

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2023-541298

(P2023-541298A)

(43)公表日 令和5年9月29日(2023.9.29)

(51)国際特許分類	F I	テーマコード(参考)
G 0 6 F 9/50 (2006.01)	G 0 6 F 9/50 1 5 0 D	5 B 1 7 5
G 0 6 F 16/182 (2019.01)	G 0 6 F 16/182	
G 0 6 F 16/28 (2019.01)	G 0 6 F 16/28	
G 0 6 F 9/46 (2006.01)	G 0 6 F 9/50 1 5 0 A	
	G 0 6 F 9/46 4 3 0	
	審査請求 有 予備審査請求 未請求 (全63頁)	

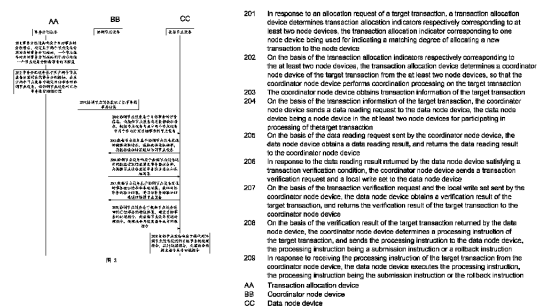
(21)出願番号 特願2023-517375(P2023-517375)
 (86)(22)出願日 令和3年10月26日(2021.10.26)
 (85)翻訳文提出日 令和5年3月15日(2023.3.15)
 (86)国際出願番号 PCT/CN2021/126408
 (87)国際公開番号 WO2022/111188
 (87)国際公開日 令和4年6月2日(2022.6.2)
 (31)優先権主張番号 202011362629.2
 (32)優先日 令和2年11月27日(2020.11.27)
 (33)優先権主張国・地域又は機関 中国(CN)
 (81)指定国・地域 AP(BW,GH,GM,KE,LR,LS,MW,MZ,NA,RW,SD,SL,ST,SZ,TZ,UG,ZM,ZW),EA(AM,AZ,BY,KG,KZ,RU,TJ,TM),EP(AL,A T,BE,BG,CH,CY,CZ,DE,DK,EE,ES,FI,FR,GB,GR,HR,HU,IE,IS,IT,LT,LU,LV,MC,最終頁に続く

(71)出願人 517392436
 騰訊科技(深セン)有限公司
 中華人民共和国518057 広東省深セン市南山区高新区科技中一路騰訊大廈35層
 (74)代理人 100110364
 弁理士 実広 信哉
 (74)代理人 100150197
 弁理士 松尾 直樹
 (72)発明者 李 海 翔
 中華人民共和国518057 広東省深セン市南山区高新区科技中一路騰訊大廈35層
 Fターム(参考) 5B175 AA01 EA03

(54)【発明の名称】 トランザクション処理方法、システム、装置、機器、及びプログラム

(57)【要約】

トランザクション処理方法、システム、装置、機器、記憶媒体及びプログラム製品は、データベースの技術分野に属する。方法は、ターゲットトランザクションの割り当て要求にตอบสนองして、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定するステップ(201)と、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、少なくとも2つのノード機器のうちターゲットトランザクションの協調ノード機器を決定し、協調ノード機器によってターゲットトランザクションを協調処理するステップ(202)と、を含む。このような方式に基づいて、個々の協調ノード機器はいずれも分散化された機器としてトランザクションを協調処理することを可能にすることにより、トランザクションをノード間で処理することを可能にし、トランザクションの処理効率を向上させることに有利であり、トランザクション処理の信頼性が比較的高く、データベースシステムのシステム性能を高めることに有利である。



【特許請求の範囲】**【請求項 1】**

トランザクション処理方法であって、前記方法はトランザクション割り当て機器に応用され、前記トランザクション割り当て機器は分散型データベースシステム中にあり、前記分散型データベースシステムにおいて同一の記憶システムを共有する少なくとも2つのノード機器がさらに含まれ、前記方法は、

ターゲットトランザクションの割り当て要求に応答して、前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定するステップであって、1つのノード機器に対応するトランザクション割り当て指標は前記1つのノード機器に新たなトランザクションを割り当てるマッチング度を指示することに用いられる、ステップと、

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの協調ノード機器を決定し、前記協調ノード機器によって前記ターゲットトランザクションを協調処理するステップと、を含む、トランザクション処理方法。

【請求項 2】

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定する前記ステップは、

トランザクション割り当てモードを決定するステップであって、前記トランザクション割り当てモードはトランザクションのビジー度に基づく割り当て、機器のビジー度に基づく割り当て、及びハイブリッドのビジー度に基づく割り当てのうちのいずれか1つを含む、ステップと、

前記トランザクション割り当てモードによって指示された決定方式に基づき、前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定するステップと、を含む、請求項 1 に記載の方法。

【請求項 3】

前記トランザクション割り当てモードはハイブリッドのビジー度に基づく割り当てを含み、前記トランザクション割り当てモードによって指示された決定方式に基づき、前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定する前記ステップは、

第1ノード機器のトランザクション処理数、前記第1ノード機器の機器リソース使用率、トランザクション処理数の重み、機器リソース使用率の重み、及び重み調節パラメータに基づいて、前記第1ノード機器に対応するトランザクション割り当て指標を決定するステップであって、前記第1ノード機器は前記少なくとも2つのノード機器のうちのいずれか1つのノード機器である、ステップを含む、請求項 2 に記載の方法。

【請求項 4】

前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの協調ノード機器を決定する前記ステップの後に、前記方法は、

前記協調ノード機器の機器識別情報を前記割り当て要求を發した端末に送信するステップであって、前記端末は前記協調ノード機器の機器識別情報に基づき、前記ターゲットトランザクションのトランザクション情報を前記協調ノード機器に送信し、前記協調ノード機器によって前記トランザクション情報に基づいて前記ターゲットトランザクションを協調処理することに用いられる、ステップをさらに含む、請求項 1 ~ 3 のいずれか一項に記載の方法。

【請求項 5】

前記分散型データベースシステムはキー値型のデータ記憶フォーマットとセグメントページ型のデータ記憶フォーマットをサポートする、請求項 1 ~ 3 のいずれか一項に記載の方法。

【請求項 6】

トランザクション処理方法であって、前記方法は協調ノード機器に応用され、前記協調

10

20

30

40

50

ノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定され、前記方法は、

前記ターゲットトランザクションのトランザクション情報を取得するステップと、

前記ターゲットトランザクションのトランザクション情報に基づいて、データノード機器にデータ読み取り要求を送信するステップであって、前記データノード機器は前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの処理に関与することに用いられるノード機器である、ステップと、

前記データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信するステップと、

前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定し、前記データノード機器に前記処理命令を送信するステップであって、前記処理命令はコミット命令、又はアポート命令であり、前記データノード機器は前記処理命令を実行することに用いられる、ステップと、を含む、トランザクション処理方法。

【請求項7】

前記データ読み取り結果は第2論理ライフサイクルを携えており、前記第2論理ライフサイクルは前記データノード機器によって前記データ読み取り要求が携えている前記ターゲットトランザクションの第1論理ライフサイクルに基づき決定され、前記第1論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成され、

前記データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求を送信する前記ステップは、

前記第1論理ライフサイクルのタイムスタンプ下限と前記第2論理ライフサイクルのタイムスタンプ下限における最大値を前記ターゲットトランザクションの第3論理ライフサイクルのタイムスタンプ下限として用い、前記第1論理ライフサイクルのタイムスタンプ上限と前記第2論理ライフサイクルのタイムスタンプ上限における最小値を前記ターゲットトランザクションの第3論理ライフサイクルのタイムスタンプ上限として用いるステップと、

前記第3論理ライフサイクルが有効であることに応答して、前記データノード機器に前記第3論理ライフサイクルを携えているトランザクション認証要求を送信するステップであって、前記第3論理ライフサイクルが有効であることは、前記第3論理ライフサイクルのタイムスタンプ下限が前記第3論理ライフサイクルのタイムスタンプ上限よりも小さいことを指示することに用いられる、ステップと、を含む、請求項6に記載の方法。

【請求項8】

前記データノード機器の数は少なくとも2つであり、前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定する前記ステップは、

前記少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果において認証に合格しなかったことを指示することに用いられる認証結果が存在することに応答して、前記アポート命令を前記ターゲットトランザクションの処理命令として用いるステップと、

前記少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果はいずれも認証に合格したことを指示することに応答して、前記少なくとも2つの認証結果が携えている論理ライフサイクルの共通部分をターゲット論理ライフサイクルとして用いるステップと、

前記ターゲット論理ライフサイクルが有効であることに応答して、前記コミット命令を前記ターゲットトランザクションの処理命令として用い、前記ターゲット論理ライフサイ

10

20

30

40

50

クルが無効であることに応答して、前記アポート命令を前記ターゲットトランザクションの処理命令として用いるステップと、を含む、請求項 6 又は 7 のいずれか一項に記載の方法。

【請求項 9】

トランザクション処理方法であって、前記方法はデータノード機器に応用され、前記データノード機器は同一の記憶システムを共有する少なくとも 2 つのノード機器のうちターゲットトランザクションの処理に関与することに用いられるノード機器であり、前記方法は、

協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得し、前記データ読み取り結果を前記協調ノード機器に返信するステップであって、前記協調ノード機器は前記少なくとも 2 つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定される、ステップと、

前記協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得し、前記ターゲットトランザクションの認証結果を前記協調ノード機器に返信するステップと、

前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受信したことに応答して、前記処理命令を実行するステップであって、前記処理命令はコミット命令、又はアポート命令である、ステップと、を含む、トランザクション処理方法。

【請求項 10】

前記データ読み取り要求は前記ターゲットトランザクションの第 1 論理ライフサイクルを携えており、前記第 1 論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成され、協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得する前記ステップは、

前記第 1 論理ライフサイクルに基づいて、前記データ読み取り要求によって指示された読み取り対象のデータ項目の可視バージョンデータを決定するステップと、

前記可視バージョンデータの作成タイムスタンプ、及び前記第 1 論理ライフサイクルに基づいて、前記ターゲットトランザクションの第 2 論理ライフサイクルを決定するステップと、

前記第 2 論理ライフサイクル、及び前記可視バージョンデータが運ばれた結果を前記データ読み取り結果として用いるステップと、を含む、請求項 9 に記載の方法。

【請求項 11】

前記トランザクション認証要求は前記ターゲットトランザクションの第 3 論理ライフサイクルを携えており、前記第 3 論理ライフサイクルは前記協調ノード機器によって前記第 1 論理ライフサイクル、及び前記第 2 論理ライフサイクルに基づいて決定された有効論理ライフサイクルであり、前記協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得する前記ステップは、

前記第 3 論理ライフサイクルのタイムスタンプ下限と前記第 2 論理ライフサイクルのタイムスタンプ下限における最大値を前記ターゲットトランザクションの第 4 論理ライフサイクルのタイムスタンプ下限として用い、前記第 3 論理ライフサイクルのタイムスタンプ上限と前記第 2 論理ライフサイクルのタイムスタンプ上限における最小値を前記ターゲットトランザクションの第 4 論理ライフサイクルのタイムスタンプ上限として用いるステップと、

前記第 4 論理ライフサイクルが有効であることに応答して、前記ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び前記第 4 論理ライフサイクルに基づいて、前記ターゲットトランザクションの第 5 論理ライフサイクルを決定するステップと、

前記第 5 論理ライフサイクルが有効であることに応答して、認証に合格したことを指示することに用いられる認証結果を前記ターゲットトランザクションの認証結果として用い、前記第 5 論理ライフサイクルが無効であることに応答して、認証に合格しなかったこと

10

20

30

40

50

を指示することに用いられる認証結果を前記ターゲットトランザクションの認証結果として用いるステップと、を含む、請求項 10 に記載の方法。

【請求項 12】

1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は前記1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプを含み、前記1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプは前記1つの書き込み対象のデータ項目を読み取った各読み取りトランザクションの論理コミットタイムスタンプにおける最大値を指示することに用いられ、前記ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び前記第4論理ライフサイクルに基づいて、前記ターゲットトランザクションの第5論理ライフサイクルを決定する前記ステップは、

10

前記各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び前記第4論理ライフサイクルに基づいて、前記ターゲットトランザクションの第5論理ライフサイクルを決定するステップであって、前記第5論理ライフサイクルのタイムスタンプ下限は前記各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプにおける最大値よりも大きい、ステップを含む、請求項 11 に記載の方法。

【請求項 13】

1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は前記1つの書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプを含み、前記ターゲット読み取りトランザクションはローカル認証に合格するか、又はコミット段階にある読み取りトランザクションであり、前記ターゲット読み取りトランザクションの終了タイムスタンプは前記ターゲット読み取りトランザクションの論理ライフサイクルのタイムスタンプ上限であり、前記ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び前記第4論理ライフサイクルに基づいて、前記ターゲットトランザクションの第5論理ライフサイクルを決定する前記ステップは、

20

前記各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプ、及び前記第4論理ライフサイクルに基づいて、前記ターゲットトランザクションの第5論理ライフサイクルを決定するステップであって、前記第5論理ライフサイクルのタイムスタンプ下限は前記各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプにおける最大値よりも大きい、ステップを含む、請求項 11 に記載の方法。

30

【請求項 14】

トランザクション処理システムであって、前記トランザクション処理システムは協調ノード機器と、データノード機器とを含み、前記協調ノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定され、前記データノード機器は前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの処理に関与することに用いられるノード機器であり、

40

前記協調ノード機器は、前記ターゲットトランザクションのトランザクション情報を取得することと、前記ターゲットトランザクションのトランザクション情報に基づいて、前記データノード機器にデータ読み取り要求を送信することと、に用いられ、

前記データノード機器は、前記協調ノード機器から送信された前記データ読み取り要求に基づいて、データ読み取り結果を取得し、前記データ読み取り結果を前記協調ノード機器に返信することに用いられ、

前記協調ノード機器はさらに、前記データノード機器から返信された前記データ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信することに用いられ、

前記データノード機器はさらに、前記協調ノード機器から送信された前記トランザクシ

50

ョン認証要求、及び前記ローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得し、前記ターゲットトランザクションの認証結果を前記協調ノード機器に返信することに用いられ、

前記協調ノード機器はさらに、前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定し、前記データノード機器に前記処理命令を送信することに用いられ、前記処理命令はコミット命令、又はアポート命令であり、

前記データノード機器はさらに、前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受信したことに応答して、前記処理命令を実行することに用いられる、トランザクション処理システム。

10

【請求項 15】

トランザクション処理装置であって、前記装置は、

ターゲットトランザクションの割り当て要求に応答して、同一の記憶システムを共有する少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定することに用いられる第1決定ユニットであって、1つのノード機器に対応するトランザクション割り当て指標は前記1つのノード機器に新たなトランザクションを割り当てるマッチング度を指示することに用いられる、第1決定ユニットと、

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの協調ノード機器を決定し、前記協調ノード機器によって前記ターゲットトランザクションを協調処理することに用いられる第2決定ユニットと、を含む、トランザクション処理装置。

20

【請求項 16】

トランザクション処理装置であって、前記装置は、

ターゲットトランザクションのトランザクション情報を取得することに用いられる取得ユニットと、

前記ターゲットトランザクションのトランザクション情報に基づいて、データノード機器にデータ読み取り要求を送信することに用いられる第1送信ユニットであって、前記データノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうち前記ターゲットトランザクションの処理に関与することに用いられるノード機器である、第1送信ユニットと、

30

前記データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求及び、ローカル書き込みセットを送信することに用いられる第2送信ユニットと、

前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定することに用いられる決定ユニットと、

前記データノード機器に前記処理命令を送信することに用いられる第3送信ユニットであって、前記処理命令はコミット命令、又はアポート命令であり、前記データノード機器は前記処理命令を実行することに用いられる、第3送信ユニットと、を含む、トランザクション処理装置。

40

【請求項 17】

トランザクション処理装置であって、前記装置は、

協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得することに用いられる第1取得ユニットであって、前記協調ノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定される、第1取得ユニットと、

前記データ読み取り結果を前記協調ノード機器に返信することに用いられる返信ユニッ

50

トと、

前記協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得することに用いられる第2取得ユニットと、

前記返信ユニットはさらに、前記ターゲットトランザクションの認証結果を前記協調ノード機器に返信することに用いられ、

前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受信したことに応答して、前記処理命令を実行することに用いられる実行ユニットであって、前記処理命令はコミット命令、又はアポート命令である、実行ユニットと、を含むトランザクション処理装置。

10

【請求項18】

コンピュータ機器であって、前記コンピュータ機器はプロセッサと、メモリとを含み、前記メモリにおいて少なくとも1つのコンピュータプログラムが記憶され、前記少なくとも1つのコンピュータプログラムは前記プロセッサによってロードされて実行され、前記コンピュータ機器に、請求項1～5のいずれか一項に記載のトランザクション処理方法、又は請求項6～8のいずれか一項に記載のトランザクション処理方法、又は請求項9～13のいずれか一項に記載のトランザクション処理方法を実現させる、コンピュータ機器。

【請求項19】

非一時的コンピュータ可読記憶媒体であって、前記非一時的コンピュータ可読記憶媒体において少なくとも1つのコンピュータプログラムが記憶され、前記少なくとも1つのコンピュータプログラムはプロセッサによってロードされて実行され、コンピュータに、請求項1～5のいずれか一項に記載のトランザクション処理方法、又は請求項6～8のいずれか一項に記載のトランザクション処理方法、又は請求項9～13のいずれか一項に記載のトランザクション処理方法を実現させる、非一時的コンピュータ可読記憶媒体。

20

【請求項20】

コンピュータプログラム製品であって、前記コンピュータプログラム製品はコンピュータ命令を含み、前記コンピュータ命令はコンピュータ可読記憶媒体に記憶され、コンピュータ機器のプロセッサは前記コンピュータ可読記憶媒体から前記コンピュータ命令を読み取り、前記プロセッサが前記コンピュータ命令を実行することにより、前記コンピュータ機器が請求項1～5のいずれか一項に記載のトランザクション処理方法、請求項6～8のいずれか一項に記載のトランザクション処理方法、又は請求項9～13のいずれか一項に記載のトランザクション処理方法を実行する、コンピュータプログラム製品。

30

【発明の詳細な説明】

【技術分野】

【0001】

本願の実施例はデータベースの技術分野に関し、特にトランザクション処理方法、システム、装置、機器、記憶媒体及びプログラム製品に関する。

【0002】

本願は、2020年11月27日に提出された、出願番号が第202011362629.2号、発明の名称が「トランザクション処理方法、機器及びコンピュータ可読記憶媒体」の中国特許出願の優先権を主張し、その全内容は引用により本願に組み込まれている。

40

【背景技術】

【0003】

データベース技術の発展に伴って、ビッグデータ、クラウドコンピューティングなどのビジネスシナリオに適應することを可能にするために、分散型データベースシステムが徐々に普及している。複数種の分散型データベースシステムにおいて、共有記憶(sharedisk)アーキテクチャに基づく分散型データベースシステムは主流のシステムとなっている。

【0004】

50

現在、share-diskアーキテクチャに基づく分散型データベースシステムにおいて、データ項目の分布状況に基づきトランザクションを割り当てて、あるデータ項目に関するトランザクションを該データ項目にサービスされる固定ノード機器に割り当てて独立して処理する。これに基づいて、トランザクションの処理効率が大幅に制限されてしまい、トランザクション処理の信頼性が比較的悪くなる。

【発明の概要】

【発明が解決しようとする課題】

【0005】

本願の実施例はトランザクション処理方法、システム、装置、機器、記憶媒体及びプログラム製品を提供し、トランザクションの処理効率を向上させることに用いることができる。

10

【課題を解決するための手段】

【0006】

一態様では、本願の実施例はトランザクション処理方法を提供し、前記方法はトランザクション割り当て機器に応用され、前記トランザクション割り当て機器は分散型データベースシステム中にあり、前記分散型データベースシステムにおいて同一の記憶システムを共有する少なくとも2つのノード機器がさらに含まれ、前記方法は、

ターゲットトランザクションの割り当て要求に応答して、前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定するステップであって、1つのノード機器に対応するトランザクション割り当て指標は前記1つのノード機器に新たなトランザクションを割り当てるマッチング度を指示することに用いられる、ステップと、

20

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの協調ノード機器を決定し、前記協調ノード機器によって前記ターゲットトランザクションを協調処理するステップと、を含む。

【0007】

トランザクション処理方法をさらに提供し、前記方法は協調ノード機器に応用され、前記協調ノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定され、前記方法は、

30

前記ターゲットトランザクションのトランザクション情報を取得するステップと、

前記ターゲットトランザクションのトランザクション情報に基づいて、データノード機器にデータ読み取り要求を送信するステップであって、前記データノード機器は前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの処理に関与することに用いられるノード機器である、ステップと、

前記データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信するステップと、

40

前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定し、前記データノード機器に前記処理命令を送信するステップであって、前記処理命令はコミット命令、又はアポート命令であり、前記データノード機器は前記処理命令を実行することに用いられる、ステップと、を含む。

【0008】

トランザクション処理方法をさらに提供し、前記方法はデータノード機器に応用され、前記データノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションの処理に関与することに用いられるノード機器であり、前記方法は、

50

協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得し、前記データ読み取り結果を前記協調ノード機器に返信するステップであって、前記協調ノード機器は前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定される、ステップと、

前記協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得し、前記ターゲットトランザクションの認証結果を前記協調ノード機器に返信するステップと、

前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受信したことに応答して、前記処理命令を実行するステップであって、前記処理命令はコミット命令、又はアポート命令である、ステップと、を含む。

10

【0009】

別の態様では、トランザクション処理システムを提供し、前記トランザクション処理システムは協調ノード機器と、データノード機器とを含み、前記協調ノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定され、前記データノード機器は前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの処理に関与することに用いられるノード機器であり、

前記協調ノード機器は、前記ターゲットトランザクションのトランザクション情報を取得することと、前記ターゲットトランザクションのトランザクション情報に基づいて、前記データノード機器にデータ読み取り要求を送信することと、に用いられ、

20

前記データノード機器は、前記協調ノード機器から送信された前記データ読み取り要求に基づいて、データ読み取り結果を取得し、前記データ読み取り結果を前記協調ノード機器に返信することに用いられ、

前記協調ノード機器はさらに、前記データノード機器から返信された前記データ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信することに用いられ、

前記データノード機器はさらに、前記協調ノード機器から送信された前記トランザクション認証要求、及び前記ローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得し、前記ターゲットトランザクションの認証結果を前記協調ノード機器に返信することに用いられ、

30

前記協調ノード機器はさらに、前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定し、前記データノード機器に前記処理命令を送信することに用いられ、前記処理命令はコミット命令、又はアポート命令であり、

前記データノード機器はさらに、前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受信したことに応答して、前記処理命令を実行することに用いられる。

