAL反映の動作フローと同じため説明を割愛する。

50

【0104】

このように、本実施形態では、WAL反映タイミング制御手段310がシステムの負荷状況や記憶装置の使用量などを考慮してWAL反映の動作タイミングを制御することで、本システムが稼動しているコンピュータの資源を有効に使うことができる。

【0105】

[第4の実施形態]

図15は、本発明の第4の実施形態のブロック図である。図15に示すように、本実施形態は、1以上のトランザクション処理部1と、レコードデータ記憶部2と、WAL記憶部3と、WAL反映部4とを有する。

【0106】

レコードデータ記憶部2は、半導体メモリや磁気ディスク等で構成され、キーとバリュウーの組で表現されるレコードを記憶する。

【0107】

WAL記憶部3は、半導体メモリや磁気ディスク等で構成され、複数の上記レコードを1つのキーバリュウーデータに集約したデータグループ単位で、その更新履歴を記憶する。好ましくは、WAL記憶部3は、キーとバリュウーの組で表現される更新ログを記憶する。上記キーは、複数の上記レコードを集約したデータであるデータグループを一意に識別する値である。また上記バリュウーは、上記データグループの更新履歴を表す。さらに上記更新履歴は、1以上の部分更新履歴を有する。そして、1つの上記部分更新履歴は、上記データグループを構成する上記レコードの更新内容と更新の順番を表すタイムスタンプとを

【0108】

トランザクション処理部1は、図示しないアプリケーションプログラム等から入力されたトランザクションを実行する機能を有する。トランザクション処理部1は、上記データグループを参照および更新するトランザクションのコミット時、WAL記憶部3に記憶された更新履歴に基づいて、レコード単位で排他違反の有無を判定し、排他違反がなければ当該トランザクションを成功させ、その更新履歴をWAL記憶部3に記録する。好ましくは、トランザクション処理部1は、主な機能部として、WAL管理部5と排他チェック処理部6とを有する。

【0109】

WAL管理部5は、入力されたトランザクションで使用する上記データグループを構成する上記レコードの最新の値をWAL記憶部3およびレコードデータ記憶部2から取得する機能を有する。またWAL管理部5は、当該データグループに対応する上記更新ログ中の最新の上記タイムスタンプをスタート時タイムスタンプとしてWAL記憶部3から取得する機能を有する。

【0110】

排他チェック処理部6は、トランザクション処理部1で実行しているトランザクションのコミット時、当該トランザクションのスタート時に取得しておいた上記スタート時タイムスタンプと、WAL記憶部3に現に記憶されている上記データグループの最新のタイムスタンプであるコミット時タイムスタンプとを比較する機能を有する。また排他チェック処理部6は、上記の比較の結果、スタート時タイムスタンプとコミット時タイムスタンプとが相違するときは、スタート時タイムスタンプより後のタイムスタンプを有する部分更新履歴中に、当該トランザクションが使用したレコードと一致するレコードが存在するかどうかを判定する機能を有する。さらに排他チェック処理部6は、上記比較と上記判定の結果とに基づいて、上記トランザクションを成功させるか、失敗させるかを決定し、成功させると決定した場合、上記トランザクションで更新したレコードの更新内容とコミット時タイムスタンプより後のタイムスタンプとを有する新たな部分更新履歴を、上記データグループに対応してWAL記憶部3に現に記憶されている更新ログに追加する機能を有する。

【0111】

W A L 反映部 4 は、W A L 記憶部 3 に記憶されている更新履歴中のレコードの更新内容をレコードデータ記憶部 2 に反映する機能を有する。

【 0 1 1 2 】

本実施形態のトランザクションシステムは、例えば、M P U 等のプロセッサとそれに接続されたメモリ、入出力装置、通信インターフェイス等を有するコンピュータと、当該コンピュータで実行されるプログラムとで実現することができる。プログラムは、磁気ディスク装置等のコンピュータ可読記録媒体に記録されており、コンピュータの立ち上げ時に M P U 等のプロセッサに読み取られ、そのプロセッサの動作を制御することにより、そのプロセッサ上に W A L 管理部 5、排他チェック処理部 6、W A L 反映部 4 を実現する。

【 0 1 1 3 】

次に本実施形態の動作を説明する。

【 0 1 1 4 】

トランザクション処理部 1 は、入力されたトランザクションを実行する。そして、トランザクション処理部 1 は、データグループを参照および更新するトランザクションのコミット時、W A L 記憶部 3 に記憶された更新履歴に基づいて、レコード単位で排他違反の有無を判定し、排他違反がなければ当該トランザクションを成功させ、その更新履歴を W A L 記憶部 3 に記録する。その後、W A L 反映部 4 は、W A L 記憶部 3 に記憶されている更新履歴中のレコードの更新内容をレコードデータ記憶部 2 に反映する。

