

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6763580号  
(P6763580)

(45) 発行日 令和2年9月30日(2020.9.30)

(24) 登録日 令和2年9月14日(2020.9.14)

(51) Int.Cl.		F I
<b>G06F 8/656</b>	<b>(2018.01)</b>	G06F 8/656
<b>G06F 16/11</b>	<b>(2019.01)</b>	G06F 16/11
<b>G06F 16/182</b>	<b>(2019.01)</b>	G06F 16/182

請求項の数 32 (全 39 頁)

<p>(21) 出願番号 特願2018-529541 (P2018-529541)</p> <p>(86) (22) 出願日 平成28年12月19日(2016.12.19)</p> <p>(65) 公表番号 特表2019-502202 (P2019-502202A)</p> <p>(43) 公表日 平成31年1月24日(2019.1.24)</p> <p>(86) 国際出願番号 PCT/CN2016/110722</p> <p>(87) 国際公開番号 W02017/114213</p> <p>(87) 国際公開日 平成29年7月6日(2017.7.6)</p> <p>審査請求日 令和1年9月27日(2019.9.27)</p> <p>(31) 優先権主張番号 201511034171.7</p> <p>(32) 優先日 平成27年12月31日(2015.12.31)</p> <p>(33) 優先権主張国・地域又は機関 中国 (CN)</p>	<p>(73) 特許権者 510330264 アリババ・グループ・ホールディング・リミテッド ALIBABA GROUP HOLDING LIMITED 英国領、ケイマン諸島、グランド・ケイマン、ジョージ・タウン、ワン・キャピタル・プレイス、フォース・フロア、ピー・オー・ボックス 847</p> <p>(74) 代理人 110000877 龍華国際特許業務法人</p>
--	---

最終頁に続く

(54) 【発明の名称】 分散記憶システムをアップグレードするための方法および装置

(57) 【特許請求の範囲】

【請求項1】

クライアントに適用可能な、分散記憶システムをアップグレードするための方法であって、

複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信するステップと、

前記複数のデータサーバの各々により返される応答を受信するステップと、

前記応答に基づいて、成功した書き込みの数が所定の数より多いかどうかを判定するステップと、

前記成功した書き込みの前記数が前記所定の数より多い場合に、成功した書き込みで応答するそれぞれのデータサーバに対して第1フィードバック情報を送信するステップと、

前記成功した書き込みの前記数が前記所定の数より多くない場合に、成功した書き込みで応答するそれぞれのデータサーバに対して第2フィードバック情報を送信するステップと

を備え、

前記第1フィードバック情報または前記第2フィードバック情報は、前記データサーバが自らの状態を判定するために使用され、前記状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、データサーバの前記状態は、アップグレード制御サーバが前記データサーバに、アップグレード動作を実行するようローリング方式で

10

20

通知するために使用される、方法。

【請求項 2】

前記成功した書き込みの前記数が前記所定の数より多くない場合に、前記成功した書き込みを有するそれぞれのデータサーバに対して第 2 フィードバック情報を前記送信するステップであって、前記第 2 フィードバック情報は、前記データサーバが自らの状態をアップグレード不可能と判定するために使用される、前記送信するステップは、

前記成功した書き込みの前記数が前記所定の数と等しい場合に、前記成功した書き込みを有するそれぞれのデータサーバに対して前記第 2 フィードバック情報を送信するステップと、

前記成功した書き込みの前記数が前記所定の数より少ない場合に、前記複数のデータサーバ以外の少なくとも 1 つのデータサーバに対して前記書き込み対象データに関する書き込み要求を送信するステップと、

前記少なくとも 1 つのデータサーバにより返される応答を受信し、前記応答に従って、成功した書き込みの現在の数が前記成功した書き込みの以前の数と組み合わせて前記所定の数と等しくなるかどうかを判定するステップと、

前記成功した書き込みの前記現在の数が前記所定の数と等しい場合に、前記成功した書き込みを有するそれぞれのデータサーバに対して前記第 2 フィードバック情報を送信するステップと

を含む、請求項 1 に記載の方法。

【請求項 3】

前記成功した書き込みを有するそれぞれのデータサーバに対して第 2 フィードバック情報を前記送信するステップの後に、

新たな書き込み対象データが現れた場合に、以前に前記成功した書き込みがあった前記データサーバを含む前記複数のデータサーバを対象として使用するステップと、

前記複数のデータサーバに対して前記同じ書き込み対象データに関する書き込み要求を送信するステップと

を更に備える、請求項 1 または 2 に記載の方法。

【請求項 4】

前記第 1 フィードバック情報および前記第 2 フィードバック情報は、クライアント識別子を含み、次いで、前記第 1 フィードバック情報または前記第 2 フィードバック情報は、前記データサーバが自らの状態を判定するために使用され、

前記第 2 フィードバック情報を受信して前記データサーバの前記状態をアップグレード不可能なものとしてマーキングした後、前記データサーバが前記第 2 フィードバック情報内の前記クライアント識別子をアップグレード不可能なリストへ書き込むために、前記第 2 フィードバック情報を使用するステップを備え、

前記データサーバが前記第 1 フィードバック情報を受信した後、前記第 1 フィードバック情報は、前記データサーバが、前記第 1 フィードバック情報内の前記クライアント識別子を前記アップグレード不可能なリストから削除すること、および、前記アップグレード不可能なリストが空であると判定した後に自らの状態をアップグレード可能なものとしてマーキングすることを行うために使用される、請求項 3 に記載の方法。

【請求項 5】

複数のデータサーバに対して前記同じ書き込み対象データに関する書き込み要求を前記送信するステップの前に、

データサーバへアクセスするときに前記データサーバにより送信される第 2 アップグレード通知を受信するステップと、

前記第 2 アップグレード通知に従ってアップグレード準備状態へと移行するステップと

を更に備える、請求項 1 に記載の方法。

【請求項 6】

データサーバに適用される、分散記憶システムをアップグレードするための方法であっ

10

20

30

40

50

て、

クライアントから送信される第1フィードバック情報または第2フィードバック情報を受信するステップであって、前記第1フィードバック情報または前記第2フィードバック情報は、前記クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される、受信するステップと、

データサーバがアップグレードされなかった場合に、前記第1フィードバック情報または前記第2フィードバック情報に従って前記データサーバの状態を判定するステップであって、前記状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、前記状態は、アップグレード制御サーバがローリング方式で、前記データサーバを選択し、前記データサーバに、アップグレード動作を実行するよう通知するために使用される、判定するステップとを備える方法。

【請求項7】

前記第1フィードバック情報および前記第2フィードバック情報は、クライアント識別子を含み、次いで、前記データサーバがアップグレードされていない状況において、前記第1フィードバック情報または前記第2フィードバック情報に従って前記データサーバの状態を前記判定するステップは、

前記第2フィードバック情報を受信すると、前記第2フィードバック情報内の前記クライアント識別子をアップグレード不可能なリストに書き込むこと、および、自らの状態をアップグレード不可能なものとしてマーキングすることを行うステップと、

前記第1フィードバック情報を受信すると前記第1フィードバック情報内の前記クライアント識別子を前記アップグレード不可能なリストから削除すること、および、前記アップグレード不可能なリストが空であると判定した後に自らの状態をアップグレード可能なものとしてマーキングすることを行うステップとを含む、請求項6に記載の方法。

【請求項8】

前記第1フィードバック情報内の前記クライアント識別子を前記アップグレード不可能なリストから削除すること、および、前記アップグレード不可能なリストが空であると判定した後に自らの状態をアップグレード可能なものとしてマーキングすることを前記行うステップは、

前記アップグレード不可能なリストが前記第1フィードバック情報内の前記クライアント識別子を含むかどうかを判定するステップと、

前記アップグレード不可能なリストが前記第1フィードバック情報内の前記クライアント識別子を含む場合に、前記クライアント識別子を前記アップグレード不可能なリストから削除するステップと、

前記アップグレード不可能なリストが空であるかどうかを判定するステップと、

前記アップグレード不可能なリストが空であれば、自らの状態をアップグレード可能なものとしてマーキングするステップと

を含む、請求項7に記載の方法。

【請求項9】

前記第2フィードバック情報を受信すると、前記第2フィードバック情報内の前記クライアント識別子をアップグレード不可能なリストに書き込むこと、および、自らの状態をアップグレード不可能なものとしてマーキングすることを前記行うステップの後に、

前記アップグレード不可能なリスト内のクライアント識別子について、対応するクライアントの第2フィードバックメッセージが所定の数の期間内に受信されなかったかどうかを判定すること、前記対応するクライアントの前記第2フィードバックメッセージが前記所定の数の期間内に受信されなかった場合に、前記クライアント識別子を前記アップグレード不可能なリストから削除することを行うステップを更に備える、請求項7に記載の方法。

10

20

30

40

50

**【請求項 10】**

クライアントにより送信される第1フィードバック情報または第2フィードバック情報を前記受信するステップの前に、

前記アップグレード制御サーバにより送信される第1アップグレード通知を受信するステップと、

前記第1アップグレード通知に従ってアップグレード準備状態へと移行すること、および、前記クライアントが前記アップグレード準備状態へと移行するよう、前記クライアントのアクセス要求を受信した後に第2アップグレード通知を前記クライアントへ送信することを行うステップと

を更に備える、請求項6に記載の方法。

10

**【請求項 11】**

アップグレード制御サーバに適用される、分散記憶システムをアップグレードするための方法であって、

複数のデータサーバのデータサーバごとに状態を取得するステップであって、前記状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、それぞれのデータサーバは一方の状態を有し、前記データサーバの前記状態は、第1フィードバック情報または第2フィードバック情報に従って判定され、前記第1フィードバック情報または前記第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される、取得するステップと、

20

アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で通知するステップであって、前記データサーバは、前記通知に従って前記アップグレード動作を実行する、通知するステップとを備える方法。

**【請求項 12】**

アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で前記通知するステップは、

毎回アップグレード可能な状態にある少なくとも1つのデータサーバを選択すること、および、前記アップグレード可能な状態にある前記少なくとも1つのデータサーバに、前記アップグレード動作を実行するよう通知することを行うステップと、

30

前記アップグレード可能な状態にある前記少なくとも1つのデータサーバが前記アップグレード動作を完全に終了しているかどうかをモニタリングするステップと、

前記アップグレード可能な状態にある前記少なくとも1つのデータサーバが前記アップグレード動作を完全に終了している場合に、次回に備えてアップグレード可能な状態にある少なくとも1つのデータサーバを選択すること、および、前記アップグレード可能な状態にある前記少なくとも1つのデータサーバに、前記アップグレード動作を実行するよう通知することを行うステップと

を含む、請求項11に記載の方法。

**【請求項 13】**

任意のデータサーバの前記アップグレード動作の結果がアップグレードの失敗であることがモニタリングされた場合に、前記データサーバをアップグレードブラックリストに追加すること、および、前記データサーバのアップグレードを中断することを行うステップを更に備える、請求項12に記載の方法。

40

**【請求項 14】**

前記データサーバは、少なくとも2つのラック上に設置され、アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で前記通知するステップは、

毎回アップグレード可能な状態にある前記データサーバを最も多く有するラックを選択すること、および、前記ラック内の前記データサーバに、アップグレード動作を実行するよう通知することを行うステップであって、前記ラック内のそれぞれのデータサーバは

50

、前記通知に従って自らの状態をチェックする、行うステップと、

データサーバがアップグレード可能な状態にある場合に、前記データサーバを再起動およびアップグレードするステップと、

データサーバがアップグレード不可能な状態またはアップグレード終了状態にある場合に、前記データサーバの再起動およびアップグレードを拒否するステップとを含む、請求項 1 1 から 1 3 の何れか一項に記載の方法。

【請求項 1 5】

毎回アップグレード可能な状態にある前記データサーバを最も多く有するラックを選択すること、および、前記ラック内の前記データサーバに、アップグレード動作を実行するよう通知することを前記行うステップの後に、

10

前記ラック内の全ての前記データサーバが前記アップグレード動作を終了しているかどうかをモニタリングするステップと、

前記ラック内の全ての前記データサーバが前記アップグレード動作を終了している場合に、次回に備えてアップグレード可能な状態にある前記データサーバを最も多く有するラックを選択すること、および、前記ラック内の前記データサーバに、アップグレード動作を実行するよう通知することを行うステップと

を更に備える、請求項 1 4 に記載の方法。

【請求項 1 6】

データサーバごとに前記状態を取得する前に、

それぞれのデータサーバがアップグレード準備状態へと移行するよう、かつ、前記クライアントのアクセス要求を受信した後にそれぞれのデータサーバが第 2 アップグレード通知を前記クライアントへ送信することで、前記クライアントが前記アップグレード準備状態へと移行するよう、第 1 アップグレード通知をそれぞれのデータサーバに送信するステップ

20

を更に備える、請求項 1 1 に記載の方法。

【請求項 1 7】

クライアントに適用される、分散記憶システムをアップグレードするための装置であって、

複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信する要求送信モジュールと、

30

前記複数のデータサーバの各々により返される応答を受信すること、および、前記応答に基づいて、成功した書き込みの数が所定の数より多いかどうかを判定することを行う判定モジュールと、

前記成功した書き込みの前記数が前記所定の数より多い場合に、前記成功した書き込みを有するそれぞれのデータサーバに対して第 1 フィードバック情報を送信する第 1 フィードバックモジュールと、

前記成功した書き込みの前記数が前記所定の数より多くない場合に、前記成功した書き込みを有するそれぞれのデータサーバに対して第 2 フィードバック情報を送信する第 2 フィードバックモジュールと

を備え、

40

前記第 1 フィードバック情報または前記第 2 フィードバック情報は、前記データサーバが自らの状態を判定するために使用され、前記状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、前記データサーバの前記状態は、アップグレード制御サーバが前記データサーバに、アップグレード動作を実行するようローリング方式で通知するために使用される、装置。

【請求項 1 8】

前記第 2 フィードバックモジュールは、

前記成功した書き込みの現在の数が前記所定の数と等しい場合に、前記成功した書き込みを有するそれぞれのデータサーバに対して前記第 2 フィードバック情報を送信する第 2 フィードバック情報送信サブモジュールと、