【0010】

別の態様では、トランザクション処理装置を提供し、前記装置は、

40

ターゲットトランザクションの割り当て要求に応答して、同一の記憶システムを共有する少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定することに用いられる第1決定ユニットであって、1つのノード機器に対応するトランザクション割り当て指標は前記1つのノード機器に新たなトランザクションを割り当てるマッチング度を指示することに用いられる、第1決定ユニットと、

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの協調ノード機器を決定し、前記協調ノード機器によって前記ターゲットトランザクションを協調処理することに用いられる第2決定ユニットと、を含む。

【0011】

50

トランザクション処理装置をさらに提供し、前記装置は、
ターゲットトランザクションのトランザクション情報を取得することに用いられる取得
ユニットと、

前記ターゲットトランザクションのトランザクション情報に基づいて、データノード機
器にデータ読み取り要求を送信することに用いられる第1送信ユニットであって、前記デ
ータノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうち前記
ターゲットトランザクションの処理に關与することに用いられるノード機器である、第1
送信ユニットと、

前記データノード機器から返信されたデータ読み取り結果がトランザクション認証条件
を満たすことに応答して、前記データノード機器にトランザクション認証要求、及びロー
カル書き込みセットを送信することに用いられる第2送信ユニットと、

前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基
づいて、前記ターゲットトランザクションの処理命令を決定することに用いられる決定ユ
ニットと、

前記データノード機器に前記処理命令を送信することに用いられる第3送信ユニットで
あって、前記処理命令はコミット命令、又はアポート命令であり、前記データノード機
器は前記処理命令を実行することに用いられる、第3送信ユニットと、を含む。

【0012】

トランザクション処理装置をさらに提供し、前記装置は、

協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を
取得することに用いられる第1取得ユニットであって、前記協調ノード機器は同一の記憶
システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協
調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも2
つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定され
る、第1取得ユニットと、

前記データ読み取り結果を前記協調ノード機器に返信することに用いられる返信ユニ
ットと、

前記協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込み
セットに基づいて、前記ターゲットトランザクションの認証結果を取得することに用いら
れる第2取得ユニットと、

前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受
信したことに応答して、前記処理命令を実行することに用いられる実行ユニットであ
って、前記処理命令はコミット命令、又はアポート命令である、実行ユニットと、を含
み、

前記返信ユニットはさらに、前記ターゲットトランザクションの認証結果を前記協調
ノード機器に返信することに用いられる。

【0013】

別の態様では、コンピュータ機器を提供し、前記コンピュータ機器はプロセッサと、メ
モリとを含み、前記メモリにおいて少なくとも1つのコンピュータプログラムが記憶され
、前記少なくとも1つのコンピュータプログラムは前記プロセッサによってロードされて
実行され、前記コンピュータ機器に前記のいずれか一項に記載のトランザクション処理方
法を実現させる。

【0014】

別の態様では、非一時的コンピュータ可読記憶媒体をさらに提供し、前記非一時的コン
ピュータ可読記憶媒体において少なくとも1つのコンピュータプログラムが記憶され、前
記少なくとも1つのコンピュータプログラムはプロセッサによってロードされて実行され
、コンピュータに前記のいずれか一項に記載のトランザクション処理方法を実現させる。

【0015】

別の態様では、コンピュータプログラム製品、又はコンピュータプログラムをさらに提
供し、前記コンピュータプログラム製品、又はコンピュータプログラムはコンピュータ命
令を含み、前記コンピュータ命令はコンピュータ可読記憶媒体に記憶され、コンピュータ

10

20

30

40

50

機器のプロセッサは前記コンピュータ可読記憶媒体から前記コンピュータ命令を読み取り、前記プロセッサは前記コンピュータ命令を実行することにより、前記コンピュータ機器が前記のいずれか一項に記載のトランザクション処理方法を実行する。

【0016】

本願の実施例において、各ノード機器のそれぞれに対応するトランザクション割り当て指標に基づきターゲットトランザクションを協調処理することに用いられる協調ノード機器を決定し、トランザクションの割り当て過程はトランザクションに関するデータ項目を考慮する必要がなく、データ項目の分布状況を考慮する必要もない。このような方式に基づいて、個々のノード機器はいずれも分散化された機器としてトランザクションを協調処理することを可能にすることにより、トランザクションをノード間で処理することを可能にし、トランザクションの処理効率を向上させることに有利となり、トランザクション処理の信頼性が比較的高くなり、データベースシステムのシステム性能を高めることに有利である。

10

【図面の簡単な説明】

【0017】

【図1】本願の実施例が提供するトランザクション処理方法の実施環境の模式図である。

【図2】本願の実施例が提供するトランザクション処理方法のフローチャートである。

【図3】本願の実施例が提供するトランザクションログのフォーマットの模式図である。

【図4】本願の実施例が提供するトランザクションログのフォーマットの模式図である。

【図5】本願の実施例が提供するトランザクション処理装置の模式図である。

20

【図6】本願の実施例が提供するトランザクション処理装置の模式図である。

【図7】本願の実施例が提供するトランザクション処理装置の模式図である。

【図8】本願の実施例が提供するコンピュータ機器の構造模式図である。

【発明を実施するための形態】

【0018】

本願の目的、技術的解決手段及び利点をより明確にするために、以下、図面を組み合わせることで本願の実施形態をさらに詳細に記述する。

【0019】

説明する必要があるように、本願の明細書及び特許請求の範囲における用語「第1」、「第2」などは、類似の対象を区別することに用いられるものであり、特定の順序又は前後順番を記述することに用いられるものではない。理解すべきであるように、ここで記述された本願の実施例がここで図示又は記述されたそれらのもの以外の順序で実施されることを可能にするように、このように使用されるデータは適切な場合に相互に交換することができる。以下の例示的な実施例において記述された実施形態は本願と一致するすべての実施形態を代表するものではない。逆に、これらは添付の特許請求の範囲において詳述された本願のいくつかの態様と一致する装置及び方法の例に過ぎない。

30

【0020】

いくつかの実施例において、本願の実施例に関する分散型データベースシステムは、共有記憶(share-disk)アーキテクチャに基づく分散型データベースシステムであり、共有記憶アーキテクチャに基づく分散型データベースシステムにおいて少なくとも2つのノード機器が含まれ、該少なくとも2つのノード機器は自体のローカルの内部メモリエリアを有し、ネットワーク通信メカニズムによって同一の記憶システムに直接アクセスし、すなわち少なくとも2つのノード機器は同一の記憶システムを共有する。例えば、同一のHDFS(Hadoop Distributed File System、分散型ファイルシステム)を共有する。少なくとも2つのノード機器によって共有された記憶システムにおいて複数のデータテーブルを記憶することができ、個々のデータテーブルは1つ又は複数のデータ項目を記憶することに用いることができる。

40

【0021】

論理という角度から見ると、分散型データベースシステムにおけるノード機器は、協調ノード機器及びデータノード機器の2つの種類のロールに分割することができる。ここで

50

、協調ノード機器は主に生産、配信処理計画、及び分散型トランザクションの協調を担当し、データノード機器は主に協調ノード機器から送信された処理計画を受信し、相応なトランザクションを実行して協調ノード機器にトランザクションに関する関連データを返信することを担当する。

【0022】

分散型データベースシステムにおいて、最小の操作実行ユニットはトランザクションであり、トランザクションが複数のデータノード機器上のデータ項目を操作する必要があるか否かに応じて、トランザクションは分散型トランザクション及びローカルトランザクションの2つの種類に分割することができる。この2つの種類の異なるトランザクションに対しては、それぞれ異なる実行プロセスを採用して、ネットワーク通信オーバーヘッドをできるだけ減少させ、トランザクション処理効率を高めることができる。ここで、分散型トランザクションは複数のデータノード機器間で読み書き操作を実行する必要があるトランザクションを表し、すなわちトランザクションは複数のデータノード機器上のデータ項目を操作する必要がある。たとえば、トランザクションTはデータノード機器RM1、RM2、RM3上のデータ項目を操作する必要がある場合、該トランザクションTは1つの分散型トランザクションである。ローカルトランザクションは単一のデータノード機器上のデータ項目のみを操作する必要があるトランザクションを表し、例えば、トランザクションTはRM1上のデータ項目のみを操作する必要があるれば、該トランザクションTは1つのローカルトランザクションである。

10

【0023】

図1は本願の実施例が提供するトランザクション処理方法の実施環境の模式図である。図1に参照されるように、本願の実施例はshare-diskフレームワークに基づく分散型データベースシステムに適用されてもよく、該分散型データベースシステムにおいてゲートウェイサーバ101と、トランザクション割り当て機器102と、分散型記憶クラスタ103と、グローバルタイムスタンプ生成クラスタ104とが含まれてもよい。分散型記憶クラスタ103においてm(mは2以上の整数である)個のノード機器が含まれ、該m個のノード機器は同一の記憶システムを共有する。

20

【0024】

ゲートウェイサーバ101は外部の読み書き要求を受信し、かつ読み書き要求に対応する読み書きトランザクションをトランザクション割り当て機器102、又は、分散型記憶クラスタ103に配信することに用いられる。たとえば、ユーザは端末上のアプリケーションクライアント端末にログインした後、アプリケーションクライアント端末をトリガーして読み書き要求を生成し、分散型データベースシステムによって提供されたAPI(Application Programming Interface、アプリケーションプログラムプログラミングインタフェース)を呼び出して該読み書き要求に対応する読み書きトランザクションをゲートウェイサーバ101に送信する。

30

【0025】

いくつかの実施例において、ゲートウェイサーバ101は分散型記憶クラスタ103におけるいずれか1つのノード機器と同一の物理マシン上に合併されてもよく、すなわち、あるノード機器をゲートウェイサーバ101として機能させる。

40

【0026】

いくつかの実施例において、アプリケーションクライアント端末が所在する端末は分散型データベースシステムにおけるトランザクション割り当て機器102、及び、分散型記憶クラスタ103と通信接続を直接確立することを可能にする。このような場合に、分散型データベースシステムにおいてゲートウェイサーバ101が存在しなくてもよい。

【0027】

トランザクション割り当て機器102は新たなトランザクションに適切なノード機器を協調ノード機器として割り当てることに用いられる。例示的な実施例において、トランザクション割り当て機器は分散型協調システム(たとえばZooKeeper)にある。分散型協調システムはゲートウェイサーバ101、分散型記憶クラスタ103、及びグロー

50

バルタイムスタンプ生成クラスタ104のうちの少なくとも1つを管理することに用いられてもよい。選択的に、技術者は、端末上のスケジューラ(scheduler)を介して該分散型協調システムにアクセスし、それによりフロントエンドのスケジューラに基づいてバックエンドの分散型協調システムを制御して、各クラスタ又はサーバに対する管理を実現することができる。例えば、技術者は、スケジューラを介してZooKeeperを制御してあるノード機器を分散型記憶クラスタ103から削除することができ、すなわちあるノード機器を無効にすることができる。

【0028】

分散型記憶クラスタ103はデータノード機器と、協調ノード機器とを含んでもよく、個々の協調ノード機器は少なくとも1つのデータノード機器に対応してもよい。データノード機器と協調ノード機器の分割は異なるトランザクションに対するものである。ある分散型トランザクションを例とすると、分散型トランザクションを発生したノード機器は協調ノード機器と呼ばれてもよく、分散型トランザクションに関する他のノード機器はデータノード機器と呼ばれる。データノード機器又は協調ノード機器の数は1つ又は複数であってもよく、本願の実施例は分散型記憶クラスタ103におけるデータノード機器又は協調ノード機器の数を具体的に限定しない。

10

【0029】

本願の実施例が提供する分散型データベースシステムにおいては、グローバルトランザクションマネージャがないため、該システムにおいてXA(extended Architecture、X/Open組織の分散型トランザクション規範)/2PC(Two-Phase Commit、2相コミット)技術を採用してノード間のトランザクション(分散型トランザクション)をサポートし、ノード間での書き込み操作時のデータの原子性及び一致性を確保することができる。このとき、協調ノード機器は2PCアルゴリズムにおけるコーディネーターとして機能することに用いられ、該協調ノード機器と対応する各データノード機器は2PCアルゴリズムにおける関与者として機能することに用いられる。

20

【0030】

個々のデータノード機器、又は、協調ノード機器はスタンドアロン機器であってもよく、マスタースタンバイ構造(すなわちワンマスターマルチスタンバイクラスタ)を採用してもよい。図1に示すように、ノード機器(データノード機器又は協調ノード機器)がワンマスターデュアルスタンバイクラスタであることを例として例示すると、個々のノード機器においては1つのマスター、及び2つのスタンバイが含まれ、選択的に、個々のマスター又はスタンバイはいずれもエージェント(agent)機器に対応して配置される。エージェント機器はマスター又はスタンバイと物理的に独立してもよく、もちろん、エージェント機器はさらにマスター又はスタンバイにおける1つのエージェントモジュールとして用いられてもよい。ノード機器1を例とすると、ノード機器1は1つのマスターデータベース、及びエージェント機器(マスターdatabase+agent、マスターDB+agentと略称する)を含み、この他、2つのスタンバイデータベース及びエージェント機器(スタンバイdatabase+agent、スタンバイDB+agentと略称する)をさらに含む。マスターデータベースは上記したマスターであり、スタンバイデータベースは上記したスタンバイである。

30

40

【0031】

グローバルタイムスタンプ生成クラスタ104は分散型トランザクションのグローバルコミットタイムスタンプ(Global Timestamp、Gts)を生成することに用いられ、該分散型トランザクションは複数のデータノード機器に関するトランザクションを指してもよく、例えば分散型読み取りトランザクションは複数のデータノード機器に記憶されたデータの読み取りに関してもよい。さらに例えば、分散型書き込みトランザクションは複数のデータノード機器におけるデータの書き込みに関してもよい。グローバルタイムスタンプ生成クラスタ104は、論理的に1つのシングルポイントとみなされてもよいが、いくつかの実施例において、1マスター3スレーブのアーキテクチャによって

50

より高い可用性を有するサービスを提供でき、クラスタの形式を採用して該グローバルコミットタイムスタンプの生成を実現することにより、シングルポイントの故障を防止し、それによりシングルポイントのボトルネック問題を避けることができる。

【0032】

選択的に、グローバルコミットタイムスタンプは、分散型データベースシステムにおいてグローバルに一意的でかつ単調に遞増する1つのタイムスタンプ識別子であり、個々のトランザクションのグローバルコミット順序をマークする。これによりトランザクション間の実際の時間上の前後関係（トランザクションの全順序関係）を反映することに用いることを可能にし、グローバルコミットタイムスタンプは物理クロック、論理クロック又はハイブリッド物理クロックのうち少なくとも1つを採用してもよく、本願の実施例はグローバルコミットタイムスタンプのタイプを具体的に限定しない。

10

【0033】

いくつかの実施例において、該グローバルタイムスタンプ生成クラスタ104は物理的に独立してもよく、分散型協調システム（例えばZooKeeper）と合併されてもよい。

【0034】

上記図1は軽量なトランザクション処理を提供するアーキテクチャ図に過ぎず、share-diskアーキテクチャに基づく分散型データベースシステムの1つの例示的な記述である。いくつかの実施例において、上記ゲートウェイサーバ101、トランザクション割り当て機器102、分散型記憶クラスタ103、及びグローバルタイムスタンプ生成クラスタ104で構成された分散型データベースシステムは、ユーザ端末にデータサービスを提供するサーバとみなされてもよい。該サーバは独立した物理サーバであってもよく、複数の物理サーバで構成されたサーバクラスタ又は分散型システムであってもよく、クラウドサービス、クラウドデータベース、クラウドコンピューティング、クラウド関数、クラウド記憶、ネットワークサービス、クラウド通信、ミドルウェアサービス、ドメイン名サービス、セキュリティサービス、CDN（Content Delivery Network、コンテンツ配信ネットワーク）、及びビッグデータ、並びに人工知能プラットフォームなどの基礎的なクラウドコンピューティングサービスを提供するクラウドサーバであってもよい。選択的に、上記ユーザ端末はスマートフォン、タブレットコンピュータ、ノートパソコン、デスクトップコンピュータ、スマートスピーカー、スマートウォッチなどであってもよいが、これらに限定されない。端末、及びサーバは有線又は無線通信方式によって直接又は間接的に接続されてもよく、本願はここで制限しない。

20

30

【0035】

上記図1に示される実施環境に基づいて、本願の実施例はトランザクション処理方法を提供する。図2に示すように、本願の実施例が提供する方法は以下のステップ201～ステップ209を含む。

【0036】

ステップ201において、トランザクション割り当て機器は、ターゲットトランザクションの割り当て要求に応答して、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定し、1つのノード機器に対応するトランザクション割り当て指標は該1つのノード機器に新たなトランザクションを割り当てるマッチング度を指示することに用いられる。

40

【0037】

トランザクション割り当て機器、及び少なくとも2つのノード機器はいずれも分散型データベースシステム中にあり、かつ少なくとも2つのノード機器は同一の記憶システムを共有する。本願の実施例は分散型データシステムの具体的な構造を限定せず、トランザクション割り当て機器と、同一の記憶システムを共有する少なくとも2つのノード機器とを含むものであればよい。

【0038】

ターゲットトランザクションとは処理対象のトランザクションを指し、ターゲットトラ

50

ンザクションは分散型トランザクションであってもよく、ローカルトランザクションであってもよいが、本願の実施例はこれを限定しない。ターゲットトランザクションの割り当て要求はターゲットトランザクションに適切なノード機器を協調ノード機器として割り当てて、割り当てられた協調ノード機器によって該ターゲットトランザクションを協調処理するように指示することに用いられる。

【0039】

ターゲットトランザクションの割り当て要求は端末によって発され、端末によって発されたターゲットトランザクションの割り当て要求は端末によってトランザクション割り当て機器に直接送信されるか、又はゲートウェイサーバによってトランザクション割り当て機器に転送されるが、本願の実施例はこれを限定しない。端末はユーザと対応するいずれかの電子機器であってもよく、スマートフォン、タブレットコンピュータ、ノートパソコン、デスクトップコンピュータ、スマートスピーカー、又はスマートウォッチのうちの少なくとも1つを含むがこれらに限定されず、本願の実施例は端末のタイプを具体的に限定しない。選択的に、端末上にアプリケーションクライアント端末がインストールされ、該アプリケーションクライアント端末はデータサービスを提供することを可能にするいずれかのクライアント端末であってもよく、例えば、該アプリケーションクライアント端末は支払いアプリケーションクライアント端末、フードデリバリーアプリケーションクライアント端末、タクシーアプリケーションクライアント端末又はソーシャルアプリケーションクライアント端末のうちの少なくとも1つであってもよく、本願の実施例はアプリケーションクライアント端末のタイプを具体的に限定しない。

10

20

【0040】

少なくとも2つのノード機器とは、分散型データベースシステムにおいて分散化されたノード機器としてトランザクションを協調処理することを可能にするノード機器を指し、個々のノード機器はいずれも分散化されたアルゴリズムによって分散型トランザクションを協調処理することに用いることを可能にする。

【0041】

トランザクション割り当て機器は、ターゲットトランザクションの割り当て要求を受信した後、ターゲットトランザクションに適切なノード機器を協調ノード機器として割り当てて、トランザクション処理の効率を確保する必要がある。ターゲットトランザクションに適切なノード機器を協調ノード機器として割り当てる過程において、トランザクション割り当て機器は先ず少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定する。1つのノード機器に対応するトランザクション割り当て指標は該1つのノード機器に新たなトランザクションを割り当てるマッチング度を指示することに用いられる。1つのノード機器に新たなトランザクションを割り当てるマッチング度が高いほど、該1つのノード機器に新たなトランザクションを割り当てるのにより適していることが説明される。

30

【0042】

トランザクション割り当て指標は、トランザクションの観点から決定された、あるノード機器に新たなトランザクションを割り当てるのに適しているか否かを評価することに用いられる指標である。1つの可能な実現形態において、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定する過程は以下のステップ2011とステップ2012とを含む。

40

【0043】

ステップ2011：トランザクション割り当てモードを決定し、トランザクション割り当てモードはトランザクションのビジー度に基づく割り当て、機器のビジー度に基づく割り当て、及びハイブリッドのビジー度に基づく割り当てのうちのいずれか1つを含む。

【0044】

トランザクション割り当てモードはノード機器に対応するトランザクション割り当て指標を決定する決定方式を指示することに用いられる。いくつかの実施例において、トランザクション割り当てモードは開発者によって設定されてトランザクション割り当て機器に

50

アップロードされる。説明する必要があるように、異なる時期に採用されたトランザクション割り当てモードは異なる可能性がある。ここでは、ステップ 2011 において決定されるのは、ターゲットトランザクションの割り当て要求を受信したときに採用すべきトランザクション割り当てモードである。

【0045】

トランザクション割り当てモードはトランザクションのビジー度に基づく割り当て、機器のビジー度に基づく割り当て、及びハイブリッドのビジー度に基づく割り当てのうちのいずれか1つを含む。ここで、トランザクションのビジー度に基づく割り当てモードとは、ノード機器のトランザクション処理数を考慮する観点からトランザクション割り当て指標を決定することを指し、ノード機器のトランザクション処理数はノード機器のトランザクションのビジー度を反映することを可能にする。機器のビジー度に基づく割り当てモードとは、ノード機器の機器リソース使用率を考慮する観点からトランザクション割り当て指標を決定することを指し、ノード機器の機器リソース使用率はノード機器の機器のビジー度を反映することを可能にする。ハイブリッドのビジー度に基づく割り当てモードとは、ノード機器のトランザクション処理数、及びノード機器の機器リソース使用率を総合的に考慮する観点からトランザクション割り当て指標を決定することを指し、ノード機器のトランザクション処理数、及びノード機器の機器リソース使用率はノード機器のハイブリッドのビジー度を反映することを可能にする。

10

【0046】

ステップ 2012：トランザクション割り当てモードによって指示された決定方式に基づき、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定する。

20

【0047】

異なるトランザクション割り当てモードによって指示された決定方式は異なり、トランザクション割り当てモードを決定した後、トランザクション割り当てモードによって指示された決定方式に基づき、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定する。次に、異なるトランザクション割り当てモード下で、少なくとも2つのノード機器のうちの第1ノード機器に対応するトランザクション割り当て指標を決定する方式をそれぞれ紹介する。ここで、第1ノード機器は少なくとも2つのノード機器のうちのいずれか1つのノード機器である。

30

【0048】

いくつかの実施例において、トランザクション割り当てモードはトランザクションのビジー度に基づく割り当てである。このような場合に、トランザクション割り当てモードによって指示された決定方式に基づき、第1ノード機器に対応するトランザクション割り当て指標を決定する方式としては、該第1ノード機器のトランザクション処理数に基づいて、該第1ノード機器に対応するトランザクション割り当て指標を決定する。