【 0 1 1 5 】

より好ましくは、本実施形態は以下のように動作する。

【 0 1 1 6 】

トランザクション処理部 1 は、トランザクションが入力されると、そのトランザクションで使用されるデータグループを構成するレコードの最新の値を W A L 記憶部 3 およびレコードデータ記憶部 2 から取得し、また、当該データグループに対応する更新ログ中の最新のタイムスタンプをスタート時タイムスタンプとして W A L 記憶部 3 から取得し、トランザクションの実行を開始する。

【 0 1 1 7 】

トランザクション処理部 1 は、実行しているトランザクションのコミット時、そのトランザクションのスタート時に取得しておいた上記スタート時タイムスタンプと、W A L 記憶部 3 に現に記憶されているデータグループの最新のタイムスタンプであるコミット時タイムスタンプとを比較する。次にトランザクション処理部 1 は、スタート時タイムスタンプとコミット時タイムスタンプとが相違するときは、スタート時タイムスタンプより後のタイムスタンプを有する全ての部分更新履歴を W A L 記憶部 3 から読み出し、それらの部分更新履歴中に、当該トランザクションが使用したレコード（参照または更新したレコード、或いは更新したレコード）と一致するレコードが存在するか否かを判定する。

【 0 1 1 8 】

トランザクション処理部 1 は、上記の比較と判定の結果とに基づいて、上記トランザクションを成功させるか、失敗させるかを決定する。具体的には、スタート時タイムスタンプとコミット時タイムスタンプとが同じであるか、同じでなくても、スタート時タイムスタンプより後のタイムスタンプを有する全ての部分更新履歴中に、当該トランザクションが使用したレコード（参照または更新したレコード、或いは更新したレコード）と一致するレコードが存在しなければ、トランザクションを成功させる。しかし、スタート時タイムスタンプとコミット時タイムスタンプとが異なっており、しかも、スタート時タイムスタンプより後のタイムスタンプを有する部分更新履歴中に、当該トランザクションが使用したレコード（参照または更新したレコード、或いは更新したレコード）と一致するレコードが存在していれば、トランザクションを失敗させる。

【 0 1 1 9 】

トランザクション処理部 1 は、トランザクションを成功させると決定した場合、当該トランザクションで更新したレコードの更新内容とコミット時タイムスタンプより後のタイムスタンプとを有する新たな部分更新履歴を作成して、当該トランザクションが使用した

10

20

30

40

50

データグループに対応してW A L 記憶部 3 に現に記憶されている更新ログに追加する。更新ログは、キーとバリューの組で表現される 1 つのキーバリューデータであるため、C A S コマンド等を使用することで、データ操作のアトミック性を保証することができる。

【 0 1 2 0 】

他方、W A L 反映部 4 は、例えば定期的に、W A L 記憶部 3 に記憶されている更新ログ中のレコードの更新内容をレコードデータ記憶部 2 に反映する。レコードは、キーとバリューの組で表現される 1 つのキーバリューデータであるため、C A S コマンド等を使用することで、データ操作のアトミック性を保証することができる。

【 0 1 2 1 】

こうして本実施形態によれば、複数のレコードを 1 つの K V データにまとめてトランザクションを実行するシステムにおけるトランザクションの同時実行性を高めることが可能である。その理由は、レコード単位で楽観的な排他を行うためである。

【 0 1 2 2 】

以上、上記各実施形態を参照して本発明を説明したが、本発明は、上述した実施形態に限定されるものではない。本発明の構成や詳細には、本発明の範囲内で当業者が理解しうる様々な変更をすることができる。

【 0 1 2 3 】

なお、本発明は、日本国にて 2 0 1 2 年 1 0 月 2 9 日に特許出願された特願 2 0 1 2 - 2 3 7 5 9 7 の特許出願に基づく優先権主張の利益を享受するものであり、当該特許出願に記載された内容は、全て本明細書に含まれるものとする。