50

前記成功した書き込みの前記数が前記所定の数より少ない場合に、前記複数のデータサーバ以外の少なくとも1つのデータサーバに対して前記書き込み対象データに関する書き込み要求を送信する書き込み要求送信サブモジュールと、

前記少なくとも1つのデータサーバにより返される応答を受信すること、および、前記応答に従って、前記成功した書き込みの前記現在の数が前記成功した書き込みの以前の数と組み合わせて前記所定の数と等しくなるかどうかを判定すること、および、前記成功した書き込みの前記現在の数が前記所定の数と等しい場合に、前記第2フィードバック情報送信サブモジュールへと移行することを行う判定サブモジュールとを含む、請求項17に記載の装置。

【請求項19】

前記第2フィードバックモジュールに次いで、

以前に前記成功した書き込みがあった前記データサーバを含む前記複数のデータサーバを対象として使用すること、および、新たな書き込み対象データが現れた場合は前記要求送信モジュールへと移行することを行う、新たな書き込み対象データの送信モジュールを更に備える、請求項17または18に記載の装置。

【請求項20】

前記第1フィードバック情報および前記第2フィードバック情報は、クライアント識別子を含み、次いで、前記第1フィードバック情報または前記第2フィードバック情報は、前記データサーバが自らの状態を判定するために使用され、

前記第2フィードバック情報を受信した後に、前記データサーバが、前記第2フィードバック情報内の前記クライアント識別子をアップグレード不可能なリストに書き込むこと、および、前記データサーバの状態をアップグレード不可能なものとしてマーキングすることを行うために、前記第2フィードバック情報を使用することと、

前記データサーバが前記第1フィードバック情報を受信した後に前記データサーバが、前記第1フィードバック情報内の前記クライアント識別子を前記アップグレード不可能なリストから削除するために、前記第1フィードバック情報を使用すること、および、前記アップグレード不可能なリストが空であると判定した後に前記データサーバの状態をアップグレード可能なものとしてマーキングすることを行うことと

を備える、請求項19に記載の装置。

【請求項21】

複数のデータサーバに対して前記同じ書き込み対象データに関する書き込み要求を前記送信するステップの前に、

データサーバへアクセスするときに前記データサーバにより送信される第2アップグレード通知を受信する第2アップグレード通知受信モジュールと、

前記第2アップグレード通知に従ってアップグレード準備状態へと移行する第2アップグレード準備モジュールと

を更に備える、請求項17に記載の装置。

【請求項22】

データサーバに適用される、分散記憶システムをアップグレードするための装置であって、

クライアントから送信される第1フィードバック情報または第2フィードバック情報を受信するフィードバック情報受信モジュールであって、前記第1フィードバック情報または前記第2フィードバック情報は、前記クライアントが複数のデータサーバに対して同じ書き込みデータに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される、フィードバック情報受信モジュールと、

前記データサーバがアップグレードされなかった場合に、前記第1フィードバック情報または前記第2フィードバック情報に従って前記データサーバの状態を判定する状態判定モジュールであって、前記状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、前記状態は、アップグレード制御サーバがローリング方式で、前記データサーバを選択し、前記データサーバに、アップグレード動作を実行するよう通知す

10

20

30

40

50

るために使用される、状態判定モジュールとを備える装置。

【請求項 2 3】

前記状態判定モジュールは、

前記第 2 フィードバック情報内のクライアント識別子をアップグレード不可能なリストに書き込むこと、および、前記第 2 フィードバック情報を受信すると、前記データサーバの状態をアップグレード不可能なものとしてマーキングすることをを行うアップグレード不可能な状態判定サブモジュールと、

前記第 1 フィードバック情報を受信すると、前記第 1 フィードバック情報内の前記クライアント識別子を前記アップグレード不可能なリストから削除すること、および、前記アップグレード不可能なリストが空であると判定した後に、前記データサーバの状態をアップグレード可能なものとしてマーキングすることをを行うアップグレード可能な判定サブモジュールと

を含む、請求項 2 2 に記載の装置。

【請求項 2 4】

前記アップグレード可能な判定サブモジュールは、

前記アップグレード不可能なリストが前記第 1 フィードバック情報内の前記クライアント識別子を含むかどうかを判定すること、および、含む場合は第 1 削除サブモジュールへと移行することをを行うクライアント識別子判定サブモジュールであって、前記第 1 削除サブモジュールは、前記クライアント識別子を前記アップグレード不可能なリストから削除する、クライアント識別子判定サブモジュールと、

前記アップグレード不可能なリストが空であるかどうかを判定すること、および、空であればアップグレード可能なマーキングサブモジュールへと移行することをを行うアップグレード不可能なリスト判定サブモジュールであって、前記アップグレード可能なマーキングサブモジュールは、前記データサーバの状態をアップグレード可能なものとしてマーキングする、アップグレード不可能なリスト判定サブモジュールと

を含む、請求項 2 3 に記載の装置。

【請求項 2 5】

前記アップグレード不可能な状態判定サブモジュールに次いで、

前記アップグレード不可能なリスト内のクライアント識別子について、対応するクライアントの第 2 フィードバックメッセージが所定の数の期間内に受信されなかったかどうかを判定すること、および、前記対応するクライアントの前記第 2 フィードバックメッセージが前記所定の数の期間内に受信されなかった場合に、第 2 削除サブモジュールへと移行することをを行う時間判定サブモジュールであって、

前記第 2 削除サブモジュールは、前記クライアント識別子を前記アップグレード不可能なリストから削除する、時間判定サブモジュールを更に備える、請求項 2 3 に記載の装置。

【請求項 2 6】

前記フィードバック情報受信モジュールより前に、

前記アップグレード制御サーバにより送信される第 1 アップグレード通知を受信する第 1 アップグレード通知受信モジュールと、

前記第 1 アップグレード通知に従ってアップグレード準備状態へと移行すること、および、前記クライアントが前記アップグレード準備状態へと移行するよう、前記クライアントのアクセス要求を受信した後に第 2 アップグレード通知を前記クライアントへ送信することをを行う第 1 アップグレード準備モジュールと

を更に備える、請求項 2 2 に記載の装置。

【請求項 2 7】

アップグレード制御サーバに適用される、分散記憶システムをアップグレードするための装置であって、

複数のデータサーバのそれぞれのデータサーバの状態を取得する状態取得モジュール

10

20

30

40

50

であって、前記状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、それぞれのデータサーバは一方の状態を有し、前記データサーバの前記状態は、第1フィードバック情報または第2フィードバック情報に従って判定され、前記第1フィードバック情報または前記第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される、状態取得モジュールと、

アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で通知するアップグレード通知モジュールであって、前記データサーバは、前記通知に従って前記アップグレード動作を実行する、アップグレード通知モジュールと

10

を備える装置。

【請求項28】

前記アップグレード通知モジュールは、

毎回アップグレード可能な状態にある少なくとも1つのデータサーバを選択すること、および、前記アップグレード可能な状態にある前記少なくとも1つのデータサーバに、前記アップグレード動作を実行するよう通知することを行う第1選択サブモジュールと、

前記アップグレード可能な状態にある前記少なくとも1つのデータサーバが前記アップグレード動作を完全に終了しているかどうかをモニタリングすること、および、終了している場合は前記第1選択サブモジュールへと移行することを行う第1モニタリングサブモジュールと

20

を含む、請求項27に記載の装置。

【請求項29】

任意のデータサーバの前記アップグレード動作の結果がアップグレードの失敗であることがモニタリングされた場合に、前記データサーバをアップグレードブラックリストに追加すること、および、前記データサーバのアップグレードを中断することを行う中断サブモジュール

を更に備える、請求項28に記載の装置。

【請求項30】

前記データサーバは、少なくとも2つのラック上に設置され、前記アップグレード通知モジュールは、

30

毎回アップグレード可能な状態にある前記データサーバを最も多く有するラックを選択すること、および、前記ラック内の前記データサーバに、アップグレード動作を実行するよう通知することを行うことであって、前記ラック内のそれぞれのデータサーバは、前記通知に従って自らの状態をチェックする、行うことと、データサーバがアップグレード可能な状態にある場合に、前記データサーバを再起動およびアップグレードすることと、データサーバがアップグレード不可能な状態またはアップグレード終了状態にある場合に、前記データサーバの再起動およびアップデートを拒否することとを行うアップグレード通知サブモジュール

を含む、請求項27から29の何れか一項に記載の装置。

【請求項31】

40

前記アップグレード通知サブモジュールに次いで、

前記ラック内の前記データサーバが前記アップグレード動作を完全に終了しているかどうかをモニタリングすること、および、終了している場合は前記アップグレード通知サブモジュールへと移行することを行う第2モニタリングサブモジュール

を更に備える、請求項30に記載の装置。

【請求項32】

前記状態取得モジュールの前に、

それぞれのデータサーバがアップグレード準備状態へと移行するように、かつ、前記クライアントのアクセス要求を受信した後にそれぞれのデータサーバが第2アップグレード通知を前記クライアントへ送信することで、前記クライアントが前記アップグレード準

50

備状態へと移行するよう、第1アップグレード通知をそれぞれのデータサーバに送信するアップグレード通知送信モジュールを更に備える、請求項27に記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

[関連出願の相互参照]

本開示は、2015年12月31日に出願された「METHOD AND APPARATUS FOR UPGRADING DISTRIBUTED STORAGE SYSTEM」と題する中国特許出願第201511034171.7号、および2016年12月19日に出願された「METHOD AND APPARATUS FOR UPGRADING DISTRIBUTED STORAGE SYSTEM」と題する国際特許出願第PCT/CN16/110722号に基づく優先権を主張するものである。これらの出願は両方とも、全体が参照により本明細書に組み込まれる。

【0002】

本開示は、分散コンピュータ技術の分野、特に分散記憶システムをアップグレードするための方法および装置に関する。

【背景技術】

【0003】

高可用性分散記憶システムが、高可用性サービスを構築するための基盤となっている。分散記憶システムは、分散データサーバで構成され、高信頼性データおよび高可用性アクセスを提供する。高信頼性は、データ冗長マルチバックアップコードまたは消去コードにより実装され、高可用性は、迅速な例外処理およびフェイルオーバーにより実装される。分散記憶システムのバージョンアップグレードを実行するには、当該システムのそれぞれのデータサーバが再起動されてバージョンアップデートを完了する必要がある。

【0004】

現行のシステムでは、以下の解決策を用いて分散記憶システムのアップグレードが実行され得る。

【0005】

第1解決策では、高水準サービスが無効となり、分散記憶システム全体が停止する。次いで、分散記憶システムの全てのデータサーバが再起動およびアップグレードされる。全てのデータサーバがアップグレードを完了した後、高水準サービスが再開する。しかしながら、この解決策では、全ての高水準サービスが利用不可能となり、こうした非可用性が長期間にわたって続くため、高可用性を必要とするサービスには適用可能でない。

【0006】

第2解決策では、高水準サービスが停止するのではなく、むしろそれぞれのデータサーバがローリング方式で再起動される。クライアントが複数のデータサーバ(データサーバの数は、バックアップを必要とするデータサーバのデフォルト数である)に対して書き込み要求を送信して、当該データサーバに書き込み対象データを書き込む。要求が失敗すると、クライアントは、当該データサーバの残りで再試行してクライアントのアクセスを確保する。しかしながら、この解決策では、クライアントがデータサーバのリカバリを待ってから再試行した場合に、この再試行手段がサーバの応答時間に大きな影響を及ぼすことになる。なぜなら、データサーバの再起動リカバリには通常何秒もの時間を必要とし、リアルタイム性の高いデータアクセスに10~100msの遅れが生じるからである。

【0007】

第3の解決策では、高水準サービスが停止するのではなく、むしろそれぞれのデータサーバがローリング方式で再起動される。クライアントが複数のデータサーバ(当該複数のデータサーバの数は、バックアップを必要とするデータサーバのデフォルト数である)に対して書き込み要求を送信して、当該データサーバに書き込み対象データを書き込む。要求が失敗すると、クライアントは、失敗したデータサーバを無視して再試行することでク

クライアントのアクセスを確保する。しかしながら、この解決策では、たとえダウンしているデータサーバをクライアントが一時的に無視したとしても、分散記憶システムがこれらのデータサーバのローリングアップグレードを継続する。分散記憶システム内のデータサーバがあまりにも急速にローリング方式で再起動された場合、クライアントが1つのデータサーバにだけ成功した書き込みを実行することで1つのデータサーバだけが成功した書き込みを有することになる場合、および、このデータサーバのディスクまたは機械がローリング再起動中に損傷を受けた場合は、クライアントの書き込み対象データが他のデータサーバに書き込まれていないのでユーザデータが失われる。たとえローリング再起動時間が延長されても、長時間のローリング中に分散記憶システムの大きなクラスタで予期せぬディスクおよび機械の例外が生じる。従って、1つのデータサーバだけがクライアントのデータの成功した書き込みを有するといった状況が存続し、ユーザデータの消失という大きな危険をもたらす。故に、この解決策では、待ち時間の問題は解消されるが、データの信頼性は低い。

10

**【発明の概要】****【0008】**

前述の問題を考慮して、分散記憶システムをアップグレードするための方法、および、それに対応する、分散記憶システムをアップグレードするための装置を提供することでこれらの問題を克服するか、または少なくとも部分的に解消するために、本開示の実施形態が紹介される。

**【0009】**

前述の問題を解消すべく、本開示は、クライアントに適用される、分散記憶システムをアップグレードするための方法を開示する。当該方法は、複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信するステップと、当該複数のデータサーバの各々により返される応答を受信するステップと、当該応答に基づいて、成功した書き込みの数が所定の数より多いかどうかを判定するステップと、成功した書き込みの数が所定の数より多い場合に、成功した書き込みで応答するそれぞれのデータサーバに対して第1フィードバック情報を送信するステップと、成功した書き込みの数が所定の数より多くない場合に、成功した書き込みで応答するそれぞれのデータサーバに対して第2フィードバック情報を送信するステップとを備える。ここで、第1フィードバック情報または第2フィードバック情報は、データサーバが自らの状態を判定するために使用され、当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、データサーバの状態は、アップグレード制御サーバがデータサーバに、アップグレード動作を実行するようローリング方式で通知するために使用される。