【0049】

第1ノード機器のトランザクション処理数とは、第1ノード機器が単位時間あたりに処理する必要があるトランザクションの数を指す。説明する必要があるように、ここでの処理する必要があるトランザクションとは、既に該第1ノード機器に割り当てられて処理されたトランザクションを指す。第1ノード機器が単位時間あたりに処理する必要があるトランザクションの数が多いほど、該第1ノード機器に新たなトランザクションを割り当てるのに適しないことが説明される。例示的な実施例において、第1ノード機器のトランザクション処理数は該第1ノード機器によってトランザクション割り当て機器にフィードバックされてもよく、トランザクション割り当て機器によってトランザクション割り当て状況に基づき自ら決定されてもよいが、本願の実施例はこれを限定しない。

40

【0050】

本願の実施例はトランザクション割り当て指標の表現形式を限定せず、例示的には、トランザクション割り当て指標の表現形式はビジーレベル、又は数値である。

【0051】

50

例示的には、トランザクション割り当て指標の表現形式がビジーレベルである場合、該第1ノード機器のトランザクション処理数に基づいて該第1ノード機器に対応するトランザクション割り当て指標を決定する方式としては、異なるビジーレベルに異なるトランザクション処理数の範囲を設定し、該第1ノード機器のトランザクション処理数があるトランザクション処理範囲に対応するビジーレベルを該第1ノード機器に対応するビジーレベルとして用いる。例示的には、ビジーレベルは「ビジー」、「部分ビジー」及び「アイドル」を含み、「ビジー」に対応するトランザクション処理数の範囲は $[10, +)$ であり、「部分ビジー」に対応するトランザクション処理数の範囲は $[3, 10)$ であり、「アイドル」に対応するトランザクション処理数の範囲は $[0, 3)$ である。第1ノード機器のトランザクション処理数が2であれば、「アイドル」を該第1ノード機器に対応するトランザクション割り当て指標として用いる。第1ノード機器に対応するトランザクション割り当て指標が「アイドル」に近くなるほど、該第1ノード機器に新たなトランザクションを割り当てるマッチング度が高くなることが説明される。

10

【0052】

例示的には、トランザクション割り当て指標の表現形式が数値である場合、該第1ノード機器のトランザクション処理数に基づいて該第1ノード機器に対応するトランザクション割り当て指標を決定する方式としては、該第1ノード機器のトランザクション処理数を数値化処理し、数値化処理した後に得られた数値を該第1ノード機器に対応するトランザクション割り当て指標として用いる。トランザクション処理数を数値化処理する方式は経験に基づき設定されるか、又は応用シナリオに基づき柔軟に調整されるが、本願の実施例はこれを限定しない。例示的には、トランザクション処理数を数値化処理する方式としては、トランザクション処理数と参照重みとの積を計算する。このような方式では、トランザクション処理数が多いほど、数値化処理した後に得られた数値が大きくなる。第1ノード機器に対応するトランザクション割り当て指標が小さいほど、該第1ノード機器に新たなトランザクションを割り当てるマッチング度が高くなることが説明される。

20

【0053】

いくつかの実施例において、トランザクション割り当てモードは機器のビジー度に基づく割り当てである。このような場合に、トランザクション割り当てモードによって指示された決定方式に基づき、第1ノード機器に対応するトランザクション割り当て指標を決定する方式としては、該第1ノード機器の機器リソース使用率に基づいて該第1ノード機器に対応するトランザクション割り当て指標を決定する。

30

【0054】

第1ノード機器の機器リソース使用率とは、第1ノード機器が既に使用した機器リソースが総機器リソースを占める割合を指し、例示的には、機器リソースとはCPU (Central Processing Unit、中央プロセッサ) リソースを指す。第1ノード機器の機器リソース使用率が高いほど、該第1ノード機器に新たなトランザクションを割り当てるのに適しないことが説明される。説明する必要があるように、第1ノード機器の機器リソース使用率は該第1ノード機器によってリアルタイムに監視されてトランザクション割り当て機器にフィードバックされてもよく、トランザクション割り当て機器によって自ら監視されて得られてもよいが、本願の実施例はこれを限定しない。

40

【0055】

該第1ノード機器の機器リソース使用率に基づいて該第1ノード機器に対応するトランザクション割り当て指標を決定する方式は、該第1ノード機器のトランザクション処理数に基づいて該第1ノード機器に対応するトランザクション割り当て指標を決定する方式を参照すればよく、ここでは詳細な説明を省略する。

【0056】

いくつかの実施例において、トランザクション割り当てモードはハイブリッドのビジー度に基づく割り当てである。このような場合に、トランザクション割り当てモードによって指示された決定方式に基づき、第1ノード機器に対応するトランザクション割り当て指標を決定する方式としては、第1ノード機器のトランザクション処理数、第1ノード機器

50

の機器リソース使用率、トランザクション処理数の重み、機器リソース使用率の重み、及び重み調節パラメータに基づいて、第1ノード機器に対応するトランザクション割り当て指標を決定する。

【0057】

例示的な実施例において、トランザクション処理数の重み、及び機器リソース使用率の重みはトランザクション処理数、及び機器リソース使用率の2つのパラメータのパーセンテージを調節することに用いられ、実測して得ることができ、例示的には、トランザクション処理数の重み、及び機器リソース使用率の重みのデフォルト値はいずれも1である。重み調節パラメータとは、機器リソース使用率とトランザクション処理数の相対割合係数を指し、機器リソース使用率とトランザクション処理数の重み割り当てを調節するために用いられ、実測して得ることができ、例示的には、重み調節パラメータのデフォルト値は0.33である。

10

【0058】

例示的には、 p_1 を利用して機器リソース使用率の重みを表し、 p_2 を利用してトランザクション処理数の重みを表し、 w を利用して重み調節パラメータを表すとすると、第1ノード機器に対応するトランザクション割り当て指標 Q は $Q = p_1 \times \text{機器リソース使用率} + p_2 \times w \times \text{トランザクション処理数}$ として表されてもよい。

【0059】

いくつかの実施例において、ハイブリッドのビジー度に基づく割り当てモードにおいて、トランザクション処理数、及び機器リソース使用率を総合的に考慮することに加えて、さらに他の態様の要素、例えば、処理する必要があるトランザクションのうちの長いトランザクションの数などを考慮することができる。このような場合に、第1ノード機器に対応するトランザクション割り当て指標 Q は $Q = p_1 \times \text{機器リソース使用率} + p_2 \times w \times \text{トランザクション処理数} + p_3 \times \text{他の要素}$ として表されてもよい。ここで、 p_3 は他の要素の重みを表し、 p_3 は他の要素のタイプに基づき実測して得ることができ、例示的には、 p_3 のデフォルト値は1である。第1ノード機器に対応するトランザクション割り当て指標 Q が小さいほど、該第1ノード機器に新たなトランザクションを割り当てるマッチング度が高くなることが説明される。

20

【0060】

説明する必要があるように、以上は第1ノード機器の観点からのみ第1ノード機器に対応するトランザクション割り当て指標を決定する過程を紹介し、上記方式に基づき分散型データベースシステムにおける少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定し、更にステップ202を実行することを可能にする。

30

【0061】

ステップ202において、トランザクション割り当て機器は、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、少なくとも2つのノード機器のうちターゲットトランザクションの協調ノード機器を決定し、協調ノード機器によってターゲットトランザクションを協調処理する。

【0062】

ターゲットトランザクションの協調ノード機器とは、少なくとも2つのノード機器のうち新たなトランザクションを割り当てるのに適しているノード機器を指す。ターゲットトランザクションの協調ノード機器はターゲットトランザクションを協調処理することに用いられ、つまり、ターゲットトランザクションの協調ノード機器とはターゲットトランザクションのコーディネーターを指す。例示的には、ターゲットトランザクションを協調処理する過程とは、分散型データベースシステムにおいてターゲットトランザクションを発生し、次にターゲットトランザクションのデータノード機器を組織して該ターゲットトランザクションを共同で処理する過程を指す。ターゲットトランザクションのデータノード機器とは、少なくとも2つのノード機器のうちターゲットトランザクションの処理に参与することに用いられるノード機器を指し、つまり、ターゲットトランザクションのデータノード機器とはターゲットトランザクションの関与者を指す。

40

50

【 0 0 6 3 】

説明する必要があるように、本願の実施例において言及された協調ノード機器及びデータノード機器はいずれもターゲットトランザクションに対するものであり、異なるトランザクションにとって、協調ノード機器又はデータノード機器は固定されるものではなく、すなわち、同一のノード機器はいくつかのトランザクションにとって協調ノード機器に属するが、別のいくつかのトランザクションにとってはデータノード機器に属する可能性がある。

【 0 0 6 4 】

少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、少なくとも2つのノード機器のうちターゲットトランザクションの協調ノード機器を決定する方式は、トランザクション割り当て指標の表現形式の差異に基づき異なっており、本願の実施例はこれを限定せず、協調ノード機器が現在新たなトランザクションを割り当てるのに適しているノード機器であることを確保することを可能にすればよい。

10

【 0 0 6 5 】

いくつかの実施例において、トランザクション割り当て指標の表現形式はビジーレベルであり、ビジーレベルはそれぞれ「ビジー」、「部分ビジー」及び「アイドル」である。このような場合に、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、少なくとも2つのノード機器のうちターゲットトランザクションの協調ノード機器を決定する方式としては、少なくとも2つのノード機器のうち対応するトランザクション割り当て指標が「アイドル」のノード機器を候補ノード機器として用い、候補ノード機器において1つのノード機器をターゲットトランザクションの協調ノード機器として任意に選択する。

20

【 0 0 6 6 】

例示的には、対応するトランザクション割り当て指標が「アイドル」のノード機器が存在しなければ、少なくとも2つのノード機器のうち対応するトランザクション割り当て指標が「部分ビジー」のノード機器を候補ノード機器として用いて、更に候補ノード機器において1つのノード機器をターゲットトランザクションの協調ノード機器として任意に選択する。例示的には、少なくとも2つのノード機器に対応するトランザクション割り当て指標がいずれも「ビジー」であれば、ターゲットトランザクションの協調ノード機器を決定することを一時的に停止し、参照期間待機した後少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を再決定し、更にターゲットトランザクションの協調ノード機器を再決定する。参照期間は経験に基づき設定され、例えば、参照期間は1つのトランザクションを完了する実測された平均期間である。

30

【 0 0 6 7 】

いくつかの実施例において、トランザクション割り当て指標の表現形式は数値であり、且つ1つのノード機器に対応するトランザクション割り当て指標が小さいほど、該1つのノード機器に新たなトランザクションを割り当てるマッチング度が高くなることが説明される。このような場合に、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、少なくとも2つのノード機器のうちターゲットトランザクションの協調ノード機器を決定する方式としては、少なくとも2つのノード機器のうち前 s (s は1以上の整数である)個の小さなトランザクション割り当て指標に対応するノード機器を候補ノード機器として用いて、候補ノード機器において1つのノード機器をターゲットトランザクションの協調ノード機器として任意に選択する。 s の値は経験に基づき設定されるか、又は少なくとも2つのノード機器の総数に基づき柔軟に調整され、本願の実施例はこれを限定せず、例えば、 s の値は1であり、又は s の値は3などである。

40

【 0 0 6 8 】

説明する必要があるように、以上は少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、少なくとも2つのノード機器のうちターゲットトランザクションの協調ノード機器を決定する方式の例示的な記述に過ぎず、本願の実施例はこれに限定されない。例示的には、トランザクション割り当て指標の表現形式が数

50

値であり、且つ1つのノード機器に対応するトランザクション割り当て指標が大きいほど、該1つのノード機器に新たなトランザクションを割り当てるマッチング度が高くなる。ことが説明される場合については、少なくとも2つのノード機器のうち前 t (t は1以上の整数である)個の大きなトランザクション割り当て指標に対応するノード機器を候補ノード機器として用い、候補ノード機器において1つのノード機器をターゲットトランザクションの協調ノード機器として任意に選択する。

【0069】

トランザクション割り当て指標に基づいて決定されたターゲットトランザクションの協調ノード機器は少なくとも2つのノード機器のうち新たなトランザクションを割り当てるのに適しているノード機器であり、更にターゲットトランザクションを該協調ノード機器に割り当てて、該協調ノード機器によって該ターゲットトランザクションを協調処理し、ターゲットトランザクションの処理効率を確保することに有利である。

10

【0070】

関連技術において、個々のノード機器は一定数のエリア (r e g i o n) にサービスし、個々のノード機器はいずれも該ノード機器によってサービスされるエリアにおけるデータ項目の分布情報を維持し、データ項目の分布情報はデータ項目の記憶位置を指示することに用いられる。この他、トランザクション割り当て機器においてエリアのメタ情報が維持される。このようなアーキテクチャでは、関連技術において、トランザクション割り当て機器は維持されたエリアのメタ情報に基づき、ターゲットトランザクションに関するデータ項目が所在するエリアにサービスすることに用いられるノード機器を決定し、更に該ノード機器によってターゲットトランザクションを独立して処理する。このような方式では、トランザクションの処理効率を大幅に制限し、真の分散型トランザクションをサポートすることができず、良好なトランザクション属性特徴付きのグローバルで一致性のあるマルチ読み取り及び一致性のあるマルチ書き込みの能力を備えない。

20

【0071】

本願の実施例においては、ノード機器は特定の固定されたエリアにサービスしなくなり、ノード機器はデータ項目の分布情報を維持しなくなり、トランザクション割り当て機器もエリアのメタ情報を維持しなくなる。例示的には、エリアのメタ情報を分散型データベースシステムにおける共有された記憶システム全体において分布する。このような改良に基づいて、トランザクション割り当て機器はトランザクション割り当て指標に基づいてターゲットトランザクションに適切なノード機器を協調ノード機器として割り当てることを可能にし、トランザクションに関するデータ項目を考慮する必要がなく、データ項目の分布状況を考慮する必要もない。ノード機器はトランザクションにおけるSQL (S t r u c t u r e d Q u e r y L a n g u a g e、構造化照会言語) ステートメントのニーズに基づき、共有された記憶システムからデータを自動的に呼び出すことを可能にする。このような方式では、個々のノード機器はいずれも分散化されたノード機器として分散型トランザクションを協調処理することを可能にすることにより、分散型データベースシステムが分散化された分散型トランザクションの処理能力を有する。

30

【0072】

1つの可能な実現形態において、ターゲットトランザクションの協調ノード機器を決定するステップの後に、協調ノード機器の機器識別情報を該割り当て要求を発した端末に送信するステップをさらに含み、端末は協調ノード機器の機器識別情報に基づき、ターゲットトランザクションのトランザクション情報を協調ノード機器に送信し、協調ノード機器によってトランザクション情報に基づいてターゲットトランザクションを協調処理することに用いられる。

40

【0073】

協調ノード機器の機器識別情報は、該協調ノード機器を一意的に識別することに用いられ、協調ノード機器の機器識別情報を割り当て要求を発した端末に送信することによって、端末にターゲットトランザクションを協調処理することに用いられる協調ノード機器を知らせることを可能にする。端末は機器識別情報に基づきターゲットトランザクションを

50

協調処理することに用いられる協調ノード機器を知った後、ターゲットトランザクションのトランザクション情報を協調ノード機器に送信する。ターゲットトランザクションのトランザクション情報はターゲットトランザクションの関連処理操作を指示することに用いられ、例示的には、ターゲットトランザクションのトランザクション情報とはSQLステートメントを指す。

【0074】

1つの可能な実現形態において、端末はターゲットトランザクションのトランザクション情報を協調ノード機器に直接送信し、又は、端末はターゲットトランザクションのトランザクション情報及び協調ノード機器の機器識別情報をゲートウェイサーバに送信し、ゲートウェイサーバによってターゲットトランザクションのトランザクション情報を協調ノード機器に転送する。

10

【0075】

協調ノード機器はターゲットトランザクションのトランザクション情報を受信した後、トランザクション情報に基づいてターゲットトランザクションを協調処理する。協調ノード機器はトランザクション情報、たとえば、SQLステートメントを解析し、トランザクション実行計画を生成し、次に関連するデータノード機器と通信することによってターゲットトランザクションの処理を完了することを可能にする。

【0076】

例示的な実施例において、ターゲットトランザクションの協調ノード機器はトランザクション割り当て指標に基づき決定され、異なるトランザクションは異なるノード機器を利用して協調処理を行うことを可能にするため、本願の実施例が提供する方法は分散化されたトランザクション処理過程を実現することを可能にする。分散化されたトランザクション処理過程において、複数の分散型トランザクションはそれぞれ複数のノード機器によって協調処理される。ターゲットトランザクションの協調ノード機器がターゲットトランザクションを協調処理する過程において、複数の分散型トランザクションが存在すると、協調ノード機器は他のノード機器と通信を確立して、他のノード機器が他の分散型トランザクションを協調処理する過程において生じたデータ情報を取得し、更に取得したデータ情報に基づきデータ異常又はシリアル化可能な認証を行って、ターゲットトランザクションがトランザクションの一致性を満たすか否かを判断し、トランザクション処理技術が正確であることを確保する。例示的な実施例において、協調ノード機器は他のノード機器から送られたデータ情報を一時データバッファ領域にバッファし、ターゲットトランザクションが終了するとクリーンアップする。

20

30

【0077】

1つの可能な実現形態において、ターゲットトランザクションの割り当て要求を発した端末にとって、該端末は分散型データベースシステムにアクセスした後、他のトランザクションが生じると、個々の他のトランザクションはいずれも該協調ノード機器によって協調処理され、又は、個々の他のトランザクションはいずれもトランザクション割り当て機器がノード機器のトランザクション割り当て指標に基づきリアルタイムに割り当てられた適した協調ノード機器によって協調処理されるが、本願の実施例はこれを限定しない。

【0078】

ステップ203において、協調ノード機器はターゲットトランザクションのトランザクション情報を取得する。

40

【0079】

ターゲットトランザクションのトランザクション情報はターゲットトランザクションの作成端末によって協調ノード機器に直接送信されてもよく、ゲートウェイサーバによって協調ノード機器に転送されてもよいが、本願の実施例はこれを限定しない。例示的には、ターゲットトランザクションのトランザクション情報とはターゲットトランザクションを実現することに用いられるSQLステートメントを指す。

【0080】

1つの可能な実現形態において、協調ノード機器はターゲットトランザクションのトラ

50

ンザクション情報を取得した後、ターゲットトランザクションを初期化する。ターゲットトランザクションを初期化する段階はトランザクションを確立するスナップショット段階とみなされてもよい。この段階ではグローバルな一貫性スナップショットポイントを確立し、グローバルな読み取り一貫性を保証することを可能にする。

【0081】

1つの可能な実現形態において、ターゲットトランザクションを初期化する過程において、協調ノード機器は以下の2つの初期化操作のうち少なくとも1つを実行することができる。

【0082】

初期化操作1：協調ノード機器はターゲットトランザクションに1つのグローバルで一意的なトランザクション識別子TIDを割り当てる。

【0083】

該トランザクション識別子TIDは該ターゲットトランザクションを一意的に識別することに用いられる。

【0084】

初期化操作2：協調ノード機器は第1トランザクションステータスリストにおいてターゲットトランザクションの初期ステータス情報を記録する。

【0085】

本願の実施例においては、協調ノード機器によって維持されたトランザクションステータスリストを第1トランザクションステータスリストと呼び、該第1トランザクションステータスリストは分散化フレームワーク下でターゲットトランザクションのグローバルステータスを記録することに用いられるグローバルステータスリストである。

【0086】

例示的な実施例において、第1トランザクションステータスリストにおいて記録されたターゲットトランザクションのステータス情報は、ターゲットトランザクションのトランザクション識別子、ターゲットトランザクションのグローバルトランザクションステータス及びターゲットトランザクションの論理ライフサイクルを含むがこれらに限定されない。ここで、論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成される。論理ライフサイクルのタイムスタンプ下限はターゲットトランザクションの開始タイムスタンプ(Begin timestamp、Bts)と呼ばれ、論理ライフサイクルのタイムスタンプ上限はターゲットトランザクションの終了タイムスタンプ(End timestamp、Ets)と呼ばれ、つまり、論理ライフサイクルはタイムスタンプ下限Bts、及びタイムスタンプ上限Etsで構成される。

【0087】

ターゲットトランザクションの初期ステータス情報において、ターゲットトランザクションのトランザクション識別子TIDは初期化操作1において割り当てられ、ターゲットトランザクションのグローバルトランザクションステータスStatusはRunning(グローバルに動作している)であり、ターゲットトランザクションの論理ライフサイクルは第1論理ライフサイクルであり、該第1論理ライフサイクルのタイムスタンプ下限Btsはグローバルで一意的に遡増するタイムスタンプ値であり、第1論理ライフサイクルのタイムスタンプ上限Etsは+である。

【0088】

例示的な実施例において、第1論理ライフサイクルのタイムスタンプ下限Bts、及びタイムスタンプ上限Etsの取得方式としては、シリアル化可能レベルより上の分離レベルに対しては、協調ノード機器はグローバルクロックからタイムスタンプ値を取得し、シリアル化可能レベル、及びより弱い分離レベルに対しては、協調ノード機器はローカルのハイブリッド論理クロック(Hybrid Logical Clock、HLC)からタイムスタンプ値を取得する。もちろん、いくつかの実施例において、シリアル化可能レベル、及びより弱い分離レベルに対しては、協調ノード機器はグローバルクロックからタイムスタンプ値を取得することによって第1論理ライフサイクルのタイムスタンプ下限B

10

20

30

40

50

t s、及びタイムスタンプ上限 E t s を取得することも可能にする。例示的な実施例において、シリアル化可能レベル、及びより弱い分離レベルに対して、ローカルの H L C からタイムスタンプ値を取得する効率が比較的高い。

【 0 0 8 9 】

グローバルクロックとはグローバル論理クロックジェネレータによって生成されたクロックを指し、単調に逓増する特性を備え、形式的には w a l l t i m e (ウォールタイム) であってもよく、自然数 N などであってもよい。例示的には、グローバルクロックは分散型データベースシステムにおけるグローバルタイムスタンプ生成クラスタによって提供される。例示的には、s h a r e - d i s k アーキテクチャに基づく分散型データベースシステムにとって、グローバルクロックは分散型データベースシステムにおける記憶システムによって A P I の方式を介して提供される。例示的には、グローバル論理クロックジェネレータはトランザクションの開始タイムスタンプ B t s、トランザクションの終了タイムスタンプ E t s に値を付与することを可能にし、さらに W A L (W r i t e A h e a d L o g g i n g、ログ先行書き込み) のグローバル L S N (L o g S e q u e n c e N u m b e r、ログシーケンス番号) に値を付与することを可能にする。

10

【 0 0 9 0 】

いくつかの実施例において、グローバルクロックは1つの論理的概念であり、全システムに統一的な単調逓増値を提供し、物理的形態は1つのグローバルな物理クロックであってもよく、1つのグローバルな論理クロックであってもよい。グローバルクロックの実現形態は、複数種の方式を有してもよい。たとえば、グローバルクロックは G o o g l e (グーグル) の「T r u e t i m e (一種のクロックメカニズム)」メカニズムと類似する1つの分散化された分散型クロックであり、又は、グローバルクロックは複数の冗長ノード(たとえば、一致性プロトコル(P a x o s / R a f t など)で構成されたクラスタ)を用いるマスタースタンプシステムが統一的に提供した1つのクロックであり、さらには、グローバルクロックは正確な同期メカニズムとノードログアウトメカニズムとを併せ持つ一種のアルゴリズムメカニズムによって提供された1つのクロックである。