【 産業上の利用可能性 】

【 0 1 2 4 】

本発明は、データ永続化手段として K V S を用いてトランザクションを実行するシステムにおいて利用可能である。

【 符号の説明 】

【 0 1 2 5 】

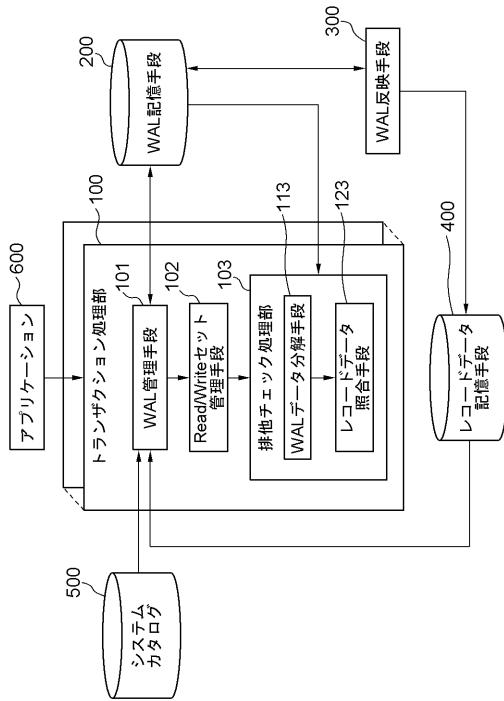
- 1 ... トランザクション処理部
- 2 ... レコードデータ記憶部
- 3 ... W A L 記憶部
- 4 ... W A L 反映部
- 5 ... W A L 管理部
- 6 ... 排他チェック処理部

10

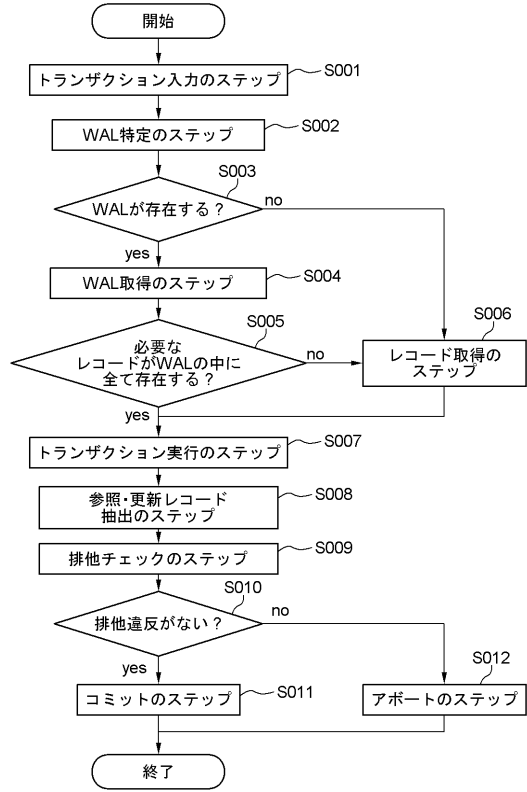
20

30

【図1】



【図2】



【図3】

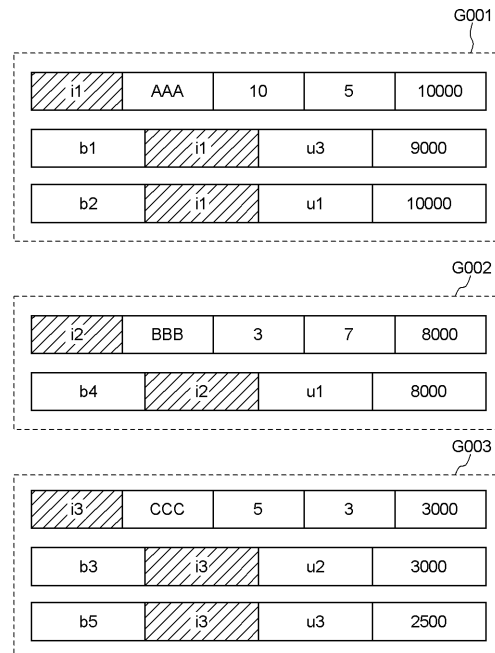
商品テーブル T001

商品ID	商品名	在庫	入札数	最高金額
i1	AAA	10	5	10000
i2	BBB	3	7	8000
i3	CCC	5	3	3000

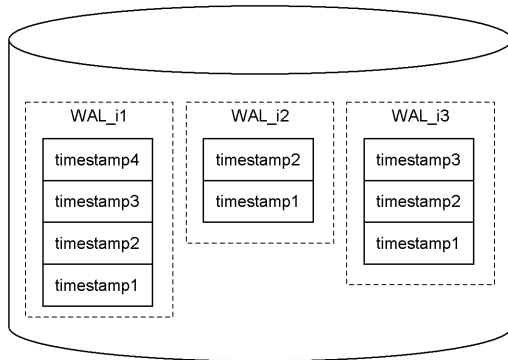
入札テーブル T002

ID	商品ID	ユーザID	入札金額
b1	i1	u3	9000
b2	i1	u1	10000
b3	i3	u2	3000
b4	i2	u1	8000
b5	i3	u3	2500