20

30

**【0010】**

本開示は更に、データサーバに適用される、分散記憶システムをアップグレードするための方法を開示する。当該方法は、クライアントから送信される第1フィードバック情報または第2フィードバック情報を受信するステップを備える。ここで、第1フィードバック情報または第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される。データサーバがアップグレードされなかった場合は、第1フィードバック情報または第2フィードバック情報に従って当該データサーバの状態を判定する。当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、当該状態は、アップグレード制御サーバがローリング方式で、データサーバを選択し、当該データサーバに、アップグレード動作を実行するよう通知するために使用される。

40

**【0011】**

本開示は更に、アップグレード制御サーバに適用される、分散記憶システムをアップグレードするための方法を開示する。当該方法は、複数のデータサーバのデータサーバごとに状態を取得するステップであって、当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、それぞれのデータサーバは一方の状態を有し、データ

50

サーバの状態は、第1フィードバック情報または第2フィードバック情報に従って判定され、第1フィードバック情報または第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される、取得するステップと、アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で通知するステップであって、この通知に従って当該データサーバはアップグレード動作を実行する、通知するステップとを備える。

【0012】

本開示は更に、クライアントに適用される、分散記憶システムをアップグレードするための装置を開示する。当該装置は、複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信するように構成される要求送信モジュールと、当該データサーバの各々により返される応答を受信すること、および、当該応答に基づいて、成功した書き込みの数が所定の数より多いかどうかを判定することを行うように構成される判定モジュールと、成功した書き込みの数が所定の数より多い場合に、成功した書き込みを有するそれぞれのデータサーバに対して第1フィードバック情報を送信するように構成される第1フィードバックモジュールと、成功した書き込みの数が所定の数より多くない場合に、成功した書き込みを有するそれぞれのデータサーバに対して第2フィードバック情報を送信するように構成される第2フィードバックモジュールとを備える。ここで、第1フィードバック情報または第2フィードバック情報は、データサーバが自らの状態を判定するために使用され、当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、データサーバの状態は、アップグレード制御サーバがデータサーバに、アップグレード動作を実行するようローリング方式で通知するために使用される。

【0013】

本開示は更に、データサーバに適用される、分散記憶システムをアップグレードするための装置を開示する。当該装置は、クライアントから送信される第1フィードバック情報または第2フィードバック情報を受信するように構成されるフィードバック情報受信モジュールであって、第1フィードバック情報または第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される、フィードバック情報受信モジュールと、データサーバがアップグレードされなかった場合に、第1フィードバック情報または第2フィードバック情報に従ってデータサーバの状態を判定するように構成される状態判定モジュールであって、当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、当該状態は、アップグレード制御サーバがローリング方式で、データサーバを選択し、当該データサーバに、アップグレード動作を実行するよう通知するために使用される、状態判定モジュールとを備える。

【0014】

本開示は更に、アップグレード制御サーバに適用される、分散記憶システムをアップグレードするための装置を開示する。当該装置は、複数のデータサーバのデータサーバごとに状態を取得するように構成される状態取得モジュールであって、当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、それぞれのデータサーバは一方の状態を有し、データサーバの状態は、第1フィードバック情報または第2フィードバック情報に従って判定され、第1フィードバック情報または第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される、状態取得モジュールと、アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で通知するように構成されるアップグレード通知モジュールであって、この通知に従って当該データサーバはアップグレード動作を実行する、アップグレード通知モジュールとを備える。本開示の実施形態には以下の利点がある。

【0015】

本開示の実施形態において、分散記憶システムがアップグレード手続きを開始した後、分散記憶システムへアクセスしているそれぞれのクライアントは、同じ書き込み対象データに関する書き込み要求を複数のデータサーバへ同時に送信する。それぞれのクライアントは次いで、書き込み対象データが幾つのデータサーバへ成功裏に書き込まれているかを解析し、成功した書き込みの数が所定の数より多いかどうかを判定し、その判定結果に従って、成功した書き込みを有するそれぞれのデータサーバに対して第1フィードバック情報または第2フィードバック情報を送信する。受信される第1フィードバック情報または第2フィードバック情報に従って、データサーバは、自らがアップグレード可能な状態にあるのか、アップグレード不可能な状態にあるのかを判定する。次いで、データサーバの状態に従って、アップグレード制御サーバはローリング方式で、データサーバを選択し、当該データサーバに、アップグレード動作を実行するよう通知してよい。前述のプロセスにより、本開示の実施形態において、アップグレード制御サーバが、ローリング方式でアップグレードするようデータサーバを制御する場合は、高水準サービスが停止する必要はない。なぜなら、クライアントがデータサーバの状態を制御し、任意のクライアントの書き込み対象データがバックアップのために少なくとも所定の数のデータサーバへ書き込まれることが保証されるからである。更には、分散記憶システムからクライアントへのより短い応答時間が保証され得ることで、データの信頼性が高まり、ユーザデータの損失という危険が著しく低下する。

10

**【図面の簡単な説明】**

【0016】

20

【図1】本開示の幾つかの実施形態に係る、クライアント側で分散記憶システムをアップグレードするための方法を示すフロー図である。

【0017】

【図2】本開示の幾つかの実施形態に係る、クライアント側で分散記憶システムをアップグレードするための方法を示すフロー図である。

【0018】

【図3】本開示の幾つかの実施形態に係る、アップグレード制御サーバ側で分散記憶システムをアップグレードするための方法を示すフロー図である。

【0019】

【図4】本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするための方法を示すフロー図である。

30

【0020】

【図5】本開示の幾つかの実施形態に係る、クライアント側で分散記憶システムをアップグレードするための装置のブロック図である。

【0021】

【図6】本開示の幾つかの実施形態に係る、クライアント側で分散記憶システムをアップグレードするための装置のブロック図である。

【0022】

【図7】本開示の幾つかの実施形態に係る、データサーバ側で分散記憶システムをアップグレードするための装置のブロック図である。

40

【0023】

【図8】本開示の幾つかの実施形態に係る、アップグレード制御サーバ側で分散記憶システムをアップグレードするための装置のブロック図である。

【0024】

【図8A】本開示の幾つかの実施形態に係る、分散記憶システムのアーキテクチャを示している。

**【発明を実施するための形態】**

【0025】

上述した本開示の目的、特徴および利点をより明瞭で理解し易くするために、本開示は、添付図および具体的な実装と併せて更に詳しく後述される。

50

## 【 0 0 2 6 】

開示される実施形態の核となる着想のうち1つは、分散記憶システムのデータサーバおよびアップグレード制御サーバに対する創造的変更が為されるというものであり、データサーバへアクセスしているクライアントに対して新たな実行論理が提供されるというものである。分散記憶システムがアップグレード手続きを開始した後、分散記憶システムへアクセスしているそれぞれのクライアントは、同じ書き込み対象データに関する書き込み要求を複数のデータサーバへ同時に送信する。次いで、書き込み対象データが成功裏に書き込まれたデータサーバの数が解析される。成功した書き込みの数が所定の数より多いかどうか判定される。この判定に従って、成功した書き込みを有するそれぞれのデータサーバに対して第1フィードバック情報または第2フィードバック情報が送信される。受信される第1フィードバック情報または第2フィードバック情報に従って、データサーバは、自らがアップグレード可能な状態にあるのか、アップグレード不可能な状態にあるのかを判定する。次いで、データサーバの状態に従って、アップグレード制御サーバはローリング方式で、データサーバを選択し、当該データサーバに、アップグレード動作を実行するよう通知してよい。前述のプロセスにより、アップグレード制御サーバが、ローリング方式でアップグレードするようアップグレード可能な状態にあるデータサーバを制御する場合は、高水準サービスが停止する必要はない。なぜなら、クライアントがデータサーバの状態を制御し、任意のクライアントの書き込み対象データがバックアップのために少なくとも所定の数のデータサーバへ書き込まれることが保証されるからである。更には、分散記憶システムからクライアントへのより短い応答時間が保証され得ることで、データの信頼性が高まり、ユーザデータの損失という危険が著しく低下する。更に、このプロセスは、サービスが影響を受けないように保証しつつ、アップグレード中の機械の予期せぬ例外を許容できる。

10

20

## 【 0 0 2 7 】

後述される本開示の実施形態は、クライアントに適用されてよい。

## 【 0 0 2 8 】

図1は、本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするための方法を示すフロー図である。具体的に言うと、当該方法は以下のステップを含んでよい。

## 【 0 0 2 9 】

ステップ110：複数のデータサーバに対して書き込み対象データに関する書き込み要求を送信する。一実施形態では、書き込み要求ごとに書き込み対象データが同じである。

30

## 【 0 0 3 0 】

本開示のこの実施形態において、クライアント(A)が分散記憶システムのアップグレード中にデータA1をデータサーバへ書き込もうとする場合は、クライアントAが複数のデータサーバに対してデータA1に関する書き込み要求を送信する。

## 【 0 0 3 1 】

本開示のこの実施形態では、複数のデータサーバの数Rが、例えば10に予め設定されてよい。本開示で具体的な数Rが限定されることはない。

## 【 0 0 3 2 】

次いで、クライアントAは、R個のデータサーバに対して書き込み対象データに関する書き込み要求を送信し、次いで、ステップ120へと進んでよい。

40

## 【 0 0 3 3 】

実際の適用においては、まず、クライアントAが分散記憶システムのスケジューリングサーバに対してR個の書き込み要求を送信してよいことが理解され得る。次いで、スケジューリングサーバは、どのR個のデータサーバに当該R個の書き込み要求が割り当てられるかを制御する。

## 【 0 0 3 4 】

異なるクライアントがR個のデータサーバに対して自らの書き込み対象データに関する書き込み要求を送信する場合は、R個のデータサーバがそれぞれ同じであってもよいし、

50

異なってもよいことに留意すべきである。

【0035】

本開示の別の実施形態において、当該方法は、ステップ110の前にステップ101および102を更に含む。

【0036】

ステップ101：データサーバへアクセスするときに当該データサーバにより送信される第2アップグレード通知を受信する。

【0037】

ステップ102：第2アップグレード通知に従って、アップグレード準備状態へと移行する。

10

【0038】

本開示のこの実施形態では、まず、分散記憶システムのアップグレード制御サーバが、それぞれのデータサーバにアップグレード準備状態への移行を通知する。次いで、それぞれのデータサーバは、クライアントがアップグレード準備状態へと移行するよう、データサーバへアクセスしているクライアントに対して第2アップグレード通知を送信してよい。

【0039】

それに応じて、アップグレード準備状態へと移行しているデータサーバにクライアントAがアクセスした後、当該データサーバは、第2アップグレード通知をクライアントに返す。第2アップグレード通知を受信した後、クライアントは、第2アップグレード通知に従ってアップグレード準備状態へと移行する。

20

【0040】

本開示のこの実施形態において、クライアントがブラウザのウェブページを介してデータサーバにアクセスする場合は、ブラウザにより開かれたウェブページを介して、アップグレードスクリプトがブラウザに送信されてよい。当該スクリプトを受信した後、ブラウザは、クライアントがアップグレード準備状態へと移行するよう当該スクリプトを実行してよい。クライアントがモバイルアプリケーション（例えばALIPAYアプリケーション）を介してデータサーバにアクセスする場合は、当該アプリケーションにアップグレード処理論理が予め追加されてよい。アップグレード処理論理は、クライアントがアップグレード準備状態へと移行するよう、アプリケーションが第2アップグレード通知を受信した後に有効となる。もちろん、本開示で具体的なやり方が限定されることはない。

30

【0041】

クライアントがアップグレード準備状態へと移行するとき、クライアントにより実行される、本開示の一実施形態におけるステップ110、120、130および140等のステップへとフローが移行してよい。

【0042】

ステップ120：データサーバの各々により返される書き込み応答を受信し、当該応答に従って、成功した書き込みの数が所定の数より多いかどうかを判定する。成功した書き込みの数が所定の数より多い場合は、ステップ130へと移行する。成功した書き込みの数が所定の数より多くない場合は、ステップ140へと移行する。

40

【0043】

本開示のこの実施形態において、正常な状況下で、それぞれのデータサーバは、クライアントにより送信される、同じ書き込み対象データに関する書き込み要求を受信した後、相応してクライアントへ応答を返す。もちろん、アップグレード中の特定のデータサーバまたはダウンしている特定のデータサーバがクライアントに応答を返すことはない。

【0044】

本開示のこの実施形態において、クライアントは、R個のデータサーバに対して書き込み対象データA1に関する書き込み要求を送信した後、受信される応答を定期的にチェックしてよい。例えば、特定の1つまたは複数のデータサーバの1つまたは複数の応答が受信されたかどうか定期的にチェックされる。指定された期間内に特定の1つまたは複数

50

のデータサーバの1つまたは複数の応答が受信されなかった場合、このことは、書き込み対象データA1が特定の1つまたは複数のデータサーバへ成功裏に書き込まれていないことを示す。指定された期間内に特定の1つまたは複数のデータサーバの1つまたは複数の応答が受信された場合は、当該応答が成功した書き込み応答であるのか、失敗した書き込み応答であるのかが判定されてよい。

**【0045】**

応答が成功した書き込み応答であれば、このことは、当該成功した書き込み応答に対応するデータサーバが書き込み対象データA1を成功裏にバックアップしたこと、すなわち成功裏に書き込んだことを示す。応答が失敗した書き込み応答であれば、当該失敗した書き込み応答に対応するデータサーバが書き込み対象データA1をバックアップするのに失敗したこと、すなわち書き込むのに失敗したことを示す。

10

**【0046】**

R個のデータサーバの応答については、MがNより大きいかどうかをクライアントAが定期的にチェックすることに留意すべきである(パラメータMおよびNについては、本明細書で更に説明される)。期間は例えば1ms(マイクロ秒)である。もちろん、期間は必要に従って決定されてよく、本開示により限定されることはない。次いで、前述の決定に基づいて、データサーバに対する書き込み対象データA1の成功した書き込みの数が算出されてよい。

**【0047】**

本開示のこの実施形態では、成功した書き込みの所定の数が予め設定されてよい。所定数は、成功した書き込みの必要最低限の数であり、クライアントの書き込み対象データのバックアップに成功したデータサーバの数と理解されてもよい。所定数は例えばNである。ここで、 $N < R$ であり、NおよびRは両方とも正の整数である。実際には、デフォルト設定により $N = 3$ および $R = 10$ であってよい。もちろん、NおよびRの値は実際の必要に従って設定されてよく、本開示で限定されることはない。