20

【 0 0 9 1 】

例示的な実施例において、1つのトランザクションの B t s、及び E t s の構成はいずれも 8 バイトからなる。8 バイトは 2 つの部分に分けられ、第 1 部分は物理タイムスタンプの値(すなわち U n i x (一種のオペレーティングシステム)タイムスタンプ、ミリ秒まで正確にする)であり、グローバル時間(g t s を用いて表す)を識別することに用いられ、第 2 部分はあるミリ秒内の単調逓増カウントであり、グローバル時間上の相対時間(すなわち局所時間、 l t s を用いて表す)を識別することに用いられる。例示的には、8 バイトにおける上位 4 4 ビットは第 1 部分であり、このように合計で 2^{44} 個の符号なし整数を表すことができる。従って理論的には合計で約 5 5 7 . 8 ([数 1]) 年の物理タイムスタンプを表すことができ、8 バイトにおける下位 2 0 ビットは第 2 部分であり、このように、ミリ秒毎に 2^{20} 個(約 1 0 0 万個)のカウントがある。例示的な実施例において、2 つの部分のビット数を調整することにより、グローバル時間 g t s、及び局所時間 l t s で表される範囲を変化させることもできる。

30

【 0 0 9 2 】

【 数 1 】

$$\frac{2^{44}}{1000 \times 60 \times 60 \times 24 \times 365} = 557.8$$

40

【 0 0 9 3 】

説明する必要があるように、実際の必要に基づき、1つのトランザクションの B t s、及び E t s の構成は 8 バイトよりも多いバイトからなるか、又は 8 バイトよりも少ないバイトからなってもよい。例示的には、B t s、及び E t s の構成を 1 0 バイトからなるよ

50

うに調整することにより、局所時間 lts を増大させ、より大きな並行トランザクションの数に対処する。

【0094】

例示的には、グローバル時間 gts 、及び局所時間 lts の2つの部分で構成された2つのタイムスタンプ $T_i.bts$ 、及び $T_j.bts$ にとって、 $T_i.bts.gts < T_j.bts.gts$ 、又は $T_i.bts.gts = T_j.bts.gts$ 且つ $T_i.bts.lts < T_j.bts.lts$ となると、 $T_i.bts < T_j.bts$ となるとみなされる。

【0095】

ステップ204において、協調ノード機器はターゲットトランザクションのトランザクション情報に基づいて、データノード機器にデータ読み取り要求を送信し、データノード機器は少なくとも2つのノード機器のうちターゲットトランザクションの処理に関与することに用いられるノード機器である。

10

【0096】

協調ノード機器はターゲットトランザクションを初期化した後、ターゲットトランザクションの実行段階を開始し、トランザクションの実行段階はトランザクションセマンティクス実現操作段階とみなされてもよい。

【0097】

データノード機器は少なくとも2つのノード機器のうちターゲットトランザクションの処理に関与することに用いられるノード機器であり、データノード機器はターゲットトランザクションに関するデータ項目を取得することを可能にし、すなわち本願の実施例におけるデータノード機器はターゲットトランザクションに関連するデータノード機器である。ターゲットトランザクションのトランザクション情報において読み取る必要があるデータの関連情報を携え、協調ノード機器はターゲットトランザクションのトランザクション情報に基づき、データ読み取り要求を生成し、次にデータノード機器にデータ読み取り要求を送信することを可能にする。

20

【0098】

1つの可能な実現形態において、データ読み取り要求は `ReadRequestMessage` (読み取り要求メッセージ) として表され、`rrqm` と略称する。

【0099】

1つの可能な実現形態において、データ読み取り要求はターゲットトランザクションの第1論理ライフサイクル、ターゲットトランザクションのトランザクション識別子及び読み取り計画を携えている。第1論理ライフサイクルはタイムスタンプ下限 Bts 、及びタイムスタンプ上限 Ets を利用して表される。読み取り計画とはターゲットトランザクションに対応するデータ読み取り計画を指し、読み取る必要があるデータ項目を指示することに用いられる。例示的な実施例において、トランザクション識別子、タイムスタンプ下限 Bts 、タイムスタンプ上限 Ets 、及び読み取り計画はそれぞれ `rrqm` の4つのフィールドにおいて記録される。

30

【0100】

説明する必要があるように、データノード機器の数は1つ又は複数であってもよく、本願の実施例においてデータノード機器の数を具体的に限定しない。データノード機器の数が複数である場合については、異なるデータノード機器に送信されたデータ読み取り要求において携える読み取り計画が異なり、異なるデータノード機器において異なるデータ項目を読み取る必要があるように指示する。

40

【0101】

ステップ205において、データノード機器は協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得し、データ読み取り結果を協調ノード機器に返信する。

【0102】

データノード機器はデータ読み取り要求を取得した後、データ読み取り要求に基づいて

50

、データ読み取り結果を取得する。1つの可能な実現形態において、データ読み取り要求がターゲットトランザクションの第1論理ライフサイクルを携えている場合については、データノード機器がデータ読み取り要求に基づいて、データ読み取り結果を取得する過程は以下のステップ2051～ステップ2053を含む。

【0103】

ステップ2051：第1論理ライフサイクルに基づいて、データ読み取り要求によって指示された読み取り対象のデータ項目の可視バージョンデータを決定する。

【0104】

データノード機器はデータ読み取り要求が携えている読み取り計画に基づいて、ターゲットトランザクションが読み取る必要があるデータ項目を決定し、ターゲットトランザクションが読み取る必要があるデータ項目を読み取り対象のデータ項目として用いることを可能にする。読み取り対象のデータ項目の可視バージョンデータとは読み取り対象のデータ項目に対応する各バージョンデータのうちターゲットトランザクションに対して可視なあるバージョンデータを指す。例示的には、データノード機器においてデータバッファ領域が設定され、データバッファ領域において読み取り対象のデータ項目の各バージョンデータが存在すると、データノード機器はデータバッファ領域から読み取り対象のデータ項目の各バージョンデータを直接取得し、データバッファ領域において読み取り対象のデータ項目の各バージョンデータが存在しなければ、データノード機器は共有された記憶システムから読み取り対象のデータ項目の各バージョンデータを取得する。

【0105】

例示的な実施例において、データノード機器はデータ読み取り要求を受信した後、先ずローカルトランザクションステータスリスト(L o c a l T S)においてターゲットトランザクションのステータス情報が含まれるか否かをチェックする。ローカルトランザクションステータスリストはデータノード機器によって維持されたトランザクションステータスリストであり、ローカルトランザクションステータスリストにおいて該データノード機器が関与するコミットされていない各トランザクションのステータス情報が記録される。例示的な実施例において、データノード機器はデータ読み取り要求を受信した後、データ読み取り要求が携えているターゲットトランザクションのトランザクション識別子に基づき、ローカルトランザクションステータスリストにおいてターゲットトランザクションのステータス情報が含まれるか否かをチェックし、チェック結果は以下の2つの種類を含む。

【0106】

チェック結果1：ローカルトランザクションステータスリストにおいてターゲットトランザクションのステータス情報が含まれていない。

【0107】

このような場合に、ローカルトランザクションステータスリストにおいて該ターゲットトランザクションのステータス情報を初期化することができ、すなわちL o c a l T Sにおいてターゲットトランザクションに関連する1つのレコードを挿入し、該レコードにおける値はそれぞれ、データ読み取り要求が携えているターゲットトランザクションのトランザクション識別子r r q m . T I D、データ読み取り要求が携えているターゲットトランザクションの第1論理ライフサイクルのタイムスタンプ下限r r q m . B t s、データ読み取り要求が携えているターゲットトランザクションの第1論理ライフサイクルのタイムスタンプ上限r r q m . E t s、及びデータ読み取り要求によって指示されたターゲットトランザクションの現在のトランザクションステータスr r q m . R u n n i n gである。

【0108】

このような場合に、第1論理ライフサイクルに基づいて、データ読み取り要求によって指示された読み取り対象のデータ項目の可視バージョンデータを決定する方式としては、読み取り対象のデータ項目の第1論理ライフサイクルに対する可視バージョンデータを決定する。

10

20

30

40

50

【0109】

チェック結果2：ローカルランザクションステータスリストにおいてターゲットランザクションのステータス情報が含まれる。

【0110】

このような場合に、データ読み取り要求を受信する前に、ターゲットランザクションが既に該データノード機器にアクセスしたことが説明される。このとき、該データノード機器上のターゲットランザクションのステータス情報を更新することができ、更新方法としては、ターゲットランザクションの論理ライフサイクルのタイムスタンプ下限 $T.Bts$ を、照会されたタイムスタンプ下限 $T.Bts$ とデータ読み取り要求において携えているタイムスタンプ下限 $rrqm.Bts$ (すなわち第1論理ライフサイクルのタイムスタンプ下限) のうちの最大値に更新し、すなわち、 $T.Bts = \max(T.Bts, rrqm.Bts)$ とする。この他、さらにターゲットランザクションの論理ライフサイクルのタイムスタンプ上限 $T.Ets$ を、照会されたタイムスタンプ上限 $T.Ets$ とデータ読み取り要求において携えているタイムスタンプ上限 $rrqm.Ets$ (すなわち第1論理ライフサイクルのタイムスタンプ上限) における最小値に更新し、すなわち、 $T.Ets = \min(T.Ets, rrqm.Ets)$ とする。更新されたタイムスタンプ下限、及び更新されたタイムスタンプ上限で構成された論理ライフサイクルを更新された論理ライフサイクルとして用いる。

10

【0111】

このような場合に、第1論理ライフサイクルに基づいて、データ読み取り要求によって指示された読み取り対象のデータ項目の可視バージョンデータを決定する方式としては、第1論理ライフサイクルに基づいて、更新された論理ライフサイクルを決定し、読み取り対象のデータ項目の更新された論理ライフサイクルに対する可視バージョンデータを決定する。

20

【0112】

説明する必要があるように、読み取り対象のデータ項目の第1論理ライフサイクルに対する可視バージョンデータを決定する実現形態は、読み取り対象のデータ項目の更新された論理ライフサイクルに対する可視バージョンデータを決定する実現形態と類似しており、本願の実施例においては、読み取り対象のデータ項目の第1論理ライフサイクルに対する可視バージョンデータを決定することを例として説明する。

30

【0113】

1つの可能な実現形態において、読み取り対象のデータ項目の第1論理ライフサイクルに対する可視バージョンデータを決定する前に、先ず第1論理ライフサイクルに対して正当性検査を行って、第1論理ライフサイクルが有効であるか否かを判断する。例示的には、第1論理ライフサイクルに対して正当性検査を行う方式としては、第1論理ライフサイクルのタイムスタンプ下限が第1論理ライフサイクルのタイムスタンプ上限よりも小さいか否かを検査する。第1論理ライフサイクルのタイムスタンプ下限が第1論理ライフサイクルのタイムスタンプ上限以上であれば、第1論理ライフサイクルが無効であることが説明され、このとき、ローカルランザクションステータスリストにおけるターゲットランザクションのランザクションステータスを `Running` から `Aborted` (アボート) に更新する。この他、データノード機器は協調ノード機器に `Abort` (アボート) メッセージを携えているデータ読み取り結果を返信する。データ読み取り結果は `ReadReplyMessage` (読み取りフィードバックメッセージ) として表され、`rrpm` と略称する。データ読み取り結果が `Abort` メッセージが携えている場合については、`rrpm` メッセージにおける `IsAbort` フィールドは1に等しく、すなわち `rrpm.IsAbort = 1` となる。

40

【0114】

第1論理ライフサイクルのタイムスタンプ下限が第1論理ライフサイクルのタイムスタンプ上限よりも小さい場合は、第1論理ライフサイクルが有効であることが説明され、このとき、読み取り対象のデータ項目の第1論理ライフサイクルに対する可視バージョンデ

50

ータを決定する操作を実行する。

【0115】

1つの可能な実現形態において、読み取り対象のデータ項目の第1論理ライフサイクルに対する可視バージョンデータを決定する過程としては、読み取り対象のデータ項目の最新のバージョンデータの作成タイムスタンプが第1論理ライフサイクルのタイムスタンプ上限よりも小さいことに応答して、該最新のバージョンデータを可視バージョンデータとして用いて、読み取り対象のデータ項目の最新のバージョンデータの作成タイムスタンプが第1論理ライフサイクルのタイムスタンプ上限以上であることに応答して、最初の作成タイムスタンプが第1論理ライフサイクルのタイムスタンプ上限よりも小さいあるバージョンデータを決定するまで、読み取り対象のデータ項目の前のバージョンデータと第1論理ライフサイクルのタイムスタンプ上限とを継続的に比較し、該バージョンデータを可視バージョンデータとして用いる。

10

【0116】

つまり、データノード機器は読み取り対象のデータ項目 x のある論理ライフサイクルに対する可視バージョンデータを決定する過程において、先ず読み取り対象のデータ項目 x の最新のバージョンデータからチェックを開始し、もし該論理ライフサイクルのタイムスタンプ上限 $T.Ets$ が最新のバージョンデータの作成タイムスタンプ Wts よりも大きければ、該最新のバージョンデータは該論理ライフサイクルに対する可視バージョンデータである。そうでなければ、該最新のバージョンデータは該論理ライフサイクルに対する可視バージョンデータではなく、 $T.Ets > Wts$ を満たす最初のあるバージョンデータ $x.v$ を見つけるまで、前のバージョンデータを検索する必要があり、該バージョンデータ $x.v$ を該論理ライフサイクルに対する可視バージョンデータとして用いる。

20

【0117】

1つの可能な実現形態において、可視バージョンデータを決定した後、該可視バージョンデータ $x.v$ を該ターゲットトランザクションの読み取りセットにおいて記憶する。選択的に、ここでの読み取りセットはローカル読み取りセットであってもよく、グローバル読み取りセットであってもよい。本願の実施例においては、該読み取りセットがローカル読み取りセットであることを例として説明しており、グローバル読み取りセットの同期に起因する通信オーバーヘッドを回避することを可能にする。

【0118】

例示的な実施例において、1つのトランザクションの読み取りセットにおいて該トランザクションが読み取る必要があるデータ項目の可視バージョンデータが記録される。説明する必要があるように、1つの分散型読み取りトランザクションにとっては、該分散型読み取りトランザクションの読み取りセットはローカル読み取りセット及びグローバル読み取りセットに分割されてもよく、ローカル読み取りセットはデータノード機器上に存在し、グローバル読み取りセットは協調ノード機器上に存在する。もちろん、協調ノード機器はグローバル読み取りセットを各データノード機器上に定期的に同期できることにより、データノード機器上にトランザクションのグローバル読み取りセットを維持させることを可能にする。

30

【0119】

ステップ2052：可視バージョンデータの作成タイムスタンプ、及び第1論理ライフサイクルに基づいて、ターゲットトランザクションの第2論理ライフサイクルを決定する。

40

【0120】

可視バージョンデータを決定した後、データノード機器は可視バージョンデータの作成タイムスタンプ、及び第1論理ライフサイクルに基づいて、ターゲットトランザクションの第2論理ライフサイクルを決定する。

【0121】

いくつかの実施例において、可視バージョンデータとは第1論理ライフサイクルに対する可視バージョンデータを指す場合、該ステップ2052の実現形態としては、可視バージョン

50

ジョンデータの作成タイムスタンプ、及び第 1 論理ライフサイクルに直接基づいて、ターゲットトランザクションの第 2 論理ライフサイクルを決定する。可視バージョンデータとは第 1 論理ライフサイクルに基づき決定された更新された論理ライフサイクルに対する可視バージョンデータを指す場合、該ステップ 2052 の実現形態としては、可視バージョンデータの作成タイムスタンプ、及び第 1 論理ライフサイクルに基づき決定された更新された論理ライフサイクルに基づいて、ターゲットトランザクションの第 2 論理ライフサイクルを決定する。本願の実施例は、可視バージョンデータの作成タイムスタンプ、及び第 1 論理ライフサイクルに直接基づいて、ターゲットトランザクションの第 2 論理ライフサイクルを決定することを例として説明する。

【0122】

10

1 つの可能な実現形態において、可視バージョンデータの作成タイムスタンプ、及び第 1 論理ライフサイクルに直接基づいて、ターゲットトランザクションの第 2 論理ライフサイクルを決定する方式としては、第 1 論理ライフサイクルのタイムスタンプ下限を調整し、第 1 論理ライフサイクルのタイムスタンプ下限を可視バージョンデータ $x.v$ の作成タイムスタンプよりも大きくし、すなわち、 $T.Bts > x.v.Wts$ にして、書き読み異常を除去し、調整した後に得られた論理ライフサイクルを第 2 論理ライフサイクルとして用いる。

【0123】

別の可能な実現形態において、可視バージョンデータが読み取り対象のデータ項目の最新のバージョンデータである場合については、可視バージョンデータの作成タイムスタンプ、及び第 1 論理ライフサイクルに直接基づいて、ターゲットトランザクションの第 2 論理ライフサイクルを決定する方式としては、第 1 論理ライフサイクルのタイムスタンプ下限を調整し、第 1 論理ライフサイクルのタイムスタンプ下限を可視バージョンデータ $x.v$ の作成タイムスタンプよりも大きくし、すなわち、 $T.Bts > x.v.Wts$ にして、書き読み異常を除去し、可視バージョンデータに対応する書き込み対象のトランザクションがヌルではないことに応答して、第 1 論理ライフサイクルのタイムスタンプ上限を調整し、第 1 論理ライフサイクルのタイムスタンプ上限を可視バージョンデータに対応する書き込み対象のトランザクションの論理ライフサイクルのタイムスタンプ下限よりも小さくし、すなわち $T.Ets < T0.Bts$ ($T0$ は可視バージョンデータに対応する書き込み対象のトランザクションを表す) にして、読み書き競合を除去し、調整した後に得られた論理ライフサイクルを第 2 論理ライフサイクルとして用いる。

20

30

【0124】

可視バージョンデータに対応する書き込み対象のトランザクション WT は可視バージョンデータと対応するデータ項目を修正しており、且つ認証に合格したトランザクションである。例示的には、書き込み対象のトランザクションのトランザクション識別子を記録することによって書き込み対象のトランザクションを記録する。いくつかの実施例において、可視バージョンデータが最新のバージョンデータである場合については、ターゲットトランザクションのトランザクション識別子を可視バージョンデータのアクティブトランザクション集合において追加し、可視バージョンデータをターゲットトランザクションのローカル読み取りセットにおいて追加する。

40

【0125】

アクティブトランザクション集合 ($RTlist$) は該最新のバージョンデータにアクセスしたアクティブトランザクションを記録することに用いられ、読み取りトランザクションリストと呼ばれてもよく、該アクティブトランザクション集合は配列の形式であってもよく、リスト、キュー、スタックなどの形式であってもよく、本願の実施例はアクティブトランザクション集合の形式を具体的に限定せず、 $RTlist$ における個々の要素は上記最新のバージョンデータを読み取ったトランザクションのトランザクション識別子 (TID) であってもよい。

【0126】

ステップ 2053 : 第 2 論理ライフサイクル、及び可視バージョンデータを携えている

50

結果をデータ読み取り結果として用いる。

【0127】

第2論理ライフサイクル、及び可視バージョンデータを決定した後、データノード機器は第2論理ライフサイクル、及び可視バージョンデータを携えている結果をデータ読み取り結果として用い、次に第2論理ライフサイクル、及び可視バージョンデータを携えているデータ読み取り結果を協調ノード機器に返信し、協調ノード機器に第2論理ライフサイクル、及び可視バージョンデータを取得させる。例示的には、データ読み取り結果は `ReadReplyMessage` (読み取りフィードバックメッセージ) として表され、`rrpm` と略称する。例示的には、第2論理ライフサイクル、及び可視バージョンデータを携えている `rrpm` において `Bts`、`Ets`、及び `Value` フィールドが含まれ、ここで、`Bts` フィールド及び `Ets` フィールドはそれぞれ第2論理ライフサイクルのタイムスタンプ下限、及び第2論理ライフサイクルのタイムスタンプ上限を記録し、`Value` フィールドは可視バージョンデータの値を記録する。

10

【0128】

ステップ206において、協調ノード機器はデータノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信する。

【0129】

データノード機器はデータ読み取り結果を協調ノード機器に返信した後、協調ノード機器はデータ読み取り結果がトランザクション認証条件を満たすか否かを判断し、次にデータ読み取り結果がトランザクション認証条件を満たすと判定した場合、データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信し、データノード機器にターゲットトランザクションを認証させる。

20

【0130】

例示的には、協調ノード機器はデータ読み取り結果がトランザクション認証条件を満たすか否かを判断する過程において、先ずデータ読み取り結果が `Abort` (アボート) メッセージを携えているか否かを判断し、すなわち `rrpm` における `IsAbort` フィールドが1に等しいか否かをチェックする。データ読み取り結果が `Abort` メッセージを携えており、すなわち `rrpm.IsAbort = 1` とすると、データ読み取り結果がトランザクション認証条件を満たさないとみなされ、このとき、グローバルアボート段階に入る。

30

【0131】

データ読み取り結果が `Abort` メッセージを携えていなければ、第1トランザクションステータスリストにおけるターゲットトランザクションの論理ライフサイクルを更新し、更新方式としては、第1論理ライフサイクルのタイムスタンプ下限と第2論理ライフサイクルのタイムスタンプ下限における最大値をターゲットトランザクションの第3論理ライフサイクルのタイムスタンプ下限として用い、第1論理ライフサイクルのタイムスタンプ上限と第2論理ライフサイクルのタイムスタンプ上限における最小値をターゲットトランザクションの第3論理ライフサイクルのタイムスタンプ上限として用いる。すなわち、 $T.Bts = \max(T.Bts, rrpm.Bts)$ 、 $T.Ets = \min(T.Ets, rrpm.Ets)$ とし、ここで、括弧内の $T.Bts$ 、及び $T.Ets$ はそれぞれ更新前の論理ライフサイクル (すなわち第1論理ライフサイクル) のタイムスタンプ下限、及びタイムスタンプ上限であり、`rrpm.Bts`、及び `rrpm.Ets` はそれぞれデータ読み取り結果が携えている第2論理ライフサイクルのタイムスタンプ下限、及びタイムスタンプ上限である。

40

【0132】

第1トランザクションステータスリストにおけるターゲットトランザクションの論理ライフサイクルを更新した後、第1トランザクションステータスリストにおける $T.Bts$ が $T.Ets$ よりも小さいか否かをチェックし、すなわち第3論理ライフサイクルのタイムスタンプ下限が第3論理ライフサイクルのタイムスタンプ上限よりも小さいか否かをチ