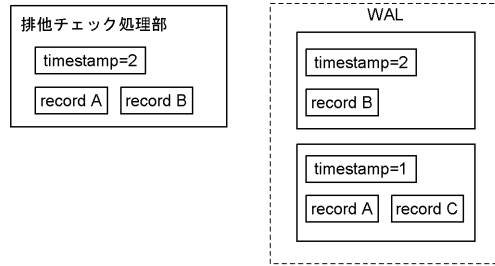
【図4】



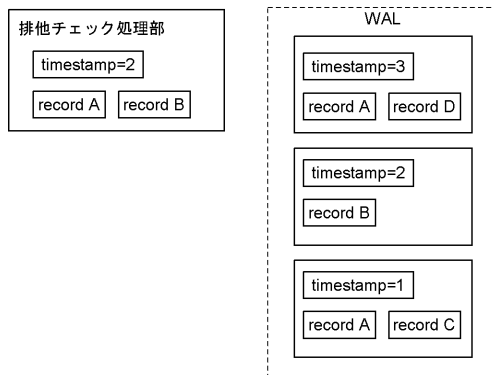
【 図 5 】



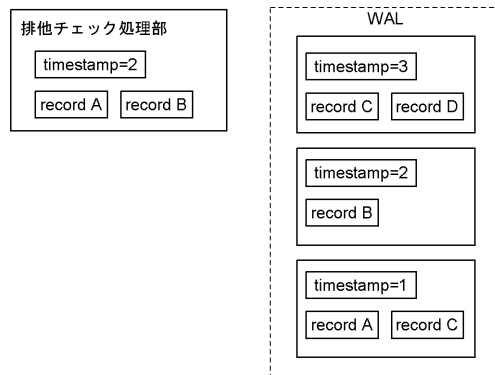
【 図 6 】



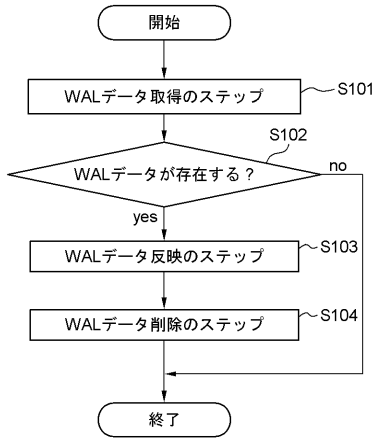
【 図 7 】



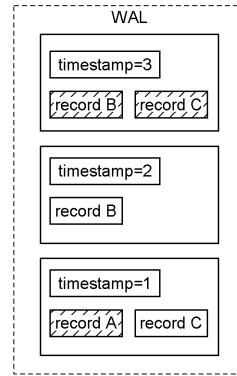
【 図 8 】



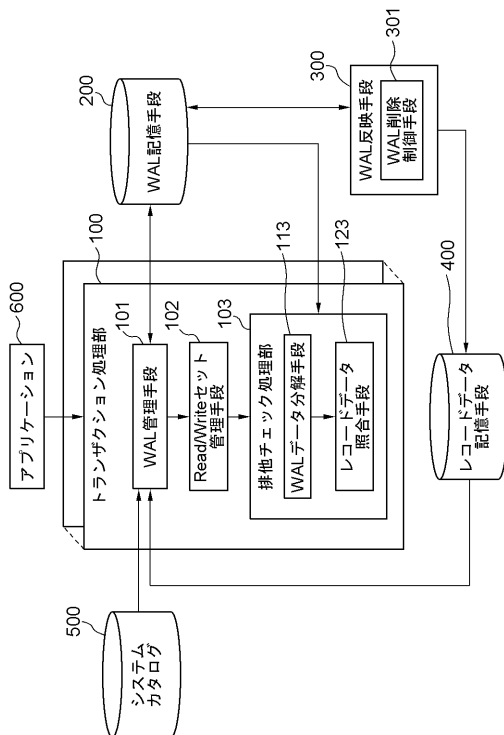
【図9】



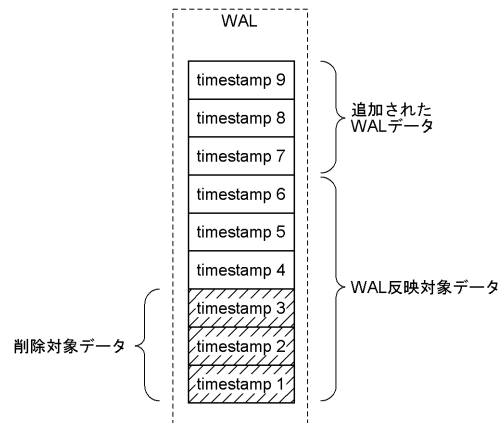
【図10】



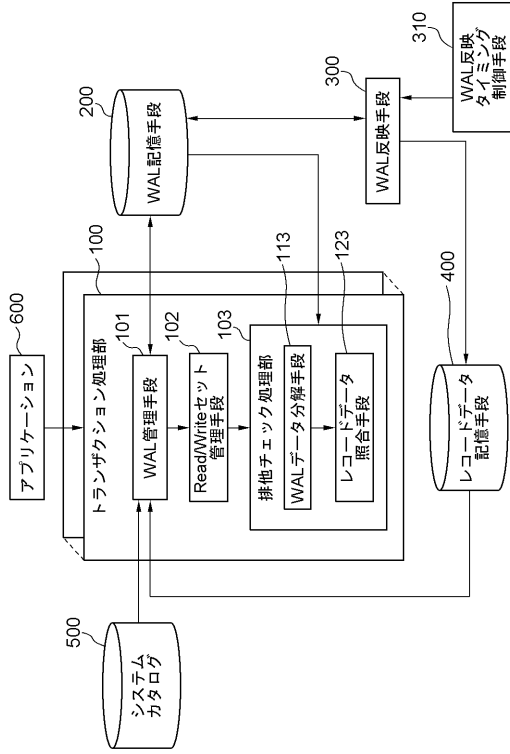
【図11】



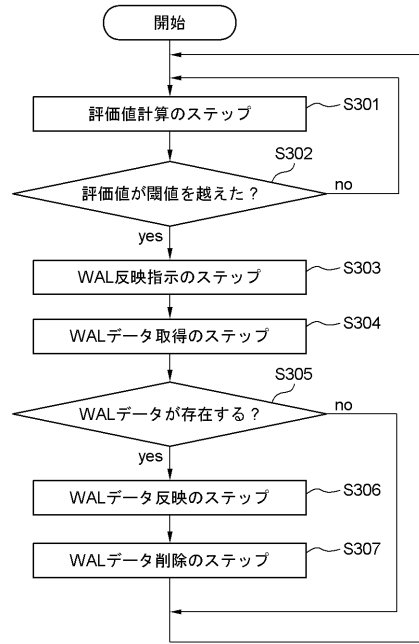
【図12】



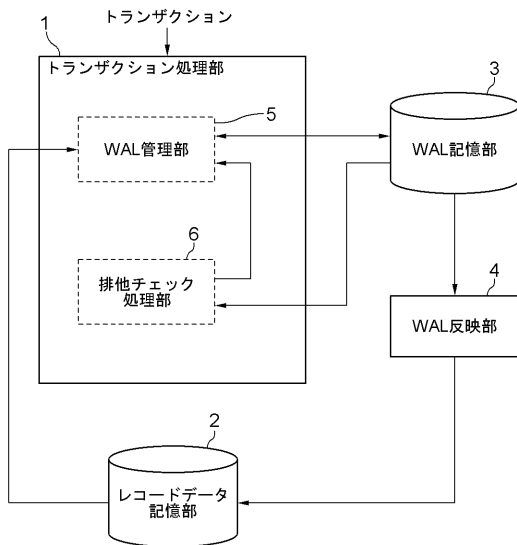
【図13】



【図14】



【図15】



フロントページの続き

(56)参考文献 国際公開第2007/063945(WO, A1)

特開2009-26029(JP, A)

特開2001-101045(JP, A)

祐成光樹 他, "ビッグデータの活用に最適なスケールアウト型 新データベース「InfoFrame Relational Store」", NEC技報, 日本電気株式会社, 2012年 9月 1日, Vol.65, No.2, p.30-33

Junichi Tatemura et al., "Microsharding: a declarative approach to support elastic OLTP workloads", ACM SIGOPS Operating Systems Review, 2012年 1月, Vol.46, No.1, p.4-11

(58)調査した分野(Int.Cl., DB名)

G06F9/46

12/00

17/30

JSTPlus(JDreamIII)