20

**【0048】**

次いで、クライアントAが現在のR個のデータサーバに対して書き込み対象データA1に関する書き込み要求を送信した後、M個のデータサーバが成功した書き込み応答を返す場合は、 $M > N$ かどうか判定される。 $M > N$ ならば、フローがステップ130へと移行する。 $M \leq N$ ならば、フローがステップ140へと移行する。 $M \geq 0$ であり、Mは整数である。

30

**【0049】**

実際の適用においては、データサーバがラック上に設置される。次いで、ディザスタリカバリ・バックアップの効果を更に高めるために、R個のデータサーバが異なるラック上に設置されてよい。R個のデータサーバは、全て異なるラック上に設置されるのが最適である。

**【0050】**

ステップ130: 成功した書き込みを有するそれぞれのデータサーバに対して第1フィードバック情報を送信する。

**【0051】**

クライアントの書き込み要求が成功裏に受信され得ること、および、内部の書き込み対象データがバックアップされ得ることから、M個のデータサーバが利用可能であると言える。次いで、クライアントは、M個のデータサーバに対して第1フィードバックメッセージを送信してよい。第1フィードバックメッセージは、例えばOKメッセージであってよく、OKメッセージは、データサーバに対してそれが自らの状態をアップグレード可能なものに設定できることを知らせる。

40

**【0052】**

実際のところ、本開示の図2では、データサーバが第1フィードバックメッセージに従って自らの状態をアップグレード可能なものに設定する方法が詳しく説明されている。

**【0053】**

50

ステップ140：成功した書き込みを有するそれぞれのデータサーバに対して第2フィードバック情報を送信する。

【0054】

第1フィードバック情報または第2フィードバック情報は、データサーバが自らの状態を判定するために使用される。当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれる。データサーバの状態は、アップグレード制御サーバがデータサーバに、アップグレード動作を実行するようローリング方式で通知するために使用される。

【0055】

M < Nならば、M個のデータサーバが利用可能となるので、クライアントは、M個のデータサーバに対して第2フィードバックメッセージを送信してよい。第2フィードバックメッセージは、例えばHOLDハートビートメッセージであってよく、HOLDメッセージは、データサーバに対してそれが自らの状態をアップグレード不可能なものに設定できることを知らせる。実際のところ、本開示の図2では、データサーバが第2フィードバックメッセージに従って自らの状態をアップグレード不可能なものに設定する方法が詳しく説明されている。

10

【0056】

M < Nという状況については、Mが少なくともNと等しいことが保証されるまで、クライアントが書き込み要求を新たなデータサーバへ再び送信することに留意すべきである。

【0057】

20

分散記憶システム内のアップグレードされなかったそれぞれのデータサーバについては、データサーバは、データサーバが第1フィードバック情報または第2フィードバック情報を受信するときに、第1フィードバック情報または第2フィードバック情報に従って自らがアップグレード可能な状態にあるのか、アップグレード不可能な状態にあるのかを判定してよい。アップグレード可能な状態は例えばOKであり、アップグレード不可能な状態は例えばHOLDである。例えば、特定のデータサーバが第2フィードバック情報を受信する場合は、データサーバが自らの状態をアップグレード不可能なものに設定してよい。特定のデータサーバが第1フィードバック情報を受信する場合は、データサーバが自らの状態をアップグレード可能なものに設定してよい。

【0058】

30

もちろん、実際の適用においては、高水準サービスが停止するのではなく、データサーバが継続的にアップグレードされるので、前述のR個のデータサーバが、アップグレードを完了するデータサーバを含んでもよい。次いで、アップグレードされたデータサーバについては、データサーバにより受信される任意の前述のフィードバック情報が更に処理されることはなく、データサーバはアップグレードされた状態のままである。

【0059】

本開示の図2では、受信される第1フィードバック情報または第2フィードバック情報に従ってデータサーバが自らの状態を判定する方法の具体的なプロセスが詳しく説明されている。

【0060】

40

更に、それぞれのデータサーバの状態は、アップグレード制御サーバがデータサーバに、アップグレード動作を実行するようローリング方式で通知するために使用されてよい。本開示の図3では、アップグレード制御サーバがデータサーバに、アップグレード動作を実行するようローリング方式で通知する方法の具体的なプロセスが説明されている。

【0061】

本開示の別の実施形態において、ステップ140は、サブステップ141～144を含む。

【0062】

サブステップ141：成功した書き込みの数が所定の数と等しい場合は、成功した書き込みを有するそれぞれのデータサーバに対して第2フィードバック情報を送信する。

50

## 【 0 0 6 3 】

例えば、前述のクライアント A は、R 個のデータサーバに対して書き込み対象データ A 1 を初めて送信する。次いで、モニタリングされた成功した書き込み応答の数が  $M = N$  であれば、成功した書き込みを有する M 個のデータサーバに対して第 2 フィードバック情報が直接送信され得る。

## 【 0 0 6 4 】

サブステップ 1 4 2 : 成功した書き込みの数が所定の数より少なければ、当該複数のデータサーバ以外の少なくとも 1 つのデータサーバに対して書き込み対象データに関する書き込み要求を送信する。

## 【 0 0 6 5 】

サブステップ 1 4 3 : 当該少なくとも 1 つのデータサーバにより返される応答を受信し、当該応答に従って、成功した書き込みの現在の数が成功した書き込みの以前の数と組み合わせて所定の数と等しくなるかどうかを判定する。成功した書き込みの現在の数が所定の数と等しい場合は、フローがサブステップ 1 4 4 へと移行する。成功した書き込みの現在の数が所定の数より少なければ、フローがサブステップ 1 4 2 へと移行する。

## 【 0 0 6 6 】

サブステップ 1 4 4 : 成功した書き込みの数が所定の数と等しい場合は、成功した書き込みを有するそれぞれのデータサーバに対して第 2 フィードバック情報を送信する。

## 【 0 0 6 7 】

クライアント A が R 個のデータサーバに対して書き込み対象データ A 1 を初めて送信するサブステップ 1 4 2 ~ 1 4 4 については、モニタリングされた成功した応答の数が  $M < N$  であれば、少なくとも  $(N - M)$  個のデータサーバに対して書き込み要求が改めて送信される。次いで  $(N - M)$  個のデータサーバの応答が受信され、 $(N - M)$  個のデータサーバにより返される成功した書き込み応答の数が、R 個のデータサーバの成功した書き込み応答の以前の数に追加されて M となる。この時点で  $M = N$  ならば、成功した書き込みを有するそれぞれのデータサーバに対して第 2 フィードバック情報が送信される。この時点で M の値は変化するが、M がこの時点で依然として N より小さければ、少なくとも  $(N - M)$  個のデータサーバに対して書き込み要求が再び送信される。このプロセスは  $M = N$  となるまで継続し、次いで、成功した書き込みを有するそれぞれのデータサーバに対して第 2 フィードバック情報が送信される。

## 【 0 0 6 8 】

更に、 $R = 10$  であり、 $N = 5$  であり、かつ、クライアント A が 10 個のデータサーバに対して書き込み対象データ A 1 を初めて送信する場合は、モニタリングされた成功した書き込み応答の数が  $M = 3$  であれば、少なくとも  $(5 - 3 = 2)$  個のデータサーバに対して書き込み要求が改めて送信される。2 つのデータサーバの成功した書き込み応答の数が更にモニタリングされ、一方の成功した書き込み応答だけが受信されたことが分かり、次いで、予め記録されている 3 に 1 が追加されて  $M = 4$  となる。次いで、少なくとも  $(5 - 4 = 1)$  個のデータサーバに対して書き込み要求が改めて送信される。前述のやり方で、モニタリングおよび判定が更に実行される。 $M = N = 5$  ならば、成功した書き込みを有するそれぞれのデータサーバに対して第 2 フィードバック情報が送信される。実際の適用において、 $M < N$  ならば、前述のプロセスにおいて第 2 フィードバックメッセージは、書き込み要求が送信されるたびに、成功した書き込みを有するデータサーバに対して直接送信されてよく、 $M = N$  になるまで待ってからまとめて送信しなくてもよい。

## 【 0 0 6 9 】

もちろん、実際の適用においては、データサーバの応答が定期的にチェックされるので、指定された期間 T、例えば  $3 * T$  において  $M < N$  であると更に判定された場合は、少なくとも  $(N - M)$  個のデータサーバに対して書き込み要求が改めて送信される。M および N に関するそのような判定は、周期ごとに為される。一般的には、複数の周期の後に  $M = N$  であると判定され得る。応答時間は一般的に短く、クライアントの応答遅延に大きな影響を及ぼすことはない。ユーザの観点からすれば、前述の状況における応答遅延は、再起

10

20

30

40

50

動を待っている際の応答遅延よりもはるかに短いことが分かる。

【0070】

このように、複数のデータベースにおけるクライアントのバックアップ対象データの成功したバックアップが、アップロードサービスの通常使用に影響を及ぼすことなく保証される。

【0071】

本開示の別の実施形態において、当該方法は、ステップ140の後にステップ150を更に含む。

【0072】

ステップ150：新たな書き込み対象データが現れた場合は、以前に成功した書き込みがあったデータサーバを含む複数のデータサーバが対象として使用され、ステップ120へと移行する。

10

【0073】

$M = N$ という状況については、成功した書き込みを有するそれぞれのデータサーバに対してクライアントがHOLDハートビート情報を返すので、これらのデータサーバが自らをアップグレード不可能な状態、例えばHOLD状態に設定することが理解され得る。これらのデータサーバにおけるクライアントの制限をより簡便に取り除くべく、本開示のこの実施形態において、クライアントはその後、以前のR個のデータサーバに対して新たな書き込み対象データに関する書き込み要求を送信する。

【0074】

20

例えば、クライアントAの書き込み対象データA1については、10個のデータサーバU1、U2・・・U10に対して書き込み要求が予め送信され、 $M = N$ と判定される。次いで、クライアントAの書き込み対象データA2については、10個のデータサーバU1、U2・・・U10に対して書き込み要求が更に送信されてよく、ステップ120での判定する段階が継続する。

【0075】

サブステップ142については、R個のデータサーバが、成功した書き込みを有するM個のデータサーバを含むこと、次いで、失敗した書き込みを有するデータサーバから残り( $R - M$ 個)のデータサーバが選択されることが理解され得る。例えば、 $R = 10$ かつ $N = 3$ であり、クライアントAの書き込み対象データA1については、10個のデータサーバU1、U2・・・U10に対して書き込み要求が初めて送信され、U1およびU2だけが成功した書き込みを有し、次いで、U11およびU12に対して2度目の書き込み要求が送信され、U11が成功した書き込みを有する。次いで、クライアントの書き込み対象データA2については、10個のデータサーバU1、U2、U11、U4、U5・・・U10が選択されてよく、当該10個のデータサーバに対して書き込み要求が送信されるが、もちろん、U1、U2、U11および他のデータサーバから成る10個のデータサーバが選択されてもよい。

30

【0076】

次いで、クライアントAの書き込み対象データA2の成功した書き込みの数が $M > N$ であれば、成功した書き込みを有するそれぞれのデータサーバに対してOKメッセージが送信されてよく、データサーバは、OKメッセージに従って自らの状態をアップグレード可能に変更するかどうかを判定してよい。

40

【0077】

このように、 $M < N$ という状況では、クライアントが書き込み対象データを同じR個のデータサーバへ再び書き込んで、当該データサーバの状態をアクティブにアップデートしてよい。その結果、当該データサーバがアップグレード不可能な状態にある時間が短縮され得る。

【0078】

本開示の別の実施形態において、第1フィードバック情報および第2フィードバック情報は、クライアント識別子を含む。

50

## 【 0 0 7 9 】

例えば、クライアント A については、第 1 フィードバック情報または第 2 フィードバック情報をサーバへ送信するとき、第 1 フィードバック情報および第 2 フィードバック情報の両方がクライアント識別子「クライアント A」を含む。

## 【 0 0 8 0 】

次いで、データサーバが自らの状態を判定するために第 1 フィードバック情報または第 2 フィードバック情報が使用されることは、第 2 フィードバック情報を受信した後に、データサーバが、第 2 フィードバック情報内のクライアント識別子をアップグレード不可能なリストに書き込むこと、および、自らの状態をアップグレード不可能なものとしてマーキングすることを行うために、第 2 フィードバック情報が使用されることを含む。

10

## 【 0 0 8 1 】

本開示のこの実施形態では、アップグレード不可能なリストがデータサーバ側で設定される。クライアント A を例として用いると、クライアント A の第 2 フィードバック情報を受信した後、データサーバは、自らの状態をアップグレード不可能な状態に設定して、クライアント A をアップグレード不可能なリストに書き込む。

## 【 0 0 8 2 】

データサーバが第 1 フィードバック情報を受信した後、第 1 フィードバック情報は、データサーバが、第 1 フィードバック情報内のクライアント識別子をアップグレード不可能なリストから削除すること、および、アップグレード不可能なリストが空であると判定した後に自らの状態をアップグレード可能なものとしてマーキングすることを行うために使用される。

20

## 【 0 0 8 3 】

本開示のこの実施形態では、このステップが実行された後、成功した書き込みを有するそれぞれのデータサーバに対してクライアント A が第 2 フィードバック情報を送信し、新たな書き込み対象データが出現し、成功した書き込みを有する以前のデータサーバを含む複数のデータサーバが対象として使用される。次いで、クライアント A により送信される第 2 フィードバックメッセージを受信するデータサーバはその後、クライアント A により送信されるメッセージを受信する。次いで、データサーバにより更に受信されるクライアント A のメッセージが第 1 フィードバックメッセージであれば、クライアント A の記録がアップグレード不可能なリストから削除されてよい。

30

## 【 0 0 8 4 】

もちろん、実際の適用においては、データサーバが特定のクライアントの第 1 フィードバックメッセージを受信し、対応するクライアント識別子は、アップグレード不可能なリストに記録されなくてもよい。その場合は、削除プロセスが実行される必要はない。

## 【 0 0 8 5 】

次いで、アップグレード不可能なリストが空であるとデータサーバが判定すれば、データサーバは、自らの状態をアップグレード可能なものに設定する。

## 【 0 0 8 6 】

もちろん、本開示のこの実施形態において、クライアントは、データサーバにより送信される、アップグレード準備状態を終了するようとの終了通知を更に受信し、当該終了通知に従ってアップグレード準備状態を終了してよい。次いで、クライアントは、通常要求送信論理に従って書き込み要求をデータサーバへ送信する。