50

エックして、第3論理ライフサイクルが有効であるか否かを判断する。第3論理ライフサイクルのタイムスタンプ下限が第3論理ライフサイクルのタイムスタンプ上限以上である場合、第3論理ライフサイクルが無効であり、このような場合に、データ読み取り結果がトランザクション認証条件を満たさないとみなされ、グローバルアポート段階に入り、第3論理ライフサイクルのタイムスタンプ下限が第3論理ライフサイクルのタイムスタンプ上限よりも小さい場合、第3論理ライフサイクルが有効であり、このような場合に、データ読み取り結果がトランザクション認証条件を満たすとみなされ、データノード機器に第3論理ライフサイクルを携えているトランザクション認証要求を送信する。

【0133】

例示的な実施例において、もし協調ノード機器がターゲットトランザクションのアポートを決めれば、第1トランザクションステータスリストにおけるターゲットトランザクションのグローバルトランザクションステータスを `G a b o r t i n g` (グローバルにアポートしている) に修正し、関連するサブノード (すなわちデータノード機器) に局所アポートを行うように通知する必要がある。

【0134】

例示的な実施例において、トランザクション認証要求を送信する前に、協調ノード機器は第1トランザクションステータスリストにおけるターゲットトランザクションのグローバルトランザクションステータスを `G v a l i d a t i n g` (グローバルに認証している) に修正する。例示的には、トランザクション認証要求は `V a l i d a t e R e q u e s t M e s s a g e` (認証要求メッセージ) として表され、`v r m` と略称する。例示的には、`v r m` において `B t s`、及び `E t s` フィールドが含まれる。ここで、`B t s` フィールド及び `E t s` フィールドはそれぞれターゲットトランザクションの第1トランザクションステータスリストにおける最新の論理ライフサイクルのタイムスタンプ下限、及びタイムスタンプ上限、すなわち第3論理ライフサイクルのタイムスタンプ下限、及びタイムスタンプ上限を記録する。

【0135】

いくつかの実施例において、データノード機器の数は複数であり、このような場合に、個々のデータノード機器はいずれも1つのデータ読み取り結果を返信し、データ読み取り結果がトランザクション認証条件を満たすことは、各データノード機器から返信された各データ読み取り結果がいずれもトランザクション認証条件を満たすことを指す。このような場合に、第3論理ライフサイクルは各データ読み取り結果を総合的に考慮することによって決定された論理ライフサイクルである。

【0136】

いくつかの実施例において、全部の所要のデータを読み取り且つ更新をローカルの内部メモリにおいて書き込んだ後、トランザクション認証条件を満たすとみなされる。つまり、協調ノード機器は第3論理ライフサイクルが有効であることに応答して、且つターゲットトランザクションのグローバル書き込みセットをローカルの内部メモリにおいて記憶し、データノード機器にトランザクション認証要求を送信する。ターゲットトランザクションのグローバル書き込みセットは端末によって生成されて協調ノード機器に伝送されるか、又は協調ノード機器によって自ら生成され、本願の実施例はこれを限定しない。

【0137】

1つのトランザクションの書き込みセットにおいて該トランザクションが更新する必要があるデータ項目が記憶され、読み取りセット構造と類似しており、同様に内部メモリチェーンテーブル構造を使用してトランザクションの書き込みセットを維持することができる。説明する必要があるように、1つの分散型書き込みトランザクションにとって、該分散型書き込みトランザクションの書き込みセットはローカル書き込みセット及びグローバル書き込みセットに分割されてもよく、ローカル書き込みセットはデータノード機器上に存在し、グローバル書き込みセットは協調ノード機器上に存在する。もちろん、協調ノード機器はグローバル書き込みセットを各データノード機器上に定期的に同期できることにより、データノード機器上にもトランザクションのグローバル書き込みセットを維持する

10

20

30

40

50

ことを可能にする。

【0138】

ターゲットトランザクションのグローバル書き込みセットを協調ノード機器のローカルの内部メモリに書き込んだ後、協調ノード機器はグローバル書き込みセットに基づいてデータノード機器のローカル書き込みセットを決定して、トランザクション認証要求をローカル書き込みセットとともにデータノード機器に送信する。データノード機器のローカル書き込みセットとはターゲットトランザクションのグローバル書き込みセットのうちデータノード機器によって書き込みを担当する必要がある書き込みセットを指す。

【0139】

ターゲットトランザクションの読み取り段階では、通信は主に協調ノード機器と関連するデータノード機器との間で発生し、データの読み取りに成功するたびに2回の通信が必要であり、協調ノード機器はデータ読み取り要求を関連するデータノード機器上に送信し、関連するデータノード機器はデータ読み取り結果を協調ノード機器に返信する。従って、データ読み取り段階では、 n (n が1よりも大きい整数である) がリモート読み取りの回数であると仮定すると、最大 $2n$ 回の通信を行う必要があり、最大通信量は $n \times$ (データ読み取り要求メッセージのサイズ+データ読み取り結果メッセージのサイズ) として表されてもよい。例示的な実施において、ターゲットトランザクションがある関連するデータノード機器の複数のデータ項目のデータを読み取る必要がある場合、複数のデータ項目のデータのデータ読み取り要求をパッケージ化して送信し、これらのデータを一括して読み取り、通信回数を節約し、データ読み取り効率を向上させる。

【0140】

ステップ207において、データノード機器は協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、ターゲットトランザクションの認証結果を取得し、ターゲットトランザクションの認証結果を協調ノード機器に返信する。

【0141】

データノード機器は協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットを受信した後、ターゲットトランザクションの正当性を認証して、ターゲットトランザクションの認証結果を取得する。この段階はトランザクションをコミットする前のトランザクション正当性認証段階である。

【0142】

データノード機器の認証過程はローカル認証過程であり、データノード機器がトランザクション認証要求、及びローカル書き込みセットに基づいて、ターゲットトランザクションの認証結果を取得する過程はデータノード機器がローカル認証操作を実行する過程である。1つの可能な実現形態において、トランザクション認証要求は第3論理ライフサイクルを携えており、第3論理ライフサイクルは協調ノード機器によって第1論理ライフサイクル、及び第2論理ライフサイクルに基づいて決定された有効論理ライフサイクルである。第3論理ライフサイクルは協調ノード機器がトランザクション認証要求を送信する前に維持したターゲットトランザクションの最新の論理ライフサイクルである。

【0143】

1つの可能な実現形態において、データノード機器はトランザクション認証要求、及びローカル書き込みセットに基づいて、ターゲットトランザクションの認証結果を取得する過程において、データノード機器は先ずローカルトランザクションステータスリストにおけるターゲットトランザクションTのステータス情報を更新し、更新方式としては、 $T.Bts = \max(T.Bts, vrm.Bts)$ 、 $T.Ets = \min(T.Ets, vrm.Ets)$ とし、ここで、括弧における $vrm.Bts$ 、及び $vrm.Ets$ はそれぞれトランザクション認証要求が携えている第3論理ライフサイクルのタイムスタンプ下限、及びタイムスタンプ上限である。本願の実施例において、トランザクション認証要求を受信する前に、データノード機器のローカルトランザクションステータスリストにおいて維持されたターゲットトランザクションの論理ライフサイクルは第2論理ライフサイク

10

20

30

40

50

ルであり、区別を容易にするために、トランザクション認証要求を受信した後且つターゲットトランザクションのステータス情報を更新した後に、ローカルトランザクションステータスリストに維持されたターゲットトランザクションの論理ライフサイクルは第4論理ライフサイクルと呼ばれる。

【0144】

つまり、データノード機器は第3論理ライフサイクルのタイムスタンプ下限と第2論理ライフサイクルのタイムスタンプ下限における最大値をターゲットトランザクションの第4論理ライフサイクルの下限として用い、第3論理ライフサイクルのタイムスタンプ上限と第2論理ライフサイクルのタイムスタンプ上限における最小値をターゲットトランザクションの第4論理ライフサイクルのタイムスタンプ上限として用いる。これにより、第4論理ライフサイクルを得る。説明する必要があるように、ここで更新されるのはデータノード機器のローカルトランザクションステータスリストにおいて維持されたターゲットトランザクションの論理ライフサイクルであり、このような更新はトランザクション並行アクセス制御に用いられ、すなわちトランザクションの一致性を確保することに用いることを可能にする。

10

【0145】

例示的な実施例において、シリアル化可能分離レベルに対しては、第4論理ライフサイクルを決定した後、第4論理ライフサイクルのタイムスタンプ下限が第4論理ライフサイクルのタイムスタンプ上限よりも小さいか否かをチェックすることによって、第4論理ライフサイクルが有効であるか否かを認証する。

20

【0146】

第4論理ライフサイクルのタイムスタンプ下限が第4論理ライフサイクルのタイムスタンプ上限以上であることに応答して、第4論理ライフサイクルが無効であることが説明され、このとき、ターゲットトランザクションのローカル認証に合格せず、データノード機器は協調ノード機器にAbortメッセージを携えている認証結果を返信する。該Abortメッセージはグローバルアポートを始めることに用いられる。協調ノード機器にターゲットトランザクションの認証結果を返信する過程は協調ノード機器にローカル認証フィードバックメッセージlvmmを送信する過程とみなしてもよく、ターゲットトランザクションの認証結果がAbortメッセージを携えている認証結果である場合については、ローカル認証フィードバックメッセージlvmmにおけるIsAbortフィールドは1に等しく、すなわちlvmm.IsAbort=1とする。

30

【0147】

第4論理ライフサイクルのタイムスタンプ下限が第4論理ライフサイクルのタイムスタンプ上限よりも小さいことに応答して、第4論理ライフサイクルが有効であることが説明され、このような場合に、ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定する。第5論理ライフサイクルとはローカル書き込みセットにおける各書き込み対象のデータ項目に対して読み書き競合認証を行う過程において更新して得られた論理ライフサイクルを指す。

【0148】

1つの可能な実現形態において、1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は該1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び該1つの書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプのうち少なくとも1つを含む。ここで、1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ(Rtsと記される)は該1つの書き込み対象のデータ項目を読み取った各読み取りトランザクションの論理コミットタイムスタンプにおける最大値を指示することに用いられ、1つの書き込み対象のデータ項目のターゲット読み取りトランザクションは該1つの書き込み対象のデータ項目に対応するローカル認証に合格するか、又はコミット段階にある読み取りトランザクションであり、ターゲット読み取りトランザクションの終了タイムスタンプはターゲット読み

40

50

取りトランザクションの論理ライフサイクルのタイムスタンプ上限である。

【0149】

例示的な実施例において、1つの書き込み対象のデータ項目のターゲット読み取りトランザクションは該1つの書き込み対象のデータ項目に対応するアクティブトランザクション集合におけるローカル認証に合格するか、又はコミット段階にある読み取りトランザクションである。該1つの書き込み対象のデータ項目に対応するアクティブトランザクション集合における各読み取りトランザクションのトランザクションステータスを検出することによって、該1つの書き込み対象のデータ項目のターゲット読み取りトランザクションを決定することができる。

【0150】

1つの可能な実現形態において、1つの書き込み対象のデータ項目の読み取りトランザクション関連情報の3つの種類の異なる場合に、ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定する過程も異なる。

【0151】

場合1：1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は該1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプを含む。

【0152】

このような場合1に、ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定する過程としては、各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定し、ここで、第5論理ライフサイクルのタイムスタンプ下限は各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプにおける最大値よりも大きい。

【0153】

第4論理ライフサイクルは第5論理ライフサイクルを決定する前に、データノード機器のローカルトランザクションステータスリストにおいて維持されたターゲットトランザクションの最新の論理ライフサイクルである。1つの可能な実現形態において、各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定する方式としては、各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプに基づいて第4論理ライフサイクルのタイムスタンプ下限を調整し、調整した後に得られた論理ライフサイクルを第5論理ライフサイクルとして用いる。

【0154】

例示的には、各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプに基づいて第4論理ライフサイクルのタイムスタンプ下限を調整する方式としては、調整されたタイムスタンプ下限 $T \cdot B t s = \max(T \cdot B t s, y \cdot R t s + 1)$ とし、ここで、括弧内の $T \cdot B t s$ は第4論理ライフサイクルのタイムスタンプ下限を表し、 $y \cdot R t s$ は各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプにおける最大値を表し、数値1は得られた第5論理ライフサイクルのタイムスタンプ下限が各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプにおける最大値よりも大きいことを確保することに用いられる。

【0155】

いくつかの実施例において、データノード機器はローカル書き込みセットを受信した後、先ずローカル書き込みセットに対応する各書き込み対象のデータ項目の書き込み対象のトランザクション $W T$ がヌルであるか否かを検出し、ある該書き込み対象のデータ項目の書き込み対象のトランザクション $W T$ がヌルではなければ、他のトランザクションが該書き込み対象のデータ項目を修正しており、且つ該トランザクションが既に認証段階に入

10

20

30

40

50

ったことが説明される。このとき、ターゲットトランザクションをアボートして書き書き競合を除去する必要がある、すなわち協調ノード機器に Abort メッセージを携えている認証結果を返信する。各書き込み対象のデータ項目の書き込み対象のトランザクション WT がいずれもヌルであれば、ターゲットトランザクションのトランザクション識別子の値を各書き込み対象のデータ項目の書き込み対象のトランザクション WT に付与し、認証段階に入ったターゲットトランザクションが各書き込み対象のデータ項目を修正する必要があることが表される。実現では、ロックされていない CAS (Compare and Swap、比較と交換) 技術を使用して書き込み対象のデータ項目 y の書き込み対象のトランザクション WT に値を付与し、性能を向上させ、又は、先ず書き込み対象のデータ項目 y の書き込み対象のトランザクション WT をロックし、他の並行トランザクションが y を並行修正することを防止し、次にロックされた書き込み対象のトランザクション WT に値を付与する。例示的には、書き込み対象のデータ項目 y 上にアドバイザリーロックを印加し、該アドバイザリーロックは書き込み対象のデータ項目 y の書き込み対象のトランザクション WT に対する修正操作の相互排他を指示することに用いられる。

10

【0156】

場合 2 : 1 つの書き込み対象のデータ項目の読み取りトランザクション関連情報は該 1 つの書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプを含む。

【0157】

このような場合 2 に、ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び第 4 論理ライフサイクルに基づいて、ターゲットトランザクションの第 5 論理ライフサイクルを決定する過程としては、各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプ、及び第 4 論理ライフサイクルに基づいて、ターゲットトランザクションの第 5 論理ライフサイクルを決定し、ここで、第 5 論理ライフサイクルのタイムスタンプ下限は各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプにおける最大値よりも大きい。

20

【0158】

1 つの可能な実現形態において、各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプ、及び第 4 論理ライフサイクルに基づいて、ターゲットトランザクションの第 5 論理ライフサイクルを決定する方式としては、各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプに基づいて第 4 論理ライフサイクルのタイムスタンプ下限を調整し、調整した後に得られた論理ライフサイクルを第 5 論理ライフサイクルとして用いる。例示的には、各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプに基づいて第 4 論理ライフサイクルのタイムスタンプ下限を調整する方式としては、調整されたタイムスタンプ下限 $T \cdot B t s = \max (T \cdot B t s, T 1 \cdot E t s + 1)$ とし、ここで、括弧内の $T \cdot B t s$ は第 4 論理ライフサイクルのタイムスタンプ下限を表し、 $T 1 \cdot E t s$ は各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプにおける最大値を表し、数値 1 は得られた第 5 論理ライフサイクルのタイムスタンプ下限が各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプにおける最大値よりも大きいことを確保することに用いられる。

30

40

【0159】

説明する必要があるように、1 つの書き込み対象のデータ項目のターゲット読み取りトランザクションの数は 1 つ又は複数である可能性があり、1 つの書き込み対象のデータ項目のターゲット読み取りトランザクションの数が複数である場合については、上記 $T 1 \cdot E t s$ とは全部の書き込み対象のデータ項目の全部のターゲット読み取りトランザクションの終了タイムスタンプにおける最大値を指す。

【0160】

このような方式に基づいてターゲットトランザクションの書き込み操作の発生をターゲ

50

ット読み取りトランザクションの読み取り操作の後に遅延させて、読み書き競合を回避することを可能にする。

【0161】

場合3：1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は該1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び該1つの書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプを含む。

【0162】

このような場合3に、ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定する過程としては、各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプに基づいて第4論理ライフサイクルを連続的に2回調整し、2回調整した後に得られた論理ライフサイクルをターゲットトランザクションの第5論理ライフサイクルとして用いる。本願の実施例は2回調整の前後順序を限定せず、例示的には、先ず各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプに基づいて第4論理ライフサイクルを調整し、次に各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプに基づいて1回調整した後に得られた論理ライフサイクルを調整する。もちろん、いくつかの実施例において、さらに先ず各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプに基づいて第4論理ライフサイクルを調整し、次に各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプに基づいて1回調整した後に得られた論理ライフサイクルを調整することができる。

10

20

【0163】

例示的には、2回調整した後に得られた論理ライフサイクルを第5論理ライフサイクルとして用いる場合については、1回調整した後に論理ライフサイクルを得た後、シリアル化可能分離レベルであると、先ず得られた論理ライフサイクルのタイムスタンプ下限がタイムスタンプ上限よりも小さいか否かを認証し、イエスであれば、次の調整を継続し、ノーであれば、直接ローカル認証が失敗するとみなされ、協調ノード機器にAbortメッセージを携えている認証結果を返信する。

30

【0164】

第5論理ライフサイクルを得た後、第5論理ライフサイクルのタイムスタンプ下限が第5論理ライフサイクルのタイムスタンプ上限よりも小さいか否かを認証することによって、第5論理ライフサイクルが有効であるか否かを判断する。第5論理ライフサイクルが有効であることに応答して、認証に合格したことを指示することに用いられる認証結果をターゲットトランザクションの認証結果として用い、第5論理ライフサイクルが無効であることに応答して、認証に合格しなかったことを指示することに用いられる認証結果をターゲットトランザクションの認証結果として用いる。認証に合格したことを指示することに用いられる認証結果をターゲットトランザクションの認証結果として用いる場合については、データノード機器のローカル認証フィードバックメッセージlv_mにおいてターゲットトランザクションのデータノード機器上で得られた最新の論理ライフサイクル(すなわち第5論理ライフサイクル)のタイムスタンプ下限B_{ts}、及びタイムスタンプ上限E_{ts}が記録される。例示的には、認証に合格しなかったことを指示することに用いられる認証結果はAbortメッセージを携えている認証結果である。

40

【0165】

例示的な実施例において、第5論理ライフサイクルが有効であると判定した場合、ターゲットトランザクションのローカル認証に合格するとみなされ、データノード機器はローカルトランザクションステータスリストにおけるターゲットトランザクションのステータス情報を更新し、ターゲットトランザクションのトランザクションステータスをValidated(認証に合格する)に更新し、すなわち、T.Status = Validated

50

e dとする。例示的な実施例において、ターゲットトランザクションのローカル認証に合格すると判定した後、データノード機器は書き込み対象のデータ項目の更新値に基づき、書き込み対象のデータ項目の新たなバージョンデータを作成する。例示的な実施例において、作成された新たなバージョンデータに該新たなバージョンデータがグローバルにコミットされていないことを指示することに用いられる第1マークを設定する。第1マークを有する新たなバージョンデータは外部に対して不可視である。

【0166】

説明する必要があるように、ターゲットトランザクションのデータノード機器におけるローカル認証に合格しなければ、データノード機器のローカルトランザクションステータスリストにおけるターゲットトランザクションのトランザクションステータスを `Aborted` (アボート) に更新し、すなわち、`T.Status = Aborted` とする必要がある。

10

【0167】

1つの可能な実現形態において、1つの書き込み対象のデータ項目のアクティブトランザクション集合において、ターゲット読み取りトランザクションを含むことに加えて、さらに動作している読み取りトランザクションを含み、ターゲットトランザクションの第5論理ライフサイクルに基づき動作している読み取りトランザクションの論理ライフサイクルを調整する必要がある、動作している読み取りトランザクションにターゲットトランザクションが新たに書き込んだデータを読ませず、それにより読み書き競合現象を回避し、トランザクションが正確に実行されることを確保する。例示的には、動作している読み取りトランザクションとはアクティブトランザクション集合におけるトランザクションステータスが `Running` (動作している) のトランザクションを指す。動作している読み取りトランザクションの論理ライフサイクルを調整する方式としては、動作している読み取りトランザクションの論理ライフサイクルのタイムスタンプ上限をターゲットトランザクションの第5論理ライフサイクルのタイムスタンプ下限よりも小さくする。動作している読み取りトランザクションが `T2` であると仮定すると、更新方式としては、`T2.Ets = min(T2.Ets, T.Bts - 1)` とする。ある動作している読み取りトランザクションの更新された論理ライフサイクルのタイムスタンプ下限がタイムスタンプ上限以上であると、該動作している読み取りトランザクションにアボートすべきであるように通知する。

20

30

【0168】

上記トランザクション認証段階から分かるように、ターゲットトランザクションの認証過程において、通信は主に協調ノード機器と関連するデータノード機器との間で発生し、通信は主に、協調ノード機器が個々の関連するデータノード機器にトランザクション認証要求、及びローカル書き込みセットを送信するステップと、関連するデータノード機器が認証結果を協調ノード機器にフィードバックするステップと、を含む。従って、ターゲットトランザクションの認証段階では、 m (m は1以上の整数である) がターゲットトランザクション T に関連するデータノード機器の数であると仮定すると、最大 $2m$ 回の通信を行う必要がある、最大通信量は $m \times$ (トランザクション認証要求メッセージのサイズ + 認証結果メッセージのサイズ) + グローバル書き込みセットのサイズとして表されてもよい。

40

【0169】

ステップ208において、協調ノード機器はデータノード機器から返信されたターゲットトランザクションの認証結果に基づいて、ターゲットトランザクションの処理命令を決定し、データノード機器に処理命令を送信し、処理命令はコミット命令、又はアボート命令である。

【0170】

協調ノード機器はデータノード機器から返信された認証結果を受信した後、受信した認証結果に基づきターゲットトランザクションがグローバル認証に合格するか否かを判断し、更にターゲットトランザクションの処理命令を決定し、データノード機器に処理命令を

50

送信する。ここで、処理命令はコミット命令、又はアボート命令である。

【0171】

1つの可能な実現形態において、データノード機器の数は1つ又は複数であり、データノード機器の数が複数である場合、個々のデータノード機器はいずれも1つの認証結果を返信する。

【0172】

データノード機器の数が少なくとも2つである場合について、データノード機器から返信されたターゲットランザクションの認証結果に基づいて、ターゲットランザクションの処理命令を決定する過程としては、少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果において認証に合格しなかったことを指示することに用いられる認証結果が存在することに応答して、アボート命令をターゲットランザクションの処理命令として用いる。少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果はいずれも認証に合格したことを指示することに応答して、少なくとも2つの認証結果が携えている論理ライフサイクルの共通部分をターゲット論理ライフサイクルとして用いて、ターゲット論理ライフサイクルが有効であることに応答して、コミット命令をターゲットランザクションの処理命令として用いて、ターゲット論理ライフサイクルが無効であることに応答して、アボート命令をターゲットランザクションの処理命令として用いる。

10

【0173】

例示的な実施例において、認証に合格しなかったことを指示することに用いられる認証結果はAbortメッセージを携えている認証結果であり、ある認証結果がAbortメッセージを携えているのではなく、論理ライフサイクル（すなわちステップ207において決定された第5論理ライフサイクル）を携えていれば、該認証結果は認証に合格したことを指示する。つまり、協調ノード機器は受信された認証結果に基づきターゲットランザクションがグローバル認証に合格することを可能にするか否かを判断する過程において、受信された認証結果においてAbortメッセージを携えている少なくとも1つの認証結果が存在し、すなわちIsAbortフィールドが1に等しいlvnが存在すると、ターゲットランザクションが全部のローカル認証に合格しないことを表明し、ターゲットランザクションのグローバル認証に合格せず、ターゲットランザクションがグローバルアボートを行う必要がある。このような場合に、アボート命令をターゲットランザクションの処理命令として用いる。協調ノード機器は第1ランザクションステータスリストにおけるターゲットランザクションのグローバルランザクションステータスをGaborting（グローバルにアボートしている）に更新する。協調ノード機器はデータノード機器にアボート命令を送信して、データノード機器にローカルアボートを行うように通知する。例示的には、処理命令はコミット/アボートメッセージcoarmの書き込みによって送信され、処理命令がアボート命令である場合、coarmにおけるIsAbortフィールドは1に等しく、すなわちcoarm.IsAbort=1とする。

20

30

【0174】

受信された認証結果においてAbortメッセージを携えている認証結果が存在しないか、又は受信された認証結果がいずれも論理ライフサイクルを携えていれば、ターゲットランザクションが全部のローカル認証に合格することが説明される。このような場合に、協調ノード機器は受信された各認証結果が携えている論理ライフサイクルの共通部分を計算して、ターゲット論理ライフサイクルを得る。ターゲット論理ライフサイクルのタイムスタンプ下限がターゲット論理ライフサイクルのタイムスタンプ上限以上であれば、ターゲット論理ライフサイクルが無効であることが説明され、ターゲットランザクションのグローバル認証に合格しなかったと判定し、ターゲットランザクションがグローバルアボートを行う必要があり、協調ノード機器はアボート命令をターゲットランザクションの処理命令として用いる。この他、協調ノード機器はさらに第1ランザクションステータスリストにおけるターゲットランザクションのグローバルランザクションステータスをGaborting（グローバルにアボートしている）に更新し、協調ノード機器

40

50

はデータノード機器にアポート命令を送信して、データノード機器にローカルアポートを行うように通知する。

【0175】

ターゲット論理ライフサイクルのタイムスタンプ下限が、ターゲット論理ライフサイクルのタイムスタンプ上限よりも小さければ、ターゲット論理ライフサイクルが有効であることが説明され、ターゲットトランザクションのグローバル認証に合格したと判定し、協調ノード機器はターゲット論理ライフサイクルから1つのタイムスタンプをランダムに選択してターゲットトランザクションの論理コミットタイムスタンプ Cts に値を付与し、例えば、ターゲット論理ライフサイクルのタイムスタンプ下限をターゲットトランザクションの論理コミットタイムスタンプとして選択する。

10

【0176】

論理コミットタイムスタンプを決定した後、協調ノード機器は第1トランザクションステータスリストにおけるターゲットトランザクション T の論理ライフサイクルのタイムスタンプ下限、及び論理ライフサイクルのタイムスタンプ上限をいずれも論理コミットタイムスタンプに更新し、すなわち、 $T.Bts = T.Ets = T.Cts$ とする。この他、第1トランザクションステータスリストにおけるターゲットトランザクションのグローバルトランザクションステータスを $Gcommitted$ (グローバル認証に合格する) に更新し、グローバルタイムスタンプ生成クラスがターゲットトランザクションにグローバルコミットタイムスタンプを割り当てるように要求し、第1トランザクションステータスリストにおけるターゲットトランザクションのグローバルコミットタイムスタンプ Gts フィールドにおいて記録する。この他、協調ノード機器はコミット命令をターゲットトランザクションの処理命令として用い、データノード機器にコミット命令を送信して、データノード機器にターゲットトランザクションをコミットするように通知する。例示的には、処理命令がコミット/アポートメッセージ $coarm$ の書き込みによって送信される場合については、処理命令がコミット命令である場合、 $coarm$ における $IsAbort$ フィールドは0に等しく、すなわち $coarm.IsAbort = 0$ とし、 $coarm$ における Cts 、及び Gts フィールドにおいてそれぞれターゲットトランザクションの論理コミットタイムスタンプ、及びターゲットトランザクションのグローバルコミットタイムスタンプが記録される。

20

【0177】

ステップ209において、データノード機器は協調ノード機器から送信されたターゲットトランザクションの処理命令を受信したことに応答して、処理命令を実行し、処理命令はコミット命令、又はアポート命令である。

30

【0178】

データノード機器は処理命令を受信した後、処理命令を実行する。データノード機器が処理命令を実行する段階はトランザクションコミット、又はアポート操作終了段階である。

【0179】

処理命令がコミット命令である場合、ターゲットトランザクションのグローバル認証に合格し、コミット段階に入ることが説明され、すなわちターゲットトランザクションのデータに対する更新をデータベースにおいて永続化し、かついくつかの後続のクリーンアップ作業を行う。例示的な実施例において、データノード機器は協調ノード機器から送信されたコミット命令を受信した後に、以下の操作A～操作Eを実行することができる。

40

【0180】

操作A：ターゲットトランザクションのローカル読み取りセットに対応する個々のデータ項目 x に対しては、データ項目 x の最大読み取りトランザクションタイムスタンプ Rts を修正し、データ項目 x の最大読み取りトランザクションタイムスタンプ Rts をターゲットトランザクションの論理コミットタイムスタンプ Cts よりも大きく又は等しくし、すなわち、 $x.Rts = \max(x.Rts, T.Cts)$ とし、ターゲットトランザクションのトランザクション識別子 TID を該データ項目 x のアクティブトランザクシ

50

ンリスト `RTList` から削除する。

【0181】

操作 B：ターゲットトランザクションのローカル書き込みセットに対応する個々のデータ項目 y に対しては、以下の操作を実行する。a) データ項目 y のオリジナルの作成タイムスタンプ `Wts` をターゲットトランザクションの論理コミットタイムスタンプ `T.Cts` に修正する。b) データ項目 y の最大読み取りトランザクションタイムスタンプをオリジナルの最大読み取りトランザクションタイムスタンプとターゲットトランザクションの論理コミットタイムスタンプにおける最大値に更新し、すなわち、 $y.Rts = \max(y.Rts, T.Cts)$ とする。c) データ項目 y をデータベースにおいて永続化し、且つデータ項目 y のマークを第 1 マークから第 2 マークに修正し、第 2 マークは外部に対して可視であることを指示することに用いられる。d) データ項目 y のアクティブトランザクションリスト `RTList` のコンテンツをクリアする。e) データ項目 y の書き込み対象のトランザクション `WT` のコンテンツをクリアする。

10

【0182】

操作 C：ターゲットトランザクションのローカル読み取りセット及びローカル書き込みセットをクリアする。

【0183】

操作 D：ローカルトランザクションステータスリストにおけるターゲットトランザクションの論理ライフサイクルのタイムスタンプ下限、及びタイムスタンプ上限をいずれもターゲットトランザクションの論理コミットタイムスタンプに更新し、すなわち、 $T.Bts = T.Ets = T.Cts$ とする。ローカルトランザクションステータスリストにおけるターゲットトランザクションのトランザクションステータスを `Committed` (コミットが完了する) に更新する。説明する必要があるように、このとき、ローカルトランザクションステータスリストは、トランザクションの一致性を確保することに用いられ、グローバルトランザクションステータスの同期に関する必要がない。

20

【0184】

操作 E：協調ノード機器にコミットに成功した `ACK (Acknowledge Character, 確認文字)` を返信する。

【0185】

協調ノード機器は全部のデータノード機器から返信されたコミットに成功した `ACK` を受信した後、第 1 トランザクションステータスリストにおけるターゲットトランザクションのグローバルトランザクションステータスを `Gcommitted` (グローバルコミットが完了する) に修正し、次に協調ノード機器はデータノード機器にステータス情報クリーンアップ命令を送信し、データノード機器にローカルトランザクションステータスリストからターゲットトランザクションのステータス情報を削除させる。

30

【0186】

処理命令がアボート命令である場合、ターゲットトランザクションのグローバル認証に合格しないことが説明され、グローバルアボート段階に入る必要があり、すなわちターゲットトランザクションをアボートし、かつ相応なクリーンアップ作業を行う。例示的には、クリーンアップ作業のコンテンツは、ターゲットトランザクションのトランザクション識別子 `TID` をターゲットトランザクションのローカル読み取りセットに対応する個々のデータ項目 x のアクティブトランザクションリスト `RTList` から削除することと、ターゲットトランザクションのローカル書き込みセットに対応する個々のデータ項目 y に対応する新たな作成データをクリーンアップし、且つデータ項目 y の書き込み対象のトランザクション `WT` のコンテンツをクリアすることと、ターゲットトランザクションのローカル読み取りセット及びローカル書き込みセットをクリアすることと、ローカルトランザクションステータスリストにおけるターゲットトランザクションのトランザクションステータスを `Aborted` (アボートが完了する) に更新することと、協調ノード機器にアボートが完了した `ACK` を返信することと、を含む。

40

【0187】

50

協調ノード機器は全部のデータノード機器から返信されたアポートが完了したACKを受信した後、第1トランザクションステータスリストにおけるターゲットトランザクションのグローバルトランザクションステータスをG a b o r t e d (グローバルアポートが完了する)に修正し、次に協調ノード機器はデータノード機器にステータス情報クリーンアップ命令を送信し、データノード機器にローカルトランザクションステータスリストからターゲットトランザクションのステータス情報を削除させる。1つの可能な実現形態において、協調ノード機器はデータノード機器にステータス情報クリーンアップ命令を一括して送信して、通信回数を減少させる。

【0188】

上記内容から明らかなように、ターゲットトランザクションのコミット/アポート段階では、通信は主に協調ノード機器と関連するデータノード機器との間で発生し、通信は主に、協調ノード機器が個々の関連するデータノード機器にコミット/アポート命令を送信するステップと、個々の関連するデータノード機器が協調ノード機器に相応なコミット/アポート完了メッセージ(ACK)を送信するステップと、を含む。従って、コミット/アポート段階では、最大2m回の通信を行い、通信量のサイズは $m \times (\text{コミット/アポート命令メッセージのサイズ} + \text{コミット/アポート完了メッセージのサイズ})$ であり、ここでm(mは1以上の整数である)はターゲットトランザクションTに関連するデータノード機器の数である。

10

【0189】

説明する必要があるように、本願の実施例はターゲットトランザクションが読み書き操作に関することを例として紹介したが、本願の実施例はこれに限定されず、ターゲットトランザクションが読み取り操作のみに関するか、又は書き込み操作のみに関する場合については、依然として本願の実施例が提供するトランザクション処理方法に基づきトランザクションに対する処理を実現することを可能にするが、本願の実施例は詳細な説明を省略する。

20

【0190】

上記ステップ201~ステップ209に基づいてトランザクションを処理する過程はトランザクションの分散化処理を実現し、並行トランザクション間の競合操作に起因するデータ異常の問題を解決することを可能にする。実現原理から見ると、本願の実施例が提供するトランザクション処理方法は主にOCC(Optimistic Concurrency Control、楽観的並行性制御)のアルゴリズムフレームワークを応用し、DTA(Dynamic Timestamp Allocation、動的割り当てタイムスタンプ)アルゴリズムを組み合わせ、ネットワークで伝送されたトランザクションデータ情報を減少させ、分散型トランザクションの認証効率を向上させ、分散型トランザクションの並行処理能力を高める。その他、さらにMVCC(Multi-Version Concurrency Control、マルチバージョン並行性制御)を組み合わせ、ロックされていないデータ読み書きを実現し、それにより局所ノード機器の並行処理能力を高める。ここで、DTAアルゴリズムはTO(Timestamp Ordering、タイムスタンプオーダリング)アルゴリズムに属し、トランザクションの論理ライフサイクルのタイムスタンプ下限、及びタイムスタンプ上限を動的に調整することができる。

30

40

【0191】

本願の実施例が提供する方法はデータ記憶フォーマットに影響されず、本願の実施例における分散型データベースシステムはキー値型のデータ記憶フォーマット(KVデータ記憶フォーマット)(たとえば、HBaseデータベースシステムにおけるデータ記憶フォーマット)をサポートするだけでなく、セグメントページ型のデータ記憶フォーマット(たとえば、PostgreSQL及びMySQL/InnoDBデータベースシステムにおけるデータ記憶フォーマット)をサポートする。

【0192】

例示的な実施例において、セグメントページ型のデータ記憶フォーマットに対しては、

50

ノード機器内にデータバッファ領域を確立し、共有された記憶システムから伝送されたデータをバッファし、それによりデータを次回取得する速度を加速させ、バッファフォーマットは下層のデータ記憶フォーマットと一致するように保持される。共有された記憶システムから伝送されたデータは、ローカルのデータバッファ領域にバッファされ、ローカルデータバッファ領域がいっぱいになるか、又はダーティデータを共有された記憶システムにフラッシュバックする必要があるか、又はバッファが無効になる（たとえば、他のノード機器上で同様のデータが修正される）ときまで、トランザクションが終了してもクリーンアップされない。

【0193】

トランザクションがコミットされる前に、個々のノード機器は共有された記憶システムにトランザクションログ（たとえば、WALログ）を算出し、トランザクションログは共有された記憶システムにLSN値を請求し、該値はグローバルで一意的で且つ遞増する1つの値である。異なるデータ記憶フォーマットでは、トランザクション処理過程において生じたトランザクションログは異なるフォーマットを有する。例示的には、データ記憶フォーマットがKVデータ記憶フォーマットである場合、トランザクションログのフォーマットは図3に示される。

10

【0194】

データベースシステムによって維持された大規模なテーブルで分割された各エリア（Region）は1つのログファイルを共有し、単一のエリアはログにおいて時間順序に応じて記憶されるが、複数のエリアは完全に時間順序に応じたものではない可能性がある。個々のログの最小ユニットはログキー（HLogKey）及びログ編集（WALEdit）の2つの部分からなる。ここで、HLogKeyはシーケンス番号（sequence id）、タイムスタンプ（timestamp）、クラスタ番号（cluster id）、エリア名（region name）及びテーブル名（table name）などからなり、WALEditは一連のキー値ペア（Key Value）からなり、1行におけるすべての列（すなわちすべてのKey Value）に対する更新操作はいずれも同一のWALEditオブジェクトにおいて含まれ、これは主に1行複数列に書き込むときの原子性を実現するためである。sequence idは、1つの記憶レベルの自己増加シーケンス番号であり、エリアのデータ復旧及び期限切れログのクリアはいずれもそれに依存し、例示的には、sequence idとはトランザクションログのLSN値を指す。

20

30

【0195】

例示的には、データ記憶フォーマットがセグメントページ型のデータ記憶フォーマットである場合、トランザクションログのフォーマットは図4に示される。個々のRegionは1つのログファイルを共有し、単一のRegionはログにおいて時間順序に応じて記憶され、且つ複数のRegionは完全に時間順序に応じたものではない可能性がある。個々のログの最小ユニットはHLogKey及びWALEditの2つの部分からなるのではなく、1つのログレコード（XLog Record）からなる。

【0196】

XLog Recordは2つの部分で構成され、第1部分はヘッダ情報であり、サイズが固定され（たとえば、24B（Bytes、バイト）、対応する構造体はXLog Recordであり、第2部分はログレコードデータ（XLog Record data）である。

40

【0197】

XLog Recordは記憶されたデータコンテンツに応じて分割され、主に以下の3つの種類に分けられる。

【0198】

第1種類：Record for backup block（バックアップブロックレコード）：full-write-page（全書き込みページ）を記憶するblock（ブロック）であり、このようなタイプのレコードはページ部分の書き込みの問題を解

50

決するためである。checkpoint (検出ポイント) が完了した後にデータページを1回目に修正し、この変更を記録してトランザクションログファイルに書き込むときにページ全体に書き込む (相応な初期化パラメータを設定する必要があり、デフォルトで開いている)。

【0199】

第2種類: Record for tuple data block (タプルデータブロックレコード): ページにおけるタプル変更を記憶することに用いられる。

【0200】

第3種類: Record for Checkpoint (チェックポイントレコード): checkpoint が発生した場合、トランザクションログファイルにおいて checkpoint 情報 (ここで Redo point (やり直しポイント) を含む) を記録する。 10

【0201】

XLog Record data は実際のデータを記憶する場所であり、以下の4つの部分からなる。

【0202】

第1部分: 0 - N個の XLogRecordBlockHeader (ログレコードブロックヘッダ) であり、各々の XLogRecordBlockHeader は1つの block data (ブロックデータ) に対応する。BKPBLOCK_HAS_IMAGE マークを設定すると、XLogRecordBlockHeader 構造体の後に XLogRecordBlockImageHeader 構造体が続き、BKPBLOCK_HAS_HOLE & BKPIIMAGE_IS_COMPRESSED マークを設定すると、XLogRecordBlockHeader 構造体の後に XLogRecordBlockCompressHeader 構造体が続き、BKPBLOCK_SAME_REL マークを設定しなければ、XLogRecordBlockHeader 構造体の後に RelFileNode が続く。例示的には、XLogRecordBlockHeader 構造体の後にさらに BlockNumber (ブロック番号) が続いてもよい。 20

【0203】

第2部分: XLogRecordDataHeader [Short | Long] (ログレコードデータヘッダ [ショート | ロング]) であり、データのサイズ < 256 Bytes とすると、Short フォーマットを使用し、そうでなければ、Long フォーマットを使用する。 30

【0204】

第3部分: block data (ブロックデータ): full-write-page data (全書き込みページデータ) 及び tuple data (タプルデータ)。full-write-page data に対しては、圧縮をイネーブルにすると、データが圧縮記憶され、圧縮された後に該 page に関連するメタデータを XLogRecordBlockCompressHeader (ログレコードブロック圧縮ヘッダ) に記憶する。 40

【0205】

第4部分: main data (主要データ): checkpoint などのログデータを記録する。 40

【0206】

例示的には、XLog Record は以下のように定義される。

ヘッダ情報 (固定サイズの XLogRecord 構造体)

XLogRecordBlockHeader 構造体

XLogRecordBlockHeader 構造体

...

XLogRecordDataHeader [Short | Long] 構造体

block data

block data

...

main data

【0207】

1つの可能な実現形態において、データ記憶フォーマットがセグメントページ型のデータ記憶フォーマットである場合については、並行トランザクションが異なるノード機器（ES）上で処理され、且つ同一のページ上の異なるデータ項目を修正する場合、ページレベルの競合が発生してデータ上書きの問題を引き起こすようになる。たとえば、トランザクションTaはノード機器ES-1上でX=2のデータ項目を修正し、トランザクションTbはノード機器ES-2上でX=3のデータ項目を修正し、X=2とX=3のデータ項目は同一のページ（page）上にあり、このとき、トランザクション処理メカニズムは、並行に動作され、トランザクションを同時実行し、トランザクションレベルでデータ異常が存在しない。しかし、ページレベルでは、ES-1とES-2のどちらでフラッシュアウトされたトランザクションログを選ぶべきであるかという選択が存在し、これは同一の物理ページに対する修正が両立できないという問題の出現をもたらす。

10

【0208】

セグメントページ型のデータ記憶フォーマットをサポートするトランザクションログにおいて、1つのセグメントページ型のリスト（list）が増加され、本セグメントログにおけるページのアドレス（たとえば、ファイル番号、テーブルスペース番号、ファイルにおける相対オフセットなど）及び個々のページ上で書き込み操作を実行しているトランザクション識別子がマーキングされる。トランザクションログが最下位層に共有された記憶システムにフラッシュアウトされる場合、認証機器はすべての並行トランザクションの共有された記憶システムのトランザクションにおけるlistにコミットされたページが重なるか否かをチェックする。もし重なっていれば、並行トランザクションが同じページを書き込んだことが表明され（書き込むのが同じデータ項目であれば、トランザクション認証段階では、トランザクション競合が存在すると検出しかつ既にアポートによって競合を解消している）、書き込むのが同じページ上の異なるデータ項目であり、このときにページレベルの競合が存在し、データ上書きイベントが発生し、そのうちの1つのノード機器ESに対応するトランザクションをアポートする必要がある。対応するトランザクションがアポートされたノード機器ESのトランザクションログをフラッシュアウトしなくなり、問題の発生を回避する。

20

30

【0209】

例示的な実施例において、上記ページレベルの競合認証を実行する主体は分散型データベースシステムにおける認証機器である。該認証機器はいずれかのノード機器と同一の物理マシン上にあってもよく、独立した機器であってもよく、本願の実施例はこれを限定しない。

【0210】

本願の実施例が提供するトランザクション処理方法に基づいて、分散型データベースシステムが分散型トランザクションをサポートするだけでなく、グローバルで一致性のあるマルチ読み取りを達成することを可能にし、分散化トランザクション処理技術によって性能にも配慮することを可能にする。良好なトランザクション属性特徴付きのグローバルで一致性のあるマルチ読み取り及び一致性のあるマルチ書き込み能力を備える。本願の実施例が提供するトランザクション処理方法に基づいて、share-diskアーキテクチャに基づく分散型データベースシステム、たとえば有名なNoSQL（Non-relational SQL、一般的に非リレーショナルデータベースを指す）下でのHBaseデータベースシステムに、分散化された分散型トランザクション処理手段を提供できることにより、HBaseと類似するデータベースシステムはエリア間、ノード間の高効率なトランザクション処理能力を備える。

40

【0211】

本願の実施例においては、各ノード機器のそれぞれに対応するトランザクション割り当

50

て指標に基づきターゲットトランザクションを協調処理することに用いられる協調ノード機器を決定し、トランザクションの割り当て過程はトランザクションに関するデータ項目を考慮する必要がなく、データ項目の分布状況を考慮する必要もない。このような方式に基づいて、個々のノード機器はいずれも分散化された機器としてトランザクションを協調処理することを可能にすることにより、トランザクションをノード間で処理することを可能にし、トランザクションの処理効率を向上させることに有利であり、トランザクション処理の信頼性が比較的高く、データベースシステムのシステム性能を高めることに有利である。

【0212】

本願の実施例はトランザクション処理システムを提供し、該トランザクション処理システムは協調ノード機器と、データノード機器とを含み、協調ノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、協調ノード機器は少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定され、データノード機器は少なくとも2つのノード機器のうちターゲットトランザクションの処理に関与することに用いられるノード機器であり、

ここで、協調ノード機器は、ターゲットトランザクションのトランザクション情報を取得することと、ターゲットトランザクションのトランザクション情報に基づいて、データノード機器にデータ読み取り要求を送信することと、に用いられ、

データノード機器は、協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得し、データ読み取り結果を協調ノード機器に返信することに用いられ、