40

## 【 0 0 8 7 】

本開示のこの実施形態では、本開示における分散記憶システムをアップグレードするための方法がクライアント側から導入され、それぞれのクライアントが、同じ書き込み対象データに関する書き込み要求を複数のデータサーバへ同時に送信し、次いで、幾つのデータサーバが成功した書き込みを有しているかが解析され、成功した書き込みの数が所定の数より多いかがどうかが判定され、その判定結果に従って、成功した書き込みを有するそれぞれのデータサーバに対して第 1 フィードバック情報または第 2 フィードバック情報が送信される。データサーバが第 1 フィードバック情報を受信した後、第 1 フィードバック

50

情報は、データサーバが、第1フィードバック情報内のクライアント識別子をアップグレード不可能なリストから削除すること、および、アップグレード不可能なリストが空であると判定した後に自らの状態をアップグレード可能なものとしてマーキングすることを行うために使用される。従って、本開示の実施形態において、アップグレード制御サーバが、ローリング方式でアップグレードするようデータサーバを制御する場合は、高水準サービスが停止する必要はない。なぜなら、クライアントがデータサーバの状態を制御し、任意のクライアントの書き込み対象データがバックアップのために少なくとも所定の数のデータサーバへ書き込まれることが保証されるからである。更には、分散記憶システムからクライアントへのより短い応答時間が保証され得ることで、データの信頼性が高まり、ユーザデータの損失という危険が著しく低下する。

10

【0088】

本開示のこの実施形態は、分散記憶システムのデータサーバ側に適用される。

【0089】

図2は、本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするための方法を示すフロー図である。具体的に言うと、当該方法は以下のステップを含んでよい。

【0090】

ステップ210：クライアントから送信される第1フィードバック情報または第2フィードバック情報を受信する。ここで、第1フィードバック情報または第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される。

20

【0091】

図1におけるクライアント側の説明を参照すると、クライアントがR個のデータサーバに対して書き込み対象データA1に関する書き込み要求を送信した後、成功した書き込みの数Mが所定の数Nより大きければ、M個のデータサーバに対して第1フィードバックメッセージが送信される。MがNと等しい場合は、M個のデータサーバに対して第2フィードバックメッセージが送信される。M<Nならば、当該複数のデータサーバ以外の少なくとも1つのデータサーバに対して書き込み対象データA1に関する書き込み要求が送信される。次いで、当該少なくとも1つのデータサーバにより返される応答が受信され、当該応答に従って、全ての成功した書き込みの現在の数Mが成功した書き込みの以前の数と組み合わせて所定の数Nと等しくなるかどうか判定される。M=Nならば、M個のデータサーバに対して第2フィードバック情報が送信される。クライアントが第1フィードバック情報および第2フィードバック情報を返す具体的なプロセスについては、図1の説明が参照されてよく、本明細書で再び説明されることはない。

30

【0092】

関連して言うと、分散記憶システム内のそれぞれのデータサーバは、それぞれのクライアントにより返されるOKメッセージ等の第1フィードバック情報を受信してよい。または、それぞれのクライアントにより返されるHOLDメッセージ等の第2フィードバックメッセージが受信されてもよい。

40

【0093】

一実施形態において、当該方法は、ステップ210の前に以下のステップを更に含む。

【0094】

ステップ201：アップグレード制御サーバにより送信される第1アップグレード通知を受信する。

【0095】

ステップ202：第1アップグレード通知に従ってアップグレード準備状態へと移行すること、および、クライアントがアップグレード準備状態へと移行するよう、クライアントのアクセス要求を受信した後に第2アップグレード通知をクライアントへ送信することを行う。

50

## 【 0 0 9 6 】

本開示のこの実施形態では、まず、分散記憶システムのアップグレード制御サーバが第1アップグレード通知をそれぞれのデータサーバへ送信する。それに応じて、それぞれのデータサーバは、第1アップグレード通知を受信し、次いで、第1アップグレード通知に従ってアップグレード準備状態へと移行する。

## 【 0 0 9 7 】

次いで、アップグレード準備状態に移行するデータサーバへのアクセスについては、特定のクライアントのアクセス要求が受信された後に、第2アップグレード通知がクライアントに返される。クライアントは次いで、第2アップグレード通知に従ってアップグレード準備状態へと移行する。

10

## 【 0 0 9 8 】

データサーバがアップグレード準備状態へと移行するときには、データサーバにより実行される、本開示のこの実施形態におけるステップ210および220等のステップへとフローが移行してよい。

## 【 0 0 9 9 】

ステップ220：データサーバがアップグレードされなかった場合は、第1フィードバック情報または第2フィードバック情報に従ってデータサーバの状態を判定する。当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれる。

## 【 0 1 0 0 】

当該状態は、アップグレード制御サーバがローリング方式で、データサーバを選択し、当該データサーバに、アップグレード動作を実行するよう通知するために使用される。

20

## 【 0 1 0 1 】

実際の適用においては、全てのデータサーバがアップグレード準備状態へと移行したばかりであれば、それぞれのデータサーバがアップグレードされることはない。しかしながら、幾つかのデータサーバが継続的にアップグレードされているときには、これらのデータサーバがアップグレードされた状態にある。例えば、成功裏にアップグレードされたデータサーバが自らの状態をDONEに設定する。この時点で、成功裏にアップグレードされたデータサーバは、第1フィードバック情報または第2フィードバック情報を受信する。しかしながら、データサーバが第1フィードバック情報または第2フィードバック情報を更に処理することはなく、アップグレードされた状態のままである。

30

## 【 0 1 0 2 】

アップグレードされていない状態にあるデータサーバだけが、第1フィードバック情報または第2フィードバック情報に従ってデータサーバ自らの状態を判定する。

## 【 0 1 0 3 】

分散クラスタ内のアップグレードされなかったそれぞれのデータサーバについては、データサーバは、データサーバが第1フィードバック情報または第2フィードバック情報を受信するときに、第1フィードバック情報または第2フィードバック情報に従って自らがアップグレード可能な状態にあるのか、アップグレード不可能な状態にあるのかを判定してよいことが理解され得る。アップグレード可能な状態は例えばOKであり、アップグレード不可能な状態は例えばHOLDである。例えば、特定のデータサーバが第2フィードバック情報を受信する場合は、データサーバが自らの状態をアップグレード不可能なものに設定してよい。特定のデータサーバが第1フィードバック情報を受信する場合は、データサーバが自らの状態をアップグレード可能なものに設定してよい。

40

## 【 0 1 0 4 】

本開示の別の実施形態において、第1フィードバック情報および第2フィードバック情報は、クライアント識別子を含む。次いで、ステップ220は、図1のステップ150に基づいてサブステップ221および222を含んでよい。

## 【 0 1 0 5 】

サブステップ221：第2フィードバック情報を受信すると、第2フィードバック情報内のクライアント識別子をアップグレード不可能なリストに書き込んで、自らの状態をア

50

アップグレード不可能なものとしてマーキングする。

【0106】

本開示のこの実施形態では、アップグレード不可能なリストがそれぞれのデータサーバに予め設定され、HOLDメッセージを送信するクライアントの識別子等の情報を記録するのに使用される。

【0107】

例えば、クライアントAが複数のデータサーバに対して書き込み対象データA1に関する書き込み要求を送信する。ここで、データサーバU1、U2、U3は成功した書き込みを有する。M=Nならば、クライアントがM個の対応するデータサーバに対してHOLDメッセージを送信する。HOLDメッセージは、HOLD命令およびクライアント識別子を含む。次いで、M個のデータサーバは、HOLDメッセージを受信するとクライアントAをアップグレード不可能なリストに記録する。もちろん、相応して現在時刻が記録されてよい。

10

【0108】

例えば、クライアントBが複数のデータサーバに対して書き込み対象データB1に関する書き込み要求を送信する。ここで、複数のデータサーバにはデータサーバU1も含まれ、M=Nである。次いで、データサーバU1のアップグレード不可能なリストがクライアントBを更に記録する。もちろん、現在時刻も記録されてよい。

【0109】

表1には、データサーバU1のアップグレード不可能なリストの記録の例が示されている。

20

【表1】

クライアント識別子	時刻
クライアントA	2015. 10. 01、10 : 00 : 00 : 00
クライアントB	2015. 10. 01、12 : 01 : 00 : 00

30

【0110】

同じクライアントにより再び送信されるHOLDメッセージについては、アップグレード不可能なリストは、クライアントに対応するクライアント識別子を1つだけ記録してもよいし、クライアントに対応するクライアント識別子を複数記録してもよい。実際の適用において、同じクライアントにより異なる時刻に送信されるHOLDメッセージについては、クライアント識別子が1つだけ記録されてよく、次いで、それぞれの時刻が時刻フィールドに記録される。

40

【0111】

HOLDメッセージが受信された後、データサーバ自らの状態は、HOLD状態へと更に変更され、データサーバがアップグレード不可能であることを示す。

【0112】

一実施形態において、当該方法は、サブステップ221の後にサブステップB21を更に含む。

【0113】

サブステップB21：アップグレード不可能なリスト内のクライアント識別子については、対応するクライアントの第2フィードバックメッセージが所定の数の期間内に受信さ

50

れなかったかどうかを判定し、対応するクライアントの第2フィードバックメッセージが所定の数の期間内に受信されなかった場合に、クライアント識別子をアップグレード不可能なリストから削除する。

【0114】

本開示のこの実施形態では、期間T1が予め設定されてよく、アップグレード不可能なリスト内のクライアント識別子に対応するクライアントについては、クライアントにより送信されるHOLDメッセージをデータサーバが所定の数の期間内に再び受信していない場合は、アップグレード不可能なリスト内のクライアント識別子の記録が削除される。

【0115】

例えば、前述のデータサーバU1の表1については、クライアントAのHOLDメッセージが3 \* T1という時間内に受信されなかった場合に、表1内のクライアントAの記録が削除される。

10

【0116】

このステップにより、データサーバがアップグレード不可能な状態へ移行した後、アップグレード不可能な状態にずっと留まるのを防ぐことができる。

【0117】

サブステップ222：第1フィードバック情報を受信すると、第1フィードバック情報内のクライアント識別子をアップグレード不可能なリストから削除し、アップグレード不可能なリストが空であると判定した後に自らの状態をアップグレード可能なものとしてマーキングする。

20

【0118】

例えば、クライアントAにより送信されるOKメッセージを前述のデータサーバU1が再び受信し、次いで、クライアントAの記録が表1から削除される。データサーバが表1内の全ての記録を削除した後に表1が空であれば、データサーバは、自らの状態をアップグレード可能なものに設定する。

【0119】

本開示の別の実施形態におけるサブステップ222は、サブステップB11 ~ B14を含む。

【0120】

サブステップB11：アップグレード不可能なリストが第1フィードバック情報内のクライアント識別子を含むかどうかを判定し、含む場合はサブステップB12へと移行する。

30

【0121】

サブステップB12：クライアント識別子をアップグレード不可能なリストから削除する。

【0122】

1つのデータサーバが使用を目的として異なるクライアントへ割り当てられてよいので、異なるクライアントの第1フィードバック情報が受信されてよい。次いで、受信される特定のクライアントの第1フィードバック情報に基づいて、アップグレード不可能なリストが第1フィードバック情報のクライアント識別子を含むかどうかを判定される。アップグレード不可能なリストが第1フィードバック情報のクライアント識別子を含む場合は、アップグレード不可能なリストからクライアント識別子が削除される。アップグレード不可能なリストが第1フィードバック情報のクライアント識別子を含まない場合は、その後の動作が実行されなくてもよい。

40

【0123】

例えば、前述のデータサーバU1については、クライアントAにより改めて送信される、クライアント識別子を含むOKメッセージが受信された場合は、表1内のクライアント識別子とのマッチングに当該クライアント識別子が使用されてよい。クライアントAが存在することが分かった場合は、次いで、クライアントAの記録が消去される。クライアントCにより送信されるOKメッセージが受信されて、クライアントCが表1に記録されて

50

いないことが分かった場合は、何の動作も実行されない。

【0124】

前述のステップにより、データサーバが簡単な論理で自らのアップグレード状態を管理するのが容易になる。

【0125】

サブステップB13：アップグレード不可能なリストが空であるかどうかを判定し、空であればサブステップB14へと移行する。アップグレード不可能なリストが空でない場合は、データサーバが自らの状態をアップグレード不可能なものに維持する。

【0126】

サブステップB14：自らの状態をアップグレード可能なものとしてマーキングする。

10

【0127】

アップグレード不可能なリストが空でない場合は、データサーバが自らの状態をアップグレード不可能なものとして維持する。

【0128】

例えば、前述の例において、データサーバU1は、まずクライアントAにより送信されるOKメッセージを受信し、クライアントAの記録を消去し、表1がまだ空ではないと判定する。従って、データサーバU1はHOLD状態のままである。次いで、クライアントBのOKメッセージが受信された場合は、表1内のクライアントBの記録がサブステップ221で削除され、この時点で表1が空であると判定される。次いで、データサーバU1は、自らのHOLD状態をOK状態に変更し、データサーバU1がアップグレード可能であることを示す。

20

【0129】

次いで、それぞれのデータサーバの状態に従って、アップグレード制御サーバはローリング方式で、データサーバを選択し、当該データサーバに、アップグレード動作を実行するよう通知してよい。アップグレード制御サーバがそれぞれのデータサーバのアップグレードを具体的に制御するプロセスについては、図3の説明が参照されてよい。

【0130】

もちろん、本開示のこの実施形態において、データサーバは、アップグレード制御サーバにより送信される、アップグレード準備状態を終了するようとの終了通知を更に受信し、当該終了通知に従ってアップグレード準備状態を終了してよい。一方で、データサーバは、クライアントがアップグレード準備状態を終了するよう、当該終了通知に従って終了通知をクライアントへ送信してよい。次いで、データサーバは、通常の処理論理に従ってクライアントの書き込み要求を処理する。