協調ノード機器はさらに、データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信することに用いられ、

データノード機器はさらに、協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、ターゲットトランザクションの認証結果を取得し、ターゲットトランザクションの認証結果を協調ノード機器に返信することに用いられ、

協調ノード機器はさらに、データノード機器から返信されたターゲットトランザクションの認証結果に基づいて、ターゲットトランザクションの処理命令を決定し、データノード機器に処理命令を送信することに用いられ、処理命令はコミット命令、又はアボート命令であり、

データノード機器はさらに、協調ノード機器から送信されたターゲットトランザクションの処理命令を受信したことに応答して、処理命令を実行することに用いられる。

【0213】

1つの可能な実現形態において、データ読み取り結果は第2論理ライフサイクルを携えており、第2論理ライフサイクルはデータノード機器によってデータ読み取り要求が携えているターゲットトランザクションの第1論理ライフサイクルに基づき決定され、第1論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成され、協調ノード機器はさらに、第1論理ライフサイクルのタイムスタンプ下限と第2論理ライフサイクルのタイムスタンプ下限における最大値をターゲットトランザクションの第3論理ライフサイクルのタイムスタンプ下限として用いることと、第1論理ライフサイクルのタイムスタンプ上限と第2論理ライフサイクルのタイムスタンプ上限における最小値をターゲットトランザクションの第3論理ライフサイクルのタイムスタンプ上限として用いることと、第3論理ライフサイクルが有効であることに応答して、データノード機器に第3論理ライフサイクルを携えているトランザクション認証要求を送信することと、第3論理ライフサイクルが有効であることは、第3論理ライフサイクルのタイムスタンプ下限が第3論理ライフサイクルのタイムスタンプ上限よりも小さいことを指示することに用いられる、ことと、に用いられる。

10

20

30

40

50

【 0 2 1 4 】

1つの可能な実現形態において、データノード機器の数は少なくとも2つであり、協調ノード機器はさらに、少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果において認証に合格しなかったことを指示することに用いられる認証結果が存在することに応答して、アポート命令をターゲットトランザクションの処理命令として用いることと、少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果はいずれも認証に合格したことを指示することに応答して、少なくとも2つの認証結果が携えている論理ライフサイクルの共通部分をターゲット論理ライフサイクルとして用いることと、ターゲット論理ライフサイクルが有効であることに応答して、コミット命令をターゲットトランザクションの処理命令として用いることと、ターゲット論理ライフサイクルが無効であることに応答して、アポート命令をターゲットトランザクションの処理命令として用いることと、に用いられる。

10

【 0 2 1 5 】

1つの可能な実現形態において、データ読み取り要求はターゲットトランザクションの第1論理ライフサイクルを携えており、第1論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成され、データノード機器は、第1論理ライフサイクルに基づいて、データ読み取り要求によって指示された読み取り対象のデータ項目の可視バージョンデータを決定することと、可視バージョンデータの作成タイムスタンプ、及び第1論理ライフサイクルに基づいて、ターゲットトランザクションの第2論理ライフサイクルを決定することと、第2論理ライフサイクル、及び可視バージョンデータを携えている結果をデータ読み取り結果として用いることと、に用いられる。

20

【 0 2 1 6 】

1つの可能な実現形態において、トランザクション認証要求はターゲットトランザクションの第3論理ライフサイクルを携えており、第3論理ライフサイクルは協調ノード機器によって第1論理ライフサイクル、及び第2論理ライフサイクルに基づいて決定された有効論理ライフサイクルであり、データノード機器はさらに、第3論理ライフサイクルのタイムスタンプ下限と第2論理ライフサイクルのタイムスタンプ下限における最大値をターゲットトランザクションの第4論理ライフサイクルのタイムスタンプ下限として用いることと、第3論理ライフサイクルのタイムスタンプ上限と第2論理ライフサイクルのタイムスタンプ上限における最小値をターゲットトランザクションの第4論理ライフサイクルのタイムスタンプ上限として用いることと、第4論理ライフサイクルが有効であることに応答して、ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定することと、第5論理ライフサイクルが有効であることに応答して、認証に合格したことを指示することに用いられる認証結果をターゲットトランザクションの認証結果として用いることと、第5論理ライフサイクルが無効であることに応答して、認証に合格しなかったことを指示することに用いられる認証結果をターゲットトランザクションの認証結果として用いることと、に用いられる。

30

【 0 2 1 7 】

1つの可能な実現形態において、1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプを含み、1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプは、1つの書き込み対象のデータ項目を読み取った各読み取りトランザクションの論理コミットタイムスタンプにおける最大値を指示することに用いられ、データノード機器はさらに、各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定することに用いられ、第5論理ライフサイクルのタイムスタンプ下限は各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプにおける最大値よりも大きい。

40

【 0 2 1 8 】

50

1つの可能な実現形態において、1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は1つの書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプを含み、ターゲット読み取りトランザクションはローカル認証に合格するか、又はコミット段階にある読み取りトランザクションであり、ターゲット読み取りトランザクションの終了タイムスタンプはターゲット読み取りトランザクションの論理ライフサイクルのタイムスタンプ上限であり、データノード機器はさらに、各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプ、及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定することに用いられ、第5論理ライフサイクルのタイムスタンプ下限は各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプにおける最大値よりも大きい。

10

【0219】

上記実施例が提供するシステムは方法実施例と同一の構想に属し、その具体的な実現過程については、詳しくは方法実施例を参照すればよいため、ここでは詳細な説明を省略する。

【0220】

図5に参照されるように、本願の実施例はトランザクション処理装置を提供し、該装置は、

ターゲットトランザクションの割り当て要求に応答して、同一の記憶システムを共有する少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定することに用いられる第1決定ユニット501であって、1つのノード機器に対応するトランザクション割り当て指標は該1つのノード機器に新たなトランザクションを割り当ててマッピング度を指示することに用いられる、第1決定ユニット501と、

20

少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、少なくとも2つのノード機器のうちターゲットトランザクションの協調ノード機器を決定し、協調ノード機器によってターゲットトランザクションを協調処理することに用いられる第2決定ユニット502と、を含む。

【0221】

1つの可能な実現形態において、第1決定ユニット501は、トランザクション割り当てモードを決定することであって、トランザクション割り当てモードはトランザクションのビジー度に基づく割り当て、機器のビジー度に基づく割り当て、及びハイブリッドのビジー度に基づく割り当てのうちのいずれか1つを含む、ことと、トランザクション割り当てモードによって指示された決定方式に基づき、少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定することと、に用いられる。

30

【0222】

1つの可能な実現形態において、トランザクション割り当てモードはハイブリッドのビジー度に基づく割り当てを含み、第1決定ユニット501はさらに、第1ノード機器のトランザクション処理数、第1ノード機器の機器リソース使用率、トランザクション処理数の重み、機器リソース使用率の重み、及び重み調節パラメータに基づいて、第1ノード機器に対応するトランザクション割り当て指標を決定することに用いられ、第1ノード機器は少なくとも2つのノード機器のうちのいずれか1つのノード機器である。

40

【0223】

1つの可能な実現形態において、該装置は、

協調ノード機器の機器識別情報を割り当て要求を発した端末に送信することに用いられる送信ユニットであって、端末は協調ノード機器の機器識別情報に基づき、ターゲットトランザクションのトランザクション情報を協調ノード機器に送信し、協調ノード機器によってトランザクション情報に基づいてターゲットトランザクションを協調処理することに用いられる、送信ユニットをさらに含む。

【0224】

1つの可能な実現形態において、分散型データベースシステムはキー値型のデータ記憶

50

フォーマットとセグメントページ型のデータ記憶フォーマットをサポートする。

【0225】

本願の実施例において、各ノード機器のそれぞれに対応するトランザクション割り当て指標に基づきターゲットトランザクションを協調処理することに用いられる協調ノード機器を決定し、トランザクションの割り当て過程はトランザクションに関するデータ項目を考慮する必要がなく、データ項目の分布状況を考慮する必要もない。このような方式に基づいて、個々のノード機器はいずれも分散化された機器としてトランザクションを協調処理することを可能にすることにより、トランザクションをノード間で処理することを可能にし、トランザクションの処理効率を向上させることに有利であり、トランザクション処理の信頼性が比較的高く、データベースシステムのシステム性能を高めることに有利である。

10

【0226】

図6に参照されるように、本願の実施例はトランザクション処理装置を提供し、該装置は、

ターゲットトランザクションのトランザクション情報を取得することに用いられる取得ユニット601と、

ターゲットトランザクションのトランザクション情報に基づいて、データノード機器にデータ読み取り要求を送信することに用いられる第1送信ユニット602であって、データノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションの処理に関与することに用いられるノード機器である、第1送信ユニット602と、

20

データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信することに用いられる第2送信ユニット603と、

データノード機器から返信されたターゲットトランザクションの認証結果に基づいて、ターゲットトランザクションの処理命令を決定することに用いられる決定ユニット604と、

データノード機器に処理命令を送信することに用いられる第3送信ユニット605であって、処理命令はコミット命令、又はアボート命令であり、データノード機器は処理命令を実行することに用いられる、第3送信ユニット605と、を含む。

30

【0227】

1つの可能な実現形態において、データ読み取り結果は第2論理ライフサイクルを携えており、第2論理ライフサイクルはデータノード機器によってデータ読み取り要求が携えているターゲットトランザクションの第1論理ライフサイクルに基づき決定され、第1論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成され、第2送信ユニット603は、第1論理ライフサイクルのタイムスタンプ下限と第2論理ライフサイクルのタイムスタンプ下限における最大値をターゲットトランザクションの第3論理ライフサイクルのタイムスタンプ下限として用いることと、第1論理ライフサイクルのタイムスタンプ上限と第2論理ライフサイクルのタイムスタンプ上限における最小値をターゲットトランザクションの第3論理ライフサイクルのタイムスタンプ上限として用いることと、第3論理ライフサイクルが有効であることに応答して、データノード機器に第3論理ライフサイクルを携えているトランザクション認証要求を送信することであって、第3論理ライフサイクルが有効であることは、第3論理ライフサイクルのタイムスタンプ下限が第3論理ライフサイクルのタイムスタンプ上限よりも小さいことを指示することに用いられる、ことと、に用いられる。

40

【0228】

1つの可能な実現形態において、データノード機器の数は少なくとも2つであり、決定ユニット604は、少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果において認証に合格しなかったことを指示することに用いられる認証結果が存在することに応答して、アボート命令をターゲットトランザクションの処理命令として用い

50

ることと、少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果はいずれも認証に合格したことを指示することに応答して、少なくとも2つの認証結果を携えている論理ライフサイクルの共通部分をターゲット論理ライフサイクルとして用いることと、ターゲット論理ライフサイクルが有効であることに応答して、コミット命令をターゲットトランザクションの処理命令として用いることと、ターゲット論理ライフサイクルが無効であることに応答して、アポート命令をターゲットトランザクションの処理命令として用いることと、に用いられる。

【0229】

本願の実施例において、各ノード機器のそれぞれに対応するトランザクション割り当て指標に基づきターゲットトランザクションを協調処理することに用いられる協調ノード機器を決定し、トランザクションの割り当て過程はトランザクションに関するデータ項目を考慮する必要がなく、データ項目の分布状況を考慮する必要もない。このような方式に基づいて、個々のノード機器はいずれも分散化された機器としてトランザクションを協調処理することを可能にすることにより、トランザクションをノード間で処理することを可能にし、トランザクションの処理効率を向上させることに有利であり、トランザクション処理の信頼性が比較的高く、データベースシステムのシステム性能を高めることに有利である。

【0230】

図7に参照されるように、本願の実施例はトランザクション処理装置を提供し、該装置は、

協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得することに用いられる第1取得ユニット701であって、協調ノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、協調ノード機器は少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定される、第1取得ユニット701と、

データ読み取り結果を協調ノード機器に返信することに用いられる返信ユニット702と、

協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、ターゲットトランザクションの認証結果を取得することに用いられる第2取得ユニット703と、

協調ノード機器から送信されたターゲットトランザクションの処理命令を受信したことに応答して、処理命令を実行することに用いられる実行ユニット704であって、処理命令はコミット命令、又はアポート命令である、実行ユニット704と、を含み、

返信ユニット702はさらに、ターゲットトランザクションの認証結果を協調ノード機器に返信することに用いられる。

【0231】

1つの可能な実現形態において、データ読み取り要求はターゲットトランザクションの第1論理ライフサイクルを携えており、第1論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成され、第1取得ユニット701は、第1論理ライフサイクルに基づいて、データ読み取り要求によって指示された読み取り対象のデータ項目の可視バージョンデータを決定することと、可視バージョンデータの作成タイムスタンプ、及び第1論理ライフサイクルに基づいて、ターゲットトランザクションの第2論理ライフサイクルを決定することと、第2論理ライフサイクル、及び可視バージョンデータを携えている結果をデータ読み取り結果として用いることと、に用いられる。

【0232】

1つの可能な実現形態において、トランザクション認証要求はターゲットトランザクションの第3論理ライフサイクルを携えており、第3論理ライフサイクルは協調ノード機器によって第1論理ライフサイクル、及び第2論理ライフサイクルに基づいて決定された有効論理ライフサイクルであり、第2取得ユニット703は、第3論理ライフサイクルのタ

10

20

30

40

50

タイムスタンプ下限と第2論理ライフサイクルのタイムスタンプ下限における最大値をターゲットトランザクションの第4論理ライフサイクルのタイムスタンプ下限として用いること、第3論理ライフサイクルのタイムスタンプ上限と第2論理ライフサイクルのタイムスタンプ上限における最小値をターゲットトランザクションの第4論理ライフサイクルのタイムスタンプ上限として用いることと、第4論理ライフサイクルが有効であることに応答して、ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定することと、第5論理ライフサイクルが有効であることに応答して、認証に合格したことを指示することに用いられる認証結果をターゲットトランザクションの認証結果として用いることと、第5論理ライフサイクルが無効であることに応答して、認証に合格しなかったことを指示することに用いられる認証結果をターゲットトランザクションの認証結果として用いることと、に用いられる。

10

【0233】

1つの可能な実現形態において、1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は該1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプを含み、該1つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプは、該1つの書き込み対象のデータ項目を読み取った各読み取りトランザクションの論理コミットタイムスタンプにおける最大値を指示することに用いられ、第2取得ユニット703はさらに、各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定することに用いられ、第5論理ライフサイクルのタイムスタンプ下限は各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプにおける最大値よりも大きい。

20

【0234】

1つの可能な実現形態において、1つの書き込み対象のデータ項目の読み取りトランザクション関連情報は該1つの書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプを含み、ターゲット読み取りトランザクションはローカル認証に合格するか、又はコミット段階にある読み取りトランザクションであり、ターゲット読み取りトランザクションの終了タイムスタンプはターゲット読み取りトランザクションの論理ライフサイクルのタイムスタンプ上限であり、第2取得ユニット703はさらに、各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプ、及び第4論理ライフサイクルに基づいて、ターゲットトランザクションの第5論理ライフサイクルを決定することに用いられ、第5論理ライフサイクルのタイムスタンプ下限は各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプにおける最大値よりも大きい。

30

【0235】

本願の実施例において、各ノード機器のそれぞれに対応するトランザクション割り当て指標に基づきターゲットトランザクションを協調処理することに用いられる協調ノード機器を決定し、トランザクションの割り当て過程はトランザクションに関するデータ項目を考慮する必要がなく、データ項目の分布状況を考慮する必要もない。このような方式に基づいて、個々のノード機器はいずれも分散化された機器としてトランザクションを協調処理することを可能にすることにより、トランザクションをノード間で処理することを可能にし、トランザクションの処理効率を向上させることに有利であり、トランザクション処理の信頼性が比較的高く、データベースシステムのシステム性能を高めることに有利である。

40

【0236】

説明する必要があるように、上記実施例が提供する装置はその機能を実現する場合、上記各機能ユニットの分割のみを例として説明したが、実際の応用において、必要に基づき上記機能を異なる機能ユニットに割り当てて完了させることができ、すなわち機器の内部構造を異なる機能ユニットに分割して、上記記述された全部又は一部の機能を完了させる

50

。また、上記実施例が提供する装置は方法実施例と同一の構想に属し、その具体的な実現過程について、詳しくは方法実施例を参照すればよく、ここでは詳細な説明を省略する。

【0237】

図8は本願の実施例が提供するコンピュータ機器の構造模式図であり、該コンピュータ機器は配置又は性能が異なるため比較的大きく異なってもよく、1つ又は複数のプロセッサ(Central Processing Units、CPU)801と、1つ又は複数のメモリ802とを含んでもよく、ここで、該1つ又は複数のメモリ802において少なくとも1つのコンピュータプログラムが記憶され、該少なくとも1つのコンピュータプログラムは該1つ又は複数のプロセッサ801によってロードされて実行され、コンピュータ機器に上記各方法実施例が提供するトランザクション処理方法を実現させる。もちろん、該コンピュータ機器はさらに、入出力を行うように、有線又は無線ネットワークインタフェース、キーボード及び入出力インタフェースなどの部材を有してもよく、該コンピュータ機器はさらに機器の機能を実現することに用いられる他の部材を含んでもよく、ここでは詳細な説明を省略する。

10

【0238】

例示的な実施例において、非一時的コンピュータ可読記憶媒体をさらに提供し、該非一時的コンピュータ可読記憶媒体において少なくとも1つのコンピュータプログラムが記憶され、該少なくとも1つのコンピュータプログラムはコンピュータ機器のプロセッサによってロードされて実行され、コンピュータに上記いずれか1つのトランザクション処理方法を実現させる。

20

【0239】

1つの可能な実現形態において、上記非一時的コンピュータ可読記憶媒体は読み取り専用メモリ(Read-Only Memory、ROM)、ランダムアクセスメモリ(Random Access Memory、RAM)、読み取り専用光ディスク(Compact Disc Read-Only Memory、CD-ROM)、磁気テープ、フロッピーディスク及び光データ記憶機器などであってもよい。

【0240】

例示的な実施例において、コンピュータプログラム製品、又はコンピュータプログラムをさらに提供し、該コンピュータプログラム製品、又はコンピュータプログラムはコンピュータ命令を含み、該コンピュータ命令はコンピュータ可読記憶媒体に記憶される。コンピュータ機器のプロセッサはコンピュータ可読記憶媒体から該コンピュータ命令を読み取り、プロセッサは該コンピュータ命令を実行することにより、該コンピュータ機器が上記いずれか1つのトランザクション処理方法を実行する。

30

【0241】

理解すべきであるように、本願における用語「少なくとも1つ」とは、1つ又は複数を指し、「複数」又は「少なくとも2つ」の意味はいずれも2つ又は2つ以上を指す。「及び/又は」は、関連対象の関連関係を記述するためのものであり、3種の関係が存在することを表すことができ、例えば、A及び/又はBは、Aが単独で存在する、AとBが同時に存在する、Bが単独で存在するという3種の状況を表すことができる。文字「/」は、一般的に前後の関連オブジェクトが「又は」の関係であることを表す。

40

【0242】

以上は本願の例示的な実施例に過ぎず、本願を制限することに用いられるものではなく、本願の精神及び原則内で行われた任意の修正、等価置換、改良などは、いずれも本願の保護範囲内に含まれるべきである。

【符号の説明】

【0243】

- 101 ゲートウェイサーバ
- 102 機器
- 103 分散型記憶クラスタ
- 104 グローバルタイムスタンプ生成クラスタ

50

- 5 0 1 第 1 決定ユニット
- 5 0 2 第 2 決定ユニット
- 6 0 1 取得ユニット
- 6 0 2 第 1 送信ユニット
- 6 0 3 第 2 送信ユニット
- 6 0 4 決定ユニット
- 6 0 5 第 3 送信ユニット
- 7 0 1 第 1 取得ユニット
- 7 0 2 返信ユニット
- 7 0 3 第 2 取得ユニット
- 7 0 4 実行ユニット
- 8 0 1 プロセッサ
- 8 0 2 メモリ

10

20

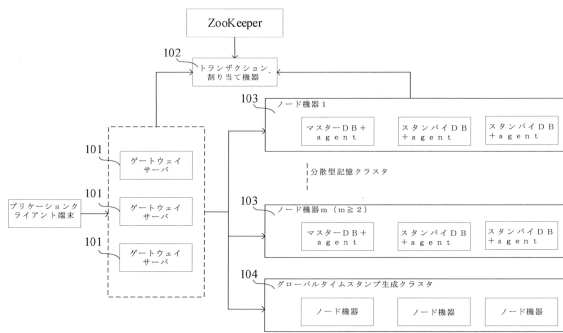
30

40

50

【 図 面 】

【 図 1 】

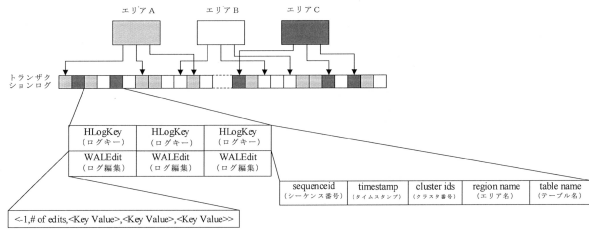


【 図 2 】

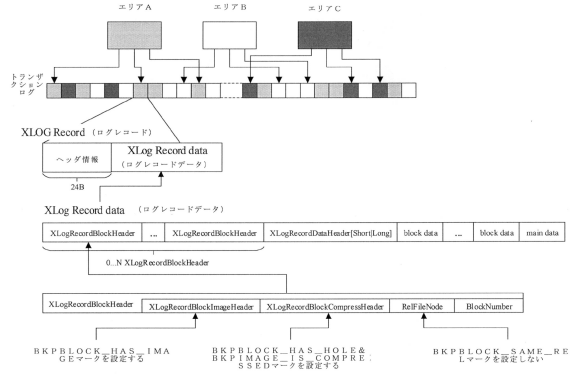


50

【 図 3 】

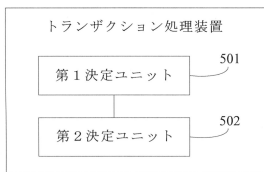


【 図 4 】

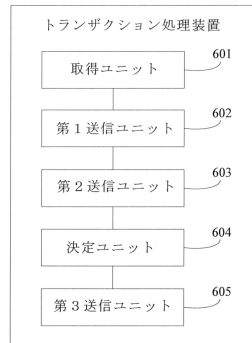


10

【 図 5 】

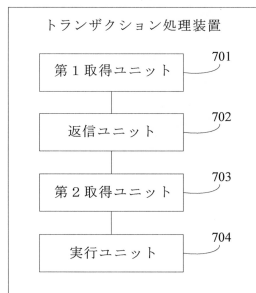


【 図 6 】

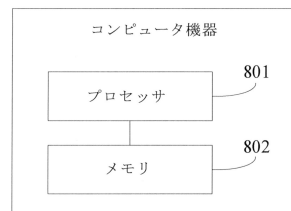


20

【 図 7 】



【 図 8 】



30

40

50

【手続補正書】

【提出日】令和5年3月15日(2023.3.15)