30

【0131】

本開示のこの実施形態では、本開示における分散記憶システムをアップグレードするための方法がデータサーバ側で導入される。クライアントにより送信される、受信される第1フィールドバック情報または第2フィールドバック情報に従って、データサーバは、自らがアップグレード可能な状態にあるのか、アップグレード不可能な状態にあるのかを判定する。第1フィールドバック情報または第2フィールドバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後、成功した書き込みの数と所定の数とを比較した結果に従って取得され、当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、当該状態は、アップグレード制御サーバがローリング方式で、データサーバを選択し、当該データサーバに、アップグレード動作を実行するために使用される。従って、クライアントがデータサーバの状態を制御することにより、任意のクライアントの書き込み対象データがバックアップのために少なくとも所定の数のデータサーバへ書き込まれることが保証される。アップグレード制御サーバが、ローリング方式でアップグレードするようそれぞれのデータサーバを制御する場合は、高水準サービスが停止する必要はない。更には、分散記憶システムからクライアントへのより短い応答時間が保証され得ることで、データの信頼性が高まり、ユーザデータの損失という危険が著しく低下する。

40

50

## 【 0 1 3 2 】

図 3 は、本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするための方法を示すフロー図である。具体的に言うと、当該方法は以下のステップを含んでよい。

## 【 0 1 3 3 】

ステップ 3 1 0 : それぞれのデータサーバの状態を取得する。当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、それぞれのデータサーバは一方の状態を有する。データサーバの状態は、第 1 フィードバック情報または第 2 フィードバック情報に従って判定され、第 1 フィードバック情報または第 2 フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される。

10

## 【 0 1 3 4 】

図 1 および図 2 の説明を参照すると、分散記憶システム内のそれぞれのデータサーバは、クライアントによりフィードバックされる第 1 フィードバックメッセージおよび / または第 2 フィードバックメッセージに従って自らの状態を判定してよい。当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれる。

## 【 0 1 3 5 】

それぞれのデータサーバは、同時に一方の状態だけを有し得る。例えば、特定のデータサーバがアップグレード可能な状態にある場合は、当該データサーバが他の状態を有することはできない。他の状況についても同じことが当てはまり、本明細書で再び説明されることはない。

20

## 【 0 1 3 6 】

次いで、本開示のこの実施形態では、アップグレード制御サーバがそれぞれのデータサーバの状態を取得できる。

## 【 0 1 3 7 】

アップグレード制御サーバがデータサーバの状態を取得するための具体的な取得方法は多数あり得るが、本開示で限定されることはない。

## 【 0 1 3 8 】

ステップ 3 2 0 : アップグレード可能な状態にある少なくとも 1 つのデータサーバに、アップグレード動作を実行するようローリング方式で通知する。この通知に従って、当該データサーバはアップグレード動作を実行する。

30

## 【 0 1 3 9 】

本開示のこの実施形態では、それぞれのデータサーバがクライアントのフィードバック情報に従って自らの状態を設定するので、アップグレード制御サーバが、アップグレード動作を実行するようアップグレード可能な状態にあるデータサーバをローリング方式で制御してよい。

## 【 0 1 4 0 】

例えば、アップグレード可能なデータサーバの群が毎回選択され、データサーバのこの群がアップグレード動作を実行するよう通知される。データサーバのこの群は、アップグレード通知を受信した後に再起動およびアップグレードされてよい。

40

## 【 0 1 4 1 】

本開示の別の実施形態において、ステップ 3 2 0 は、サブステップ 3 2 1 および 3 2 2 を含む。

## 【 0 1 4 2 】

サブステップ 3 2 1 : 毎回アップグレード可能な状態にある少なくとも 1 つのデータサーバを選択し、アップグレード可能な状態にある当該少なくとも 1 つのデータサーバに、アップグレード動作を実行するよう通知する。

## 【 0 1 4 3 】

例えば、アップグレード可能な状態にあるデータサーバには、U 1、U 2、U 3、U 1

50

0、U11・・・U20等が含まれる。次いで、アップグレード制御サーバが、K個のデータサーバ、例えば3つのデータサーバを毎回そこから選択し、当該データサーバに、アップグレード動作を実行するよう通知してよい。Kは0より大きい整数である。もちろん、Kは実際の必要に従って設定されてよく、本開示により限定されることはない。

【0144】

サブステップ322：アップグレード可能な状態にある少なくとも1つのデータサーバがアップグレード動作を完全に終了しているかどうかをモニタリングし、終了している場合はサブステップ321へと移行する。

【0145】

例えば、アップグレード制御サーバは、前述のデータサーバU1、U2およびU3に、アップグレード動作を実行するよう通知し、次いで、データサーバU1、U2およびU3は、再起動およびアップグレードされる。データサーバが成功裏にアップグレードされた後、当該データサーバは自らの状態をアップグレードされた状態、例えばDONEに変更し得る。

10

【0146】

次いで、アップグレード制御サーバは、これらのデータサーバの状態がDONEであるかどうかをモニタリングしてよく、DONEであればアップグレードは成功である。

【0147】

もちろん、実際の適用においては、1つまたは複数のデータサーバが、データサーバのこの群のアップグレード中にアップグレードに失敗することもある。アップグレード制御サーバは次いで、データサーバがアップグレードに成功したのか、またはアップグレードに失敗したのかをモニタリングする。実際の適用において、機械の再起動に失敗した、または再起動後のシステムバージョンが変わらないといった状況が起きた場合は、アップグレードに失敗したと判定されてよい。

20

【0148】

アップグレード動作を終了する前述のステップがアップグレードの成功およびアップグレードの失敗を両方とも含んでよいことに留意すべきである。アップグレード動作を完全に終了するステップは、全てのデータサーバが成功裏にアップグレードされること、および、幾つかのデータサーバだけが成功裏にアップグレードされた場合に残りのデータサーバがアップグレードに失敗することを伴う。

30

【0149】

U1、U2およびU3の状態が全てDONEであれば、次回に備えてアップグレード可能な状態にある少なくとも1つのデータサーバを選択すること、および、アップグレード可能な状態にある当該少なくとも1つのデータサーバに、アップグレード動作を実行するよう通知することを行うステップへとフローが移行する。

【0150】

アップグレード制御サーバにより通知される全てのデータサーバがアップグレード動作を完全に終了していることをアップグレード制御サーバがモニタリングした場合は、アップグレード制御サーバは、データサーバの次の群をアップグレード可能な状態にあるデータサーバからアップグレードのためにローリング方式で選択してよい。

40

【0151】

このようにして、極度に頻繁なローリングが回避され、データの消失という危険が低下する。

【0152】

本開示の別の実施形態において、当該方法は、データサーバがアップグレード動作を実行した後に以下のステップを更に含む。

【0153】

サブステップ323：任意のデータサーバのアップグレード動作の結果がアップグレードの失敗であることがモニタリングされた場合は、当該データサーバをアップグレードブラックリストに追加し、当該データサーバのアップグレードを中断する。

50

## 【 0 1 5 4 】

特定のデータサーバがアップグレードに失敗した場合は、アップグレード制御サーバが当該データサーバをアップグレードブラックリストに追加し、当該データサーバのアップグレードを中断する。これらのデータサーバは次いで、オフラインになるのを待つが、または手動での修理を待つ。

## 【 0 1 5 5 】

一実施形態において、本開示のこの実施形態における複数のデータサーバは、少なくとも2つのラック上に設置されてよい。実際の適用においては、もちろん、分散記憶システム内の様々なデータサーバが複数のラック上に設置されてよく、1つのラックは、1つのデータサーバ・サブクラス用である。

10

## 【 0 1 5 6 】

更に、ステップ320はサブステップC11を含む。

## 【 0 1 5 7 】

サブステップC11：毎回アップグレード可能な状態にあるデータサーバを最も多く有するラックを選択し、当該ラック内のデータサーバに、アップグレード動作を実行するよう通知する。この通知に従って、当該ラック内のそれぞれのデータサーバは、自らの状態をチェックする。データサーバがアップグレード可能な状態にある場合は、当該データサーバが再起動およびアップグレードされる。データサーバがアップグレード不可能な状態またはアップグレード終了状態にある場合は、当該データサーバが再起動およびアップグレードを拒否する。

20

## 【 0 1 5 8 】

本開示のこの実施形態では、データサーバがラック上に設置され、データサーバの群が1つのラック上に設置される。しかしながら、多数のクライアントについては、異なるラックの異なるデータサーバにアクセスしてよい。それぞれのラックがアップグレード可能な状態にあるデータサーバを有してもよいし、アップグレード不可能な状態にあるデータサーバを有してもよいし、アップグレードを終了するデータサーバを有してもよい。

## 【 0 1 5 9 】

アップグレード制御サーバがアップグレード通知を送信しやすくすべく、本開示のこの実施形態では、ラック単位でアップグレード通知が送信される。例えば、ラックのIPセグメントは、200.200.200.\*\*\*であり、本開示のこの実施形態におけるアップグレード制御サーバは、200.200.200.\*\*\*に関する通知を1つ生成すること、および、当該ラックのそれぞれのデータサーバが当該通知を受信できるよう、当該通知を当該ラックにブロードキャストすることを行う必要がある。

30

## 【 0 1 6 0 】

次いで、ラック内のデータサーバが前述のアップグレード通知を受信した後、データサーバはまず、自らの状態がOKであるかどうかを判定する。状態がOKであれば、データサーバは、再起動およびアップグレードされる。代わりに状態がHOLDまたはDONEであれば、データサーバは、再起動およびアップグレードを拒否する。

## 【 0 1 6 1 】

本開示のこの実施形態では、アップグレードを迅速に終了すべく、OK状態にあるデータサーバを最も多く有するラックが選択され、当該ラックに対してアップグレード通知が送信される。

40

## 【 0 1 6 2 】

例えば、アップグレード可能な状態にあるデータサーバには、U1、U2、U3、U10、U11・・・U20等が含まれる。次いで、K個のデータサーバ、例えば3つのデータサーバが毎回そこから選択されてよく、これらは、アップグレード動作を実行するよう通知される。

## 【 0 1 6 3 】

一実施形態において、当該方法は、サブステップC11の後にサブステップC12を更に含む。

50

## 【0164】

サブステップC12：ラック内のデータサーバがアップグレード動作を完全に終了しているかどうかをモニタリングし、終了している場合はアップグレード通知サブモジュールへと移行する。

## 【0165】

次いで、ラック内のアップグレード可能な状態にあるデータサーバのアップグレード動作が全て終了していることがモニタリングされた場合は、次のラックが選択されてよく、次のラックのデータサーバに対してアップグレード通知が送信される。この周期は、全てのデータサーバがアップグレードを終了するまで繰り返す。

## 【0166】

前述のプロセスにより、アップグレードに失敗したデータサーバを除く全てのデータサーバがアップグレードを終了していることをアップグレード制御サーバがモニタリングした場合、アップグレード制御サーバは、それぞれのデータサーバにアップグレード準備状態を終了するよう通知し、通常の処理論理に戻ってよい。次いで、それぞれのデータサーバは、当該データサーバへアクセスしているクライアントにアップグレード準備状態を終了するよう通知する。次いで、クライアントは、通常の処理論理に戻る。ステップ310および320が実行される必要はない。

## 【0167】

本開示のこの実施形態では、本開示における分散記憶システムをアップグレードするための方法がアップグレード制御サーバ側で導入される。アップグレード制御サーバが、ローリング方式でアップグレードするようデータサーバを制御する場合は、高水準サービスが停止する必要はない。なぜなら、クライアントがデータサーバの状態を制御し、任意のクライアントの書き込み対象データがバックアップのために少なくとも所定の数のデータサーバへ書き込まれることが保証されるからである。更には、分散記憶システムからクライアントへのより短い応答時間が保証され得ることで、データの信頼性が高まり、ユーザデータの損失という危険が著しく低下する。更に、当該方法は、サービスが影響を受けないように保証しつつ、アップグレード中の機械の予期せぬ例外を許容できる。加えて、大量のデータを移行する必要なしに迅速なアップグレードが実現される。

## 【0168】

図4は、本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするための方法を示すフロー図である。具体的に言うと、当該方法は以下のステップを含んでよい。

## 【0169】

ステップ410：クライアントが、同じ書き込み対象データに関する複数のデータサーバに対して書き込み要求を送信する。

## 【0170】

ステップ412：クライアントは、データサーバの各々により返される応答を受信し、当該応答に従って、成功した書き込みの数が所定の数より多いかどうかを判定する。成功した書き込みの数が所定の数より多い場合は、ステップ414へと移行する。成功した書き込みの数が所定の数より多くない場合は、ステップ416へと移行する。

## 【0171】

ステップ414：クライアントは、成功した書き込みを有するそれぞれのデータサーバに対して第1フィードバック情報を送信する。

## 【0172】

ステップ416：クライアントは、成功した書き込みを有するそれぞれのデータサーバに対して第2フィードバック情報を送信する。

## 【0173】

本開示の別の実施形態において、第1フィードバック情報および第2フィードバック情報は、クライアント識別子を含む。成功した書き込みを有するそれぞれのデータサーバに対して第2フィードバック情報を送信するステップの後、当該方法は更に以下のステップ

10

20

30

40

50

を含む。

【0174】

ステップ417（図示せず）：新たな書き込み対象データが現れた場合は、以前に成功した書き込みがあったデータサーバを含む当該複数のデータサーバを対象として使用され、当該複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求が送信される。

【0175】

ステップ418：データサーバは、クライアントにより送信される第1フィードバック情報または第2フィードバック情報を受信する。

【0176】

ステップ420：データサーバがアップグレードされていない状況では、第1フィードバック情報または第2フィードバック情報に従ってデータサーバの状態を判定する。当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれる。

【0177】

一実施形態において、ステップ420は、ステップ417に基づいてサブステップD11～D16を含む。

【0178】

サブステップD11：第2フィードバック情報を受信すると、第2フィードバック情報内のクライアント識別子をアップグレード不可能なリストに書き込んで、自らの状態をアップグレード不可能なものとしてマーキングする。

【0179】

サブステップD12：アップグレード不可能なリスト内のクライアント識別子については、対応するクライアントの第2フィードバックメッセージが所定の数の期間内に受信されなかったかどうかを判定する。対応するクライアントの第2フィードバックメッセージが所定の数の期間内に受信されなかった場合は、サブステップD13へと移行する。対応するクライアントの第2フィードバックメッセージが所定の数の期間内に受信された場合は、アップグレード不可能な状態が維持される。

【0180】

サブステップD13：クライアント識別子をアップグレード不可能なリストから削除する。サブステップD15へと移行する。

【0181】

サブステップD14：アップグレード不可能なリストが第1フィードバック情報内のクライアント識別子を含むかどうかを判定し、含む場合はサブステップD13へと移行する。アップグレード不可能なリストが第1フィードバック情報内のクライアント識別子を含まない場合は、フローがサブステップD15へと移行する。

【0182】

サブステップD15：アップグレード不可能なリストが空であるかどうかを判定し、空であればサブステップD16へと移行する。

【0183】

サブステップD16：自らの状態をアップグレード可能なものとしてマーキングする。

【0184】

ステップ422：アップグレード制御サーバがそれぞれのデータサーバの状態を取得する。

【0185】

ステップ424：アップグレード制御サーバは、アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で通知する。

【0186】

ステップ426：この通知に従って、データサーバは、アップグレード動作を実行する。

。

10

20

30

40

50

## 【 0 1 8 7 】

もちろん、この実施形態において、クライアント側のステップの原理については、図 1 の説明が参照されてよく、データサーバ側のステップの原理については、図 2 の説明が参照されてよく、アップグレード制御サーバ側のステップの原理については、図 3 の説明が参照されてよい。本明細書でこの説明が再び提供されることはない。

## 【 0 1 8 8 】

本開示のこの実施形態では、本開示における分散記憶システムをアップグレードするための方法が、クライアント、データサーバおよびアップグレード制御サーバを含む 3 つの側面で導入される。分散記憶システムへアクセスしているクライアントについては、それぞれのクライアントが、同じ書き込み対象データに関する書き込み要求を複数のデータサーバへ同時に送信し、次いで、幾つのデータサーバが成功した書き込みを有しているかが解析され、成功した書き込みの数が所定の数より多いかが判定され、その判定結果に従って、成功した書き込みを有するそれぞれのデータサーバに対して第 1 フィードバック情報または第 2 フィードバック情報が送信される。受信される第 1 フィードバック情報または第 2 フィードバック情報に従って、データサーバは、自らがアップグレード可能な状態にあるのか、アップグレード不可能な状態にあるのかを判定する。アップグレード制御サーバは、アップグレード動作を実行するようアップグレード可能な状態にあるデータサーバをローリング方式で制御してよい。前述のプロセスにより、アップグレード制御サーバが、ローリング方式でアップグレードするようアップグレード可能な状態にあるデータサーバを制御する場合は、高水準サービスが停止する必要はない。なぜなら、クライアントがデータサーバの状態を制御し、任意のクライアントの書き込み対象データがバックアップのために少なくとも所定の数のデータサーバへ書き込まれることが保証されるからである。更には、分散記憶システムからクライアントへのより短い応答時間が保証され得ることで、データの信頼性が高まり、ユーザデータの損失という危険が著しく低下する。更に、当該方法は、サービスが影響を受けないように保証しつつ、アップグレード中の機械の予期せぬ例外を許容できる。

## 【 0 1 8 9 】

図 5 は、本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするための装置のブロック図である。具体的に言うと、当該装置は、要求送信モジュール 5 1 0、判定モジュール 5 2 0、第 1 フィードバックモジュール 5 3 0 および第 2 フィードバックモジュール 5 4 0 を含んでよい。

## 【 0 1 9 0 】

要求送信モジュール 5 1 0 は、複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信するように構成される。

## 【 0 1 9 1 】

要求送信モジュール 5 1 0 に加えて、本開示の別の実施形態における当該装置は、データサーバへアクセスするときに当該データサーバにより送信される第 2 アップグレード通知を受信するように構成される第 2 アップグレード通知受信モジュールと、第 2 アップグレード通知に従ってアップグレード準備状態へと移行するように構成される第 2 アップグレード準備モジュールとを更に備える。

## 【 0 1 9 2 】

判定モジュール 5 2 0 は、当該データサーバの各々により返される応答を受信すること、および、当該応答に基づいて、成功した書き込みの数が所定の数より多いかを判定することを行うように構成される。

## 【 0 1 9 3 】

第 1 フィードバックモジュール 5 3 0 は、成功した書き込みの数が所定の数より多い場合に、成功した書き込みを有するそれぞれのデータサーバに対して第 1 フィードバック情報を送信するように構成される。

## 【 0 1 9 4 】

第 2 フィードバックモジュール 5 4 0 は、成功した書き込みの数が所定の数より多くな

10

20

30

40

50

い場合に、成功した書き込みを有するそれぞれのデータサーバに対して第2フィードバック情報を送信するように構成される。ここで、第1フィードバック情報または第2フィードバック情報は、データサーバが自らの状態を判定するために使用され、当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、データサーバの状態は、アップグレード制御サーバがデータサーバに、アップグレード動作を実行するようローリング方式で通知するために使用される。

【0195】

本開示の別の実施形態において、第2フィードバックモジュール540は、成功した書き込みの現在の数が所定の数と等しい場合に、成功した書き込みを有するそれぞれのデータサーバに対して第2フィードバック情報を送信するように構成される第2フィードバック情報送信サブモジュールと、成功した書き込みの数が所定の数より少ない場合に、当該複数のデータサーバ以外の少なくとも1つのデータサーバに対して書き込み対象データに関する書き込み要求を送信するように構成される書き込み要求送信サブモジュールと、当該少なくとも1つのデータサーバにより返される応答を受信し、当該応答に従って、成功した書き込みの現在の数が成功した書き込みの以前の数と組み合わせて所定の数と等しくなるかどうかを判定すること、成功した書き込みの現在の数が所定の数と等しい場合に、第2フィードバック情報送信サブモジュールへと移行すること、第2フィードバック情報内のクライアント識別子をアップグレード不可能なリストに書き込んで、自らの状態をアップグレード不可能なものとしてマーキングすることを行うように構成される判定サブモジュールとを備える。

【0196】

データサーバが第1フィードバック情報を受信した後、第1フィードバック情報は、データサーバが、第1フィードバック情報内のクライアント識別子をアップグレード不可能なリストから削除すること、および、アップグレード不可能なリストが空であると判定した後に自らの状態をアップグレード可能なものとしてマーキングすることを行うために使用される。

【0197】

図6は、本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするための装置のブロック図である。具体的に言うと、当該装置は、フィードバック情報受信モジュール610および状態判定モジュール620を含んでよい。

【0198】

フィードバック情報受信モジュール610は、クライアントから送信される第1フィードバック情報または第2フィードバック情報を受信するように構成される。ここで、第1フィードバック情報または第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される。

【0199】

フィードバック情報受信モジュール610に加えて、本開示の別の実施形態における当該装置は、アップグレード制御サーバにより送信される第1アップグレード通知を受信するように構成される第1アップグレード通知受信モジュールと、第1アップグレード通知に従ってアップグレード準備状態へと移行すること、および、クライアントがアップグレード準備状態へと移行するよう、クライアントのアクセス要求を受信した後に第2アップグレード通知をクライアントへ送信することを行うように構成される第1アップグレード準備モジュールとを更に備える。

【0200】

状態判定モジュール620は、データサーバがアップグレードされていない状況において、第1フィードバック情報または第2フィードバック情報に従ってデータサーバの状態を判定するように構成される。当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、当該状態は、アップグレード制御サーバがローリング方式で、データサーバを選択し、当該データサーバに、アップグレード動作を実行するよう通知

10

20

30

40

50

するために使用される。

【0201】

本開示の別の実施形態において、状態判定モジュール620は、アップグレード可能な判定サブモジュールおよびアップグレード不可能な状態判定サブモジュールを備える。

【0202】

アップグレード可能な判定サブモジュールは、第1フィードバック情報を受信すると、第1フィードバック情報内のクライアント識別子をアップグレード不可能なリストから削除すること、および、アップグレード不可能なリストが空であると判定した後に自らの状態をアップグレード可能なものとしてマーキングすることを行うように構成される。

【0203】

実施形態において、アップグレード可能な判定サブモジュールは、アップグレード不可能なリストが第1フィードバック情報内のクライアント識別子を含むかどうかを判定すること、および、含む場合は第1削除サブモジュールへと移行することを行うように構成されるクライアント識別子判定サブモジュールであって、当該第1削除サブモジュールは、クライアント識別子をアップグレード不可能なリストから削除するように構成される、クライアント識別子判定サブモジュールと、アップグレード不可能なリストが空であるかどうかを判定すること、および、空であればアップグレード可能なマーキングサブモジュールへと移行することを行うように構成されるアップグレード不可能なリスト判定サブモジュールであって、当該アップグレード可能なマーキングサブモジュールは、自らの状態をアップグレード可能なものとしてマーキングするように更に構成される、アップグレード不可能なリスト判定サブモジュールとを備える。

【0204】

アップグレード不可能な状態判定サブモジュールは、第2フィードバック情報を受信すると第2フィードバック情報内のクライアント識別子をアップグレード不可能なリストに書き込むこと、および、自らの状態をアップグレード可能なものとしてマーキングすることを行うように構成される。

【0205】

アップグレード不可能な状態判定サブモジュールに加えて、本開示の別の実施形態における当該装置は、アップグレード不可能なリスト内のクライアント識別子について、対応するクライアントの第2フィードバックメッセージが所定の数の期間内に受信されなかったかどうかを判定すること、および、対応するクライアントの第2フィードバックメッセージが所定の数の期間内に受信されなかった場合に、第2削除サブモジュールへと移行することを行うように構成される時間判定サブモジュールであって、当該第2削除サブモジュールは、クライアント識別子をアップグレード不可能なリストから削除するように構成される、時間判定サブモジュールを更に備える。

【0206】

図7は、本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするための装置のブロック図である。具体的に言うと、当該装置は状態取得モジュール710およびアップグレード通知モジュール720を含んでよい。

【0207】

状態取得モジュール710は、データサーバごとに状態を取得するように構成される。当該状態にはアップグレード可能な状態およびアップグレード不可能な状態が含まれ、それぞれのデータサーバは一方の状態を有する。データサーバの状態は、第1フィードバック情報または第2フィードバック情報に従って判定され、第1フィードバック情報または第2フィードバック情報は、クライアントが複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信した後に、成功した書き込みの数と所定の数とを比較した結果に従って取得される。

【0208】

状態取得モジュール710に加えて、本開示の別の実施形態における当該装置は、それぞれのデータサーバがアップグレード準備状態へと移行するよう、第1アップグレード通

10

20

30

40

50

知をそれぞれのデータサーバへ送信するように構成されるアップグレード通知送信モジュールを更に備える。それぞれのデータサーバは、クライアントがアップグレード準備状態へと移行するよう、クライアントのアクセス要求を受信した後に第2アップグレード通知をクライアントへ送信する。

【0209】

アップグレード通知モジュール720は、アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で通知するように構成される。この通知に従って、当該データサーバはアップグレード動作を実行する。

【0210】

本開示の別の実施形態において、アップグレード通知モジュールは、毎回アップグレード可能な状態にある少なくとも1つのデータサーバを選択すること、および、アップグレード可能な状態にある当該少なくとも1つのデータサーバに、アップグレード動作を実行するよう通知することを行うように構成される第1選択サブモジュールと、アップグレード可能な状態にある当該少なくとも1つのデータサーバがアップグレード動作を完全に終了しているかどうかをモニタリングすること、および、終了している場合は第1選択サブモジュールへと移行することを行うように構成される第1モニタリングサブモジュールとを備える。

【0211】

第1モニタリングサブモジュールに次いで、本開示の別の実施形態における当該装置は、任意のデータサーバのアップグレード動作の結果がアップグレードの失敗であることがモニタリングされた場合に、データサーバをアップグレードブラックリストに追加すること、および、データサーバのアップグレードを中断することを行うように構成される中断サブモジュールを更に備える。

【0212】

本開示の別の実施形態において、データサーバは、少なくとも2つのラック上に設置される。次いで、アップグレード通知モジュール720は、毎回アップグレード可能な状態にあるデータサーバを最も多く有するラックを選択すること、および、当該ラック内のデータサーバに、アップグレード動作を実行するよう通知することを行うように構成されるアップグレード通知サブモジュールを備える。この通知に従って、当該ラック内のそれぞれのデータサーバは、自らの状態をチェックする。データサーバがアップグレード可能な状態にある場合は、当該データサーバが再起動およびアップグレードされる。データサーバがアップグレード不可能な状態またはアップグレード終了状態にある場合は、当該データサーバが再起動およびアップグレードを拒否する。

【0213】

アップグレード通知サブモジュールに加えて、本開示の別の実施形態における当該装置は、ラック内のデータサーバがアップグレード動作を完全に終了しているかどうかをモニタリングすること、および、終了している場合はアップグレード通知サブモジュールへと移行することを行うように構成されるモニタリングサブモジュールを更に備える。

【0214】

図8および図8Aを参照されたい。これらの図は、本開示の幾つかの実施形態に係る、分散記憶システムをアップグレードするためのシステムのブロック図である。具体的に言うと、当該システムは、以下のモジュール、すなわちクライアント810、データサーバ820およびアップグレード制御サーバ830を含んでよい。

【0215】

図8Aは、一実施形態における分散記憶システムのアーキテクチャを示す概略図である。本開示のこの実施形態において、それぞれのクライアント810は、分散記憶システム内のR個のデータサーバ820に対して書き込み要求を送信してよく、アップグレード制御サーバ830は、全てのデータサーバのアップグレードプロセスを制御する。

【0216】

10

20

30

40

50

図8は、クライアント810、データサーバ820およびアップグレード制御サーバ830間の接続関係を示している。

【0217】

クライアント810は、複数のデータサーバに対して同じ書き込み対象データに関する書き込み要求を送信するように構成される要求送信モジュール811と、当該データサーバの各々により返される応答を受信すること、および、当該応答に基づいて、成功した書き込みの数が所定の数より多いかどうかを判定することを行うように構成される判定モジュール812と、成功した書き込みの数が所定の数より多い場合に、成功した書き込みを有するそれぞれのデータサーバに対して第1フィードバック情報を送信するように構成される第1フィードバックモジュール813と、成功した書き込みの数が所定の数より多くない場合に、成功した書き込みを有するそれぞれのデータサーバに対して第2フィードバック情報を送信するように構成される第2フィードバックモジュール814とを備える。