【手続補正2】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

トランザクション割り当て機器が実行するトランザクション処理方法であって、前記トランザクション割り当て機器は分散型データベースシステム中にあり、前記分散型データベースシステムにおいて同一の記憶システムを共有する少なくとも2つのノード機器がさらに含まれ、前記方法は、

ターゲットトランザクションの割り当て要求に応答して、前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定するステップであって、1つのノード機器に対応するトランザクション割り当て指標は前記1つのノード機器に新たなトランザクションを割り当てるマッチング度を指示することに用いられる、ステップと、

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの協調ノード機器を決定し、前記協調ノード機器によって前記ターゲットトランザクションを協調処理するステップと、を含む、トランザクション処理方法。

【請求項2】

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定する前記ステップは、

トランザクション割り当てモードを決定するステップであって、前記トランザクション割り当てモードはトランザクションのビジー度に基づく割り当て、機器のビジー度に基づく割り当て、及びハイブリッドのビジー度に基づく割り当てのうちのいずれか1つを含む、ステップと、

前記トランザクション割り当てモードによって指示された決定方式に基づき、前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定するステップと、を含む、請求項1に記載の方法。

【請求項3】

前記トランザクション割り当てモードはハイブリッドのビジー度に基づく割り当てを含み、前記トランザクション割り当てモードによって指示された決定方式に基づき、前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定する前記ステップは、

第1ノード機器のトランザクション処理数、前記第1ノード機器の機器リソース使用率、トランザクション処理数の重み、機器リソース使用率の重み、及び重み調節パラメータに基づいて、前記第1ノード機器に対応するトランザクション割り当て指標を決定するステップであって、前記第1ノード機器は前記少なくとも2つのノード機器のうちのいずれか1つのノード機器である、ステップを含む、請求項2に記載の方法。

【請求項4】

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの協調ノード機器を決定し、前記協調ノード機器によって前記ターゲットトランザクションを協調処理する前記ステップの後に、前記方法は、

前記協調ノード機器の機器識別情報を前記割り当て要求を発した端末に送信するステップであって、前記端末は前記協調ノード機器の機器識別情報に基づき、前記ターゲットトランザクションのトランザクション情報を前記協調ノード機器に送信し、前記協調ノード

10

20

30

40

50

機器によって前記トランザクション情報に基づいて前記ターゲットトランザクションを協調処理することに用いられる、ステップをさらに含む、請求項 1 ~ 3 のいずれか一項に記載の方法。

【請求項 5】

前記分散型データベースシステムはキー値型のデータ記憶フォーマットとセグメントページ型のデータ記憶フォーマットをサポートする、請求項 1 ~ 3 のいずれか一項に記載の方法。

【請求項 6】

協調ノード機器が実行するトランザクション処理方法であって、前記協調ノード機器は同一の記憶システムを共有する少なくとも 2 つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも 2 つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定され、前記方法は、

前記ターゲットトランザクションのトランザクション情報を取得するステップと、

前記ターゲットトランザクションのトランザクション情報に基づいて、データノード機器にデータ読み取り要求を送信するステップであって、前記データノード機器は前記少なくとも 2 つのノード機器のうち前記ターゲットトランザクションの処理に関与することに用いられるノード機器である、ステップと、

前記データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信するステップと、

前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定し、前記データノード機器に前記処理命令を送信するステップであって、前記処理命令はコミット命令、又はアポート命令であり、前記データノード機器は前記処理命令を実行することに用いられる、ステップと、を含む、トランザクション処理方法。

【請求項 7】

前記データ読み取り結果は第 2 論理ライフサイクルを携えており、前記第 2 論理ライフサイクルは前記データノード機器によって前記データ読み取り要求が携えている前記ターゲットトランザクションの第 1 論理ライフサイクルに基づき決定され、前記第 1 論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成され、

前記データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信する前記ステップは、

前記第 1 論理ライフサイクルのタイムスタンプ下限と前記第 2 論理ライフサイクルのタイムスタンプ下限における最大値を前記ターゲットトランザクションの第 3 論理ライフサイクルのタイムスタンプ下限として用い、前記第 1 論理ライフサイクルのタイムスタンプ上限と前記第 2 論理ライフサイクルのタイムスタンプ上限における最小値を前記ターゲットトランザクションの第 3 論理ライフサイクルのタイムスタンプ上限として用いるステップと、

前記第 3 論理ライフサイクルが有効であることに応答して、前記データノード機器に前記第 3 論理ライフサイクルを携えているトランザクション認証要求を送信するステップであって、前記第 3 論理ライフサイクルが有効であることは、前記第 3 論理ライフサイクルのタイムスタンプ下限が前記第 3 論理ライフサイクルのタイムスタンプ上限よりも小さいことを指示することに用いられる、ステップと、を含む、請求項 6 に記載の方法。

【請求項 8】

前記データノード機器の数は少なくとも 2 つであり、前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定し、前記データノード機器に前記処理命令を送信する前記ステップは、

10

20

30

40

50

前記少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果において認証に合格しなかったことを指示することに用いられる認証結果が存在することに対応して、前記アポート命令を前記ターゲットトランザクションの処理命令として用いるステップと、

前記少なくとも2つのデータノード機器から返信された少なくとも2つの認証結果はいずれも認証に合格したことを指示することに対応して、前記少なくとも2つの認証結果が携えている論理ライフサイクルの共通部分をターゲット論理ライフサイクルとして用いるステップと、

前記ターゲット論理ライフサイクルが有効であることに応答して、前記コミット命令を前記ターゲットトランザクションの処理命令として用い、前記ターゲット論理ライフサイクルが無効であることに応答して、前記アポート命令を前記ターゲットトランザクションの処理命令として用いるステップと、を含む、請求項6又は7のいずれか一項に記載の方法。

10

【請求項9】

データノード機器が実行するトランザクション処理方法であって、前記データノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションの処理に関与することに用いられるノード機器であり、前記方法は、

協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得し、前記データ読み取り結果を前記協調ノード機器に返信するステップであって、前記協調ノード機器は前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定される、ステップと、

20

前記協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得し、前記ターゲットトランザクションの認証結果を前記協調ノード機器に返信するステップと、

前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受信したことに応答して、前記処理命令を実行するステップであって、前記処理命令はコミット命令、又はアポート命令である、ステップと、を含む、トランザクション処理方法。

【請求項10】

前記データ読み取り要求は前記ターゲットトランザクションの第1論理ライフサイクルを携えており、前記第1論理ライフサイクルはタイムスタンプ下限、及びタイムスタンプ上限で構成され、

30

協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得し、前記データ読み取り結果を前記協調ノード機器に返信する前記ステップは、

前記第1論理ライフサイクルに基づいて、前記データ読み取り要求によって指示された読み取り対象のデータ項目の可視バージョンデータを決定するステップと、

前記可視バージョンデータの作成タイムスタンプ、及び前記第1論理ライフサイクルに基づいて、前記ターゲットトランザクションの第2論理ライフサイクルを決定するステップと、

前記第2論理ライフサイクル、及び前記可視バージョンデータが運ばれた結果を前記データ読み取り結果として用いるステップと、を含む、請求項9に記載の方法。

40

【請求項11】

前記トランザクション認証要求は前記ターゲットトランザクションの第3論理ライフサイクルを携えており、前記第3論理ライフサイクルは前記協調ノード機器によって前記第1論理ライフサイクル、及び前記第2論理ライフサイクルに基づいて決定された有効論理ライフサイクルであり、

前記協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得する前記ステップは、

前記第3論理ライフサイクルのタイムスタンプ下限と前記第2論理ライフサイクルのタイムスタンプ下限における最大値を前記ターゲットトランザクションの第4論理ライフサ

50

イクルのタイムスタンプ下限として用い、前記第 3 論理ライフサイクルのタイムスタンプ上限と前記第 2 論理ライフサイクルのタイムスタンプ上限における最小値を前記ターゲットトランザクションの第 4 論理ライフサイクルのタイムスタンプ上限として用いるステップと、

前記第 4 論理ライフサイクルが有効であることに応答して、前記ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び前記第 4 論理ライフサイクルに基づいて、前記ターゲットトランザクションの第 5 論理ライフサイクルを決定するステップと、

前記第 5 論理ライフサイクルが有効であることに応答して、認証に合格したことを指示することに用いられる認証結果を前記ターゲットトランザクションの認証結果として用い、前記第 5 論理ライフサイクルが無効であることに応答して、認証に合格しなかったことを指示することに用いられる認証結果を前記ターゲットトランザクションの認証結果として用いるステップと、を含む、請求項 10 に記載の方法。

【請求項 12】

1 つの書き込み対象のデータ項目の読み取りトランザクション関連情報は前記 1 つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプを含み、前記 1 つの書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプは前記 1 つの書き込み対象のデータ項目を読み取った各読み取りトランザクションの論理コミットタイムスタンプにおける最大値を指示することに用いられ、

前記第 4 論理ライフサイクルが有効であることに応答して、前記ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び前記第 4 論理ライフサイクルに基づいて、前記ターゲットトランザクションの第 5 論理ライフサイクルを決定する前記ステップは、

前記各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプ、及び前記第 4 論理ライフサイクルに基づいて、前記ターゲットトランザクションの第 5 論理ライフサイクルを決定するステップであって、前記第 5 論理ライフサイクルのタイムスタンプ下限は前記各書き込み対象のデータ項目の最大読み取りトランザクションタイムスタンプにおける最大値よりも大きい、ステップを含む、請求項 11 に記載の方法。

【請求項 13】

1 つの書き込み対象のデータ項目の読み取りトランザクション関連情報は前記 1 つの書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプを含み、前記ターゲット読み取りトランザクションはローカル認証に合格するか、又はコミット段階にある読み取りトランザクションであり、前記ターゲット読み取りトランザクションの終了タイムスタンプは前記ターゲット読み取りトランザクションの論理ライフサイクルのタイムスタンプ上限であり、

前記第 4 論理ライフサイクルが有効であることに応答して、前記ローカル書き込みセットに対応する各書き込み対象のデータ項目の読み取りトランザクション関連情報及び前記第 4 論理ライフサイクルに基づいて、前記ターゲットトランザクションの第 5 論理ライフサイクルを決定する前記ステップは、

前記各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプ、及び前記第 4 論理ライフサイクルに基づいて、前記ターゲットトランザクションの第 5 論理ライフサイクルを決定するステップであって、前記第 5 論理ライフサイクルのタイムスタンプ下限は前記各書き込み対象のデータ項目のターゲット読み取りトランザクションの終了タイムスタンプにおける最大値よりも大きい、ステップを含む、請求項 11 に記載の方法。

【請求項 14】

トランザクション処理システムであって、前記トランザクション処理システムは協調ノード機器と、データノード機器とを含み、前記協調ノード機器は同一の記憶システムを共有する少なくとも 2 つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも 2 つのノード機

10

20

30

40

50

器のそれぞれに対応するトランザクション割り当て指標に基づき決定され、前記データノード機器は前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの処理に参与することに用いられるノード機器であり、

前記協調ノード機器は、前記ターゲットトランザクションのトランザクション情報を取得することと、前記ターゲットトランザクションのトランザクション情報に基づいて、前記データノード機器にデータ読み取り要求を送信することと、に用いられ、

前記データノード機器は、前記協調ノード機器から送信された前記データ読み取り要求に基づいて、データ読み取り結果を取得し、前記データ読み取り結果を前記協調ノード機器に返信することに用いられ、

前記協調ノード機器はさらに、前記データノード機器から返信された前記データ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求、及びローカル書き込みセットを送信することに用いられ、

前記データノード機器はさらに、前記協調ノード機器から送信された前記トランザクション認証要求、及び前記ローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得し、前記ターゲットトランザクションの認証結果を前記協調ノード機器に返信することに用いられ、

前記協調ノード機器はさらに、前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定し、前記データノード機器に前記処理命令を送信することに用いられ、前記処理命令はコミット命令、又はアボート命令であり、

前記データノード機器はさらに、前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受信したことに応答して、前記処理命令を実行することに用いられる、トランザクション処理システム。

【請求項15】

トランザクション処理装置であって、前記装置は、

ターゲットトランザクションの割り当て要求に応答して、同一の記憶システムを共有する少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標を決定することに用いられる第1決定ユニットであって、1つのノード機器に対応するトランザクション割り当て指標は前記1つのノード機器に新たなトランザクションを割り当てるマッチング度を指示することに用いられる、第1決定ユニットと、

前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づいて、前記少なくとも2つのノード機器のうち前記ターゲットトランザクションの協調ノード機器を決定し、前記協調ノード機器によって前記ターゲットトランザクションを協調処理することに用いられる第2決定ユニットと、を含む、トランザクション処理装置。

【請求項16】

トランザクション処理装置であって、前記装置は、

ターゲットトランザクションのトランザクション情報を取得することに用いられる取得ユニットと、

前記ターゲットトランザクションのトランザクション情報に基づいて、データノード機器にデータ読み取り要求を送信することに用いられる第1送信ユニットであって、前記データノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうち前記ターゲットトランザクションの処理に参与することに用いられるノード機器である、第1送信ユニットと、

前記データノード機器から返信されたデータ読み取り結果がトランザクション認証条件を満たすことに応答して、前記データノード機器にトランザクション認証要求及び、ローカル書き込みセットを送信することに用いられる第2送信ユニットと、

前記データノード機器から返信された前記ターゲットトランザクションの認証結果に基づいて、前記ターゲットトランザクションの処理命令を決定することに用いられる決定ユニットと、

10

20

30

40

50

前記データノード機器に前記処理命令を送信することに用いられる第3送信ユニットであって、前記処理命令はコミット命令、又はアポート命令であり、前記データノード機器は前記処理命令を実行することに用いられる、第3送信ユニットと、を含む、トランザクション処理装置。

【請求項17】

トランザクション処理装置であって、前記装置は、

協調ノード機器から送信されたデータ読み取り要求に基づいて、データ読み取り結果を取得することに用いられる第1取得ユニットであって、前記協調ノード機器は同一の記憶システムを共有する少なくとも2つのノード機器のうちターゲットトランザクションを協調処理することに用いられるノード機器であり、前記協調ノード機器は前記少なくとも2つのノード機器のそれぞれに対応するトランザクション割り当て指標に基づき決定される、第1取得ユニットと、

前記データ読み取り結果を前記協調ノード機器に返信することに用いられる返信ユニットと、

前記協調ノード機器から送信されたトランザクション認証要求、及びローカル書き込みセットに基づいて、前記ターゲットトランザクションの認証結果を取得することに用いられる第2取得ユニットと、

前記協調ノード機器から送信された前記ターゲットトランザクションの処理命令を受信したことに応答して、前記処理命令を実行することに用いられる実行ユニットであって、前記処理命令はコミット命令、又はアポート命令である、実行ユニットと、

前記返信ユニットはさらに、前記ターゲットトランザクションの認証結果を前記協調ノード機器に返信することに用いられる、トランザクション処理装置。

【請求項18】

コンピュータ機器であって、前記コンピュータ機器はプロセッサと、メモリとを含み、前記メモリにおいて少なくとも1つのコンピュータプログラムが記憶され、前記少なくとも1つのコンピュータプログラムは前記プロセッサによってロードされて実行され、前記コンピュータ機器に、請求項1～5のいずれか一項に記載のトランザクション処理方法、又は請求項6～8のいずれか一項に記載のトランザクション処理方法、又は請求項9～13のいずれか一項に記載のトランザクション処理方法を実現させる、コンピュータ機器。

【請求項19】

コンピュータプログラムであって、プロセッサが前記コンピュータプログラムを実行することにより、コンピュータ機器が請求項1～5のいずれか一項に記載のトランザクション処理方法、請求項6～8のいずれか一項に記載のトランザクション処理方法、又は請求項9～13のいずれか一項に記載のトランザクション処理方法を実行するように構成された、コンピュータプログラム。

10

20

30

40

50

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No.

PCT/CN2021/126408

A. CLASSIFICATION OF SUBJECT MATTER G06F 9/48(2006.01)i According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) G06F		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) CNPAT, CNKI, WPI, EPODOC: 事务, 分配, 分发, 节点, 指标, 利用率, 使用率, 空闲, 验证, 本地写集, 回滚, transaction, assign, node, index, idle, verify, abort, utilization ratio		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
PX	CN 112162846 A (TENCENT TECHNOLOGY SHENZHEN CO., LTD.) 01 January 2021 (2021-01-01) claims 1-15	1-20
Y	CN 111597015 A (TENCENT TECHNOLOGY SHENZHEN CO., LTD. et al.) 28 August 2020 (2020-08-28) description paragraphs 61-98, 140-144, 229-390	1-20
Y	CN 110287022 A (BEIJING DA MI TECHNOLOGY CO., LTD.) 27 September 2019 (2019-09-27) description paragraphs 2, 48-119	1-20
A	CN 108958942 A (ZHENGZHOU YUNHAI INFORMATION TECHNOLOGY CO., LTD.) 07 December 2018 (2018-12-07) entire document	1-20
A	US 2011153566 A1 (MICROSOFT CORPORATION) 23 June 2011 (2011-06-23) entire document	1-20
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: “A” document defining the general state of the art which is not considered to be of particular relevance “E” earlier application or patent but published on or after the international filing date “L” document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) “O” document referring to an oral disclosure, use, exhibition or other means “P” document published prior to the international filing date but later than the priority date claimed		“T” later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention “X” document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone “Y” document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art “&” document member of the same patent family
Date of the actual completion of the international search 07 January 2022		Date of mailing of the international search report 25 January 2022
Name and mailing address of the ISA/CN China National Intellectual Property Administration (ISA/CN) No. 6, Xitucheng Road, Jimenqiao, Haidian District, Beijing 100088, China Facsimile No. (86-10)62019451		Authorized officer Telephone No.

Form PCT/ISA/210 (second sheet) (January 2015)

10

20

30

40

50

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.
PCT/CN2021/126408

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)	Publication date (day/month/year)
CN	112162846	A	01 January 2021	None	
CN	111597015	A	28 August 2020	None	
CN	110287022	A	27 September 2019	None	
CN	108958942	A	07 December 2018	None	
US	2011153566	A1	23 June 2011	US 8396831 B2	12 March 2013

10

20

30

40

50

国际检索报告

国际申请号

PCT/CN2021/126408

A. 主题的分类		
G06F 9/48(2006.01)i		
按照国际专利分类(IPC)或者同时按照国家分类和IPC两种分类		
B. 检索领域		
检索的最低限度文献(标明分类系统和分类号)		
G06F		
包含在检索领域中的除最低限度文献以外的检索文献		
在国际检索时查阅的电子数据库(数据库的名称, 和使用的检索词(如使用))		
CNPAT, CNKI, WPI, EPDOC; 事务, 分配, 分发, 节点, 指标, 利用率, 使用率, 空闲, 验证, 本地写集, 回滚, transaction, assign, node, index, idle, verify, abort, utilization ratio		
C. 相关文件		
类型*	引用文件, 必要时, 指明相关段落	相关的权利要求
PX	CN 112162846 A (腾讯科技深圳有限公司) 2021年1月1日 (2021 - 01 - 01) 权利要求1-15	1-20
Y	CN 111597015 A (腾讯科技深圳有限公司 等) 2020年8月28日 (2020 - 08 - 28) 说明书第61-98、140-144、229-390段	1-20
Y	CN 110287022 A (北京大米科技有限公司) 2019年9月27日 (2019 - 09 - 27) 说明书第2、48-119段	1-20
A	CN 108958942 A (郑州云海信息技术有限公司) 2018年12月7日 (2018 - 12 - 07) 全文	1-20
A	US 2011153566 A1 (MICROSOFT CORPORATION) 2011年6月23日 (2011 - 06 - 23) 全文	1-20
<input type="checkbox"/> 其余文件在C栏的续页中列出。 <input checked="" type="checkbox"/> 见同族专利附件。		
* 引用文件的具体类型: “A” 认为不特别相关的表示了现有技术一般状态的文件 “E” 在国际申请日的当天或之后公布的在先申请或专利 “L” 可能对优先权要求构成怀疑的文件, 或为确定另一篇引用文件的公布日而引用的或者因其他特殊理由而引用的文件(如具体说明的) “O” 涉及口头公开、使用、展览或其他方式公开的文件 “P” 公布日先于国际申请日但迟于所要求的优先权日的文件		“T” 在申请日或优先权日之后公布, 与申请不相抵触, 但为了理解发明之理论或原理的在后文件 “X” 特别相关的文件, 单独考虑该文件, 认定要求保护的发明不是新颖的或不具有创造性 “Y” 特别相关的文件, 当该文件与另一篇或者多篇该类文件结合并且这种结合对于本领域技术人员为显而易见时, 要求保护的发明不具有创造性 “&” 同族专利的文件
国际检索实际完成的日期		国际检索报告邮寄日期
2022年1月7日		2022年1月25日
ISA/CN的名称和邮寄地址		受权官员
中国国家知识产权局(ISA/CN) 中国北京市海淀区蓟门桥西土城路6号 100088		巩瑜
传真号 (86-10)62019451		电话号码 86-(10)-53961382

PCT/ISA/210 表(第2页) (2015年1月)

10

20

30

40

50

国际检索报告
关于同族专利的信息

国际申请号

PCT/CN2021/126408

检索报告引用的专利文件			公布日 (年/月/日)	同族专利	公布日 (年/月/日)
CN	112162846	A	2021年1月1日	无	
CN	111597015	A	2020年8月28日	无	
CN	110287022	A	2019年9月27日	无	
CN	108958942	A	2018年12月7日	无	
US	2011153566	A1	2011年6月23日	US	8396831 B2 2013年3月12日

10

20

30

40

PCT/ISA/210 表(同族专利附件) (2015年1月)

50

フロントページの続き

MK,MT,NL,NO,PL,PT,RO,RS,SE,SI,SK,SM,TR),OA(BF,BJ,CF,CG,CI,CM,GA,GN,GQ,GW,KM,ML,MR,N
E,SN,TD,TG),AE,AG,AL,AM,AO,AT,AU,AZ,BA,BB,BG,BH,BN,BR,BW,BY,BZ,CA,CH,CL,CN,CO,CR,CU,
CZ,DE,DJ,DK,DM,DO,DZ,EC,EE,EG,ES,FI,GB,GD,GE,GH,GM,GT,HN,HR,HU,ID,IL,IN,IR,IS,IT,JO,JP,K
E,KG,KH,KN,KP,KR,KW,KZ,LA,LC,LK,LR,LS,LU,LY,MA,MD,ME,MG,MK,MN,MW,MX,MY,MZ,NA,N
G,NI,NO,NZ,OM,PA,PE,PG,PH,PL,PT,QA,RO,RS,RU,RW,SA,SC,SD,SE,SG,SK,SL,ST,SV,SY,TH,TJ,TM,
TN,TR,TT,TZ,UA,UG,US,UZ,VC,VN,WS,ZA,ZM,ZW

(特許庁注：以下のものは登録商標)

1 . U N I X