10

【0218】

データサーバ820は、クライアントの書き込み要求を受信すること、および、当該クライアントに応答を返すことを行うように構成されるデータ記憶モジュール821と、当該クライアントにより送信される第1フィードバック情報または第2フィードバック情報を受信するように構成されるフィードバック情報受信モジュール822と、データサーバ自らがアップグレードされなかった場合に、第1フィードバック情報または第2フィードバック情報に従ってデータサーバ自らの状態を判定するように構成される状態判定モジュール823と、アップグレード制御サーバの通知に従ってアップグレード動作を実行するように構成されるアップグレードモジュール824とを備える。

20

【0219】

アップグレード制御サーバ830は、それぞれのデータサーバの状態を取得するように構成される状態取得モジュール831と、アップグレード可能な状態にある少なくとも1つのデータサーバに、アップグレード動作を実行するようローリング方式で通知するように構成されるアップグレード通知モジュール832とを備える。

【0220】

本開示のこの実施形態において、クライアントのモジュールについては、図5の説明が参照されてよい。データサーバのモジュールについては、図6の説明が参照されてよい。アップグレード制御サーバのモジュールについては、図7の説明が参照されてよい。これらの構成要素により実行される方法は、前述の図で説明される方法と同様であり、本明細書で再び説明されることはない。

30

【0221】

本開示のこの実施形態では、本開示における分散記憶システムをアップグレードするための方法が、クライアント、データサーバおよびアップグレード制御サーバを含む3つの側面で導入される。分散記憶システムへアクセスしているクライアントについては、それぞれのクライアントが、同じ書き込み対象データに関する書き込み要求を複数のデータサーバへ同時に送信し、次いで、幾つのデータサーバが成功した書き込みを有しているかが解析され、成功した書き込みの数が所定の数より多いかが判定され、その判定結果に従って、成功した書き込みを有するそれぞれのデータサーバに対して第1フィードバック情報または第2フィードバック情報が送信される。受信される第1フィードバック情報または第2フィードバック情報に従って、データサーバは、自らがアップグレード可能な状態にあるのか、アップグレード不可能な状態にあるのかを判定する。アップグレード制御サーバは、アップグレード動作を実行するようアップグレード可能な状態にあるデータサーバをローリング方式で制御してよい。前述のプロセスにより、アップグレード制御サーバが、ローリング方式でアップグレードするようアップグレード可能な状態にあるデータサーバを制御する場合は、高水準サービスが停止する必要はない。なぜなら、クライアントがデータサーバの状態を制御し、任意のクライアントの書き込み対象データがバックアップのために少なくとも所定の数のデータサーバへ書き込まれることが保証されるからである。更には、分散記憶システムからクライアントへのより短い応答時間が保証され

40

50

得ることで、データの信頼性が高まり、ユーザデータの損失という危険が著しく低下する。更に、当該方法は、サービスが影響を受けないように保証しつつ、アップグレード中の機械の予期せぬ例外を許容できる。

【0222】

装置の実施形態は、方法の実施形態と同様なので、比較的簡潔に説明される。関連する部分については、方法の実施形態の説明の部分が参照されてよい。

【0223】

本明細書内の実施形態は、それぞれの実施形態が他の実施形態とは異なる部分を強調する形で漸次説明される。相互参照により、当該実施形態の同一の部分または同様の部分が取得されてよい。

【0224】

当業者であれば、本開示の実施形態が方法、装置またはコンピュータプログラム製品として提供され得ることを理解するはずである。従って、本開示の実施形態は、完全なハードウェアの実施形態、完全なソフトウェアの実施形態、またはソフトウェアとハードウェアとを組み合わせた実施形態として実装されてよい。更に、本開示の実施形態は、内部にコンピュータ使用可能プログラムコードを含む1つまたは複数のコンピュータ使用可能記憶媒体（これらに限定されるわけではないが、磁気ディスクメモリ、CD ROM、光学メモリ等が含まれる）上で実装されるコンピュータプログラム製品の形を取ってよい。

【0225】

典型的な構成において、コンピュータデバイスは、1つまたは複数のプロセッサ（CPU）、入出力インタフェース、ネットワークインタフェースおよびメモリを含む。メモリには、リードオンリメモリ（ROM）またはフラッシュメモリ（FLASH RAM）といった、非永続的メモリ、ランダムアクセスメモリ（RAM）および/または不揮発性メモリ等の形を取った非一時的コンピュータ可読媒体が含まれてよい。メモリは、コンピュータ可読媒体の例である。コンピュータ可読媒体には、任意の方法または技術により情報の記憶を実現できる永続的および非永続的な移動可能媒体並びに移動不可能媒体が含まれる。情報は、コンピュータ可読命令、データ構造、プログラムのモジュール、または他のデータであってよい。コンピュータの記憶媒体の例としては、これらに限定されるわけではないが、相変化メモリ（PRAM）、スタティックランダムアクセスメモリ（SRAM）、ダイナミックランダムアクセスメモリ（DRAM）もしくは他の種類のランダムアクセスメモリ（RAM）、リードオンリメモリ（ROM）、電氣的消去可能プログラマブルリードオンリメモリ（EEPROM）、フラッシュメモリもしくは他のメモリ技術、リードオンリ・コンパクトディスク・リードオンリメモリ（CD ROM）、デジタル多用途ディスク（DVD）もしくは他の光記憶装置、磁気カセット、磁気テープ、磁気ディスク記憶装置もしくは他の磁気記憶デバイス、または、コンピューティングデバイスによりアクセス可能な情報を記憶するために使用され得る任意の他の非伝送媒体が挙げられる。本明細書における定義を考慮すると、コンピュータ可読媒体に、変調されたデータ信号および搬送波等の非持続的コンピュータ可読媒体（一時的媒体）は含まれない。

【0226】

本開示の実施形態は、本開示の実施形態の方法、端末デバイス（システム）およびコンピュータプログラム製品に係るフロー図および/またはブロック図を参照しながら説明される。フロー図および/またはブロック図におけるそれぞれのフローおよび/またはそれぞれのブロック、並びに、フロー図および/またはブロック図におけるフロー同士および/またはブロック同士の組み合わせがコンピュータプログラム命令により実装され得ることを理解すべきである。これらのコンピュータプログラム命令は、機械を生成すべく、汎用コンピュータ、専用コンピュータ、組込みプロセッサ、または任意の他のプログラマブルデータ処理端末デバイスのプロセッサに対して提供されてよい。その結果、任意の他のプログラマブルデータ処理端末デバイスのコンピュータまたはプロセッサにより実行される命令は、フロー図における1つもしくは複数のプロセス、および/または、ブロック図における1つもしくは複数のブロックで指定された機能を実装するための装置を生成する

10

20

30

40

50

## 【0227】

これらのコンピュータプログラム命令は、コンピュータまたは別のプログラマブルデータ処理端末デバイスに特定のやり方で動作するよう指示できるコンピュータ可読メモリに記憶されてもよい。その結果、コンピュータ可読メモリに記憶される命令は、フロー図の1つもしくは複数のフロー、および/または、ブロック図の1つもしくは複数のブロックで指定された機能を実装する命令手段を含む製造品を生産する。

## 【0228】

これらのコンピュータプログラム命令は、コンピュータまたは任意の他のプログラマブルデータ処理端末デバイス上にロードされてもよい。その結果、一連の動作ステップがコンピュータまたは任意の他のプログラマブル端末デバイス上で実行されることにより、コンピュータ実装プロセスが生成される。従って、コンピュータまたは任意の他のプログラマブル端末デバイス上で実行される命令は、フロー図における1つもしくは複数のプロセス、および/または、ブロック図における1つもしくは複数のブロックで指定された機能を実装するためのステップを提供する。

## 【0229】

本開示の様々な実施形態が説明されてきたが、ひとたび基本的な創造的概念を説明しておけば、当業者はこれらの実施形態の他の変形例および改良例を作ることができる。従って、添付の請求項は、当該実施形態、並びに、本開示の範囲内に入る全ての変形例および改良例を含むものとして解釈されるように意図されている。

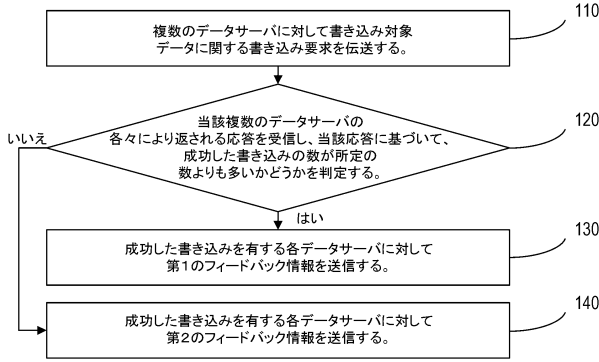
## 【0230】

最後に、本願における第1および第2といった関係用語は、単に1つのエンティティまたは動作を別のエンティティまたは動作と区別するために使用されているに過ぎず、これらのエンティティまたは動作がこうした実際の関係または順序を持つ必要があるわけでもなければ、それを持つと示唆しているわけでもないことを更に留意すべきである。更に「含む(include)」、「含む(comprise)」という用語、またはこれらの他の変形語は、非排他的な包含物をカバーするように意図されている。その結果、一連の要素を含むプロセス、方法、物品または端末デバイスは、当該要素を含むだけではなく、明記されていない他の要素も含むか、または、当該プロセス、当該方法、当該物品もしくは当該端末デバイスに本来備わっている要素を更に含む。「1つの・・・を含む(including one・・・)」という記述により更なる制限なく定義される要素は、当該要素を含むプロセス、方法、物品または端末デバイスにおける追加的な同一要素の存在を除外するものではない。

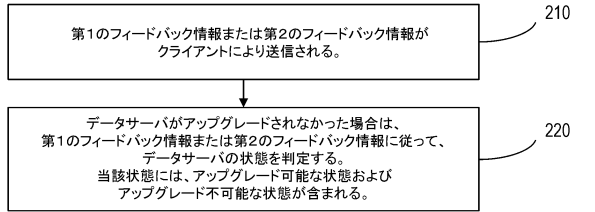
## 【0231】

本開示で提供される、分散記憶システムをアップグレードするための方法、分散記憶装置をアップグレードするための装置、および分散記憶装置をアップグレードするためのシステムが以上に詳しく紹介された。本明細書には、具体的な例を参照しながら本開示の原理および実装が記載されている。上記の実施形態の説明は、単に本開示の方法および本質的な着想を理解する手助けをするために提供されているに過ぎない。当業者であれば、本開示の着想に従って具体的な実装および実施形態に対する変更を加えることができる。上記を考慮すると、本明細書の内容は本開示を限定するものと解釈されるべきではない。

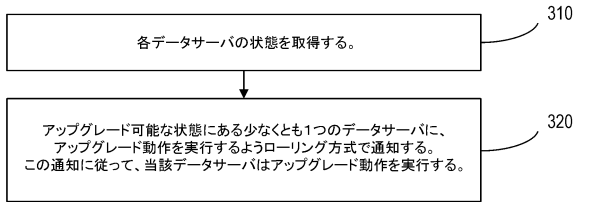
【図1】



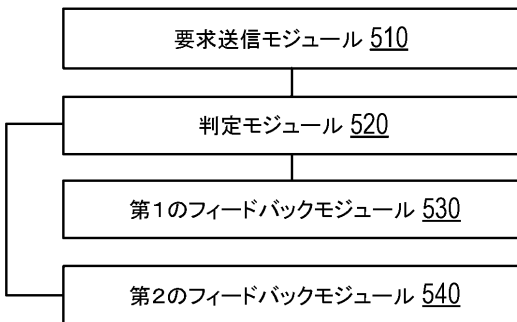
【図2】



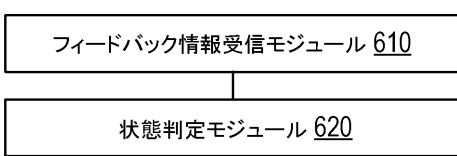
【図3】



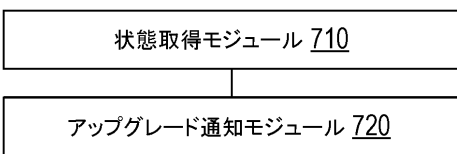
【図5】



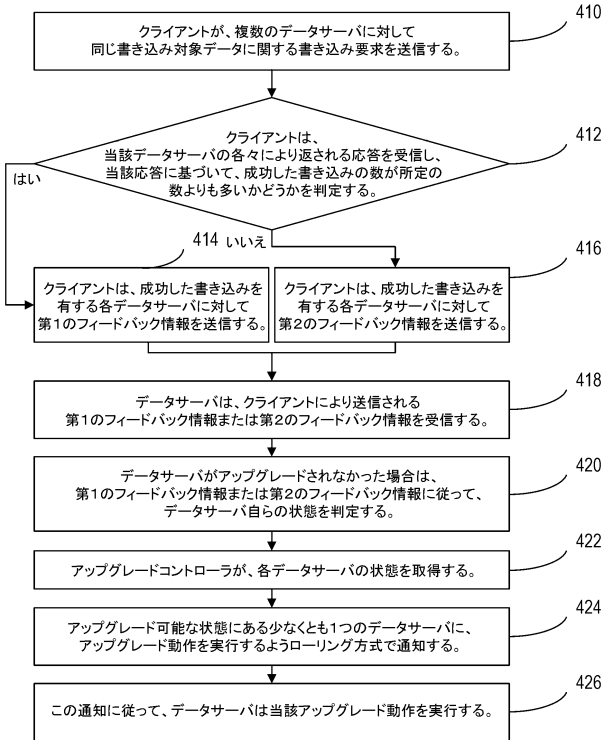
【図6】



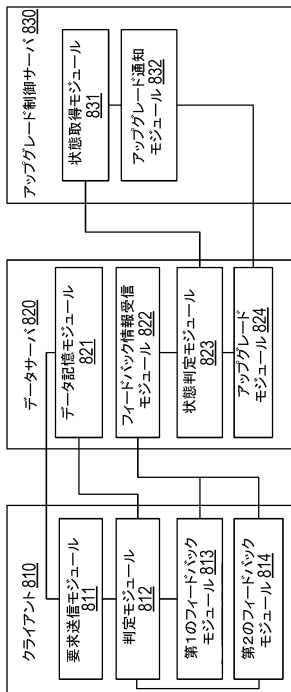
【図7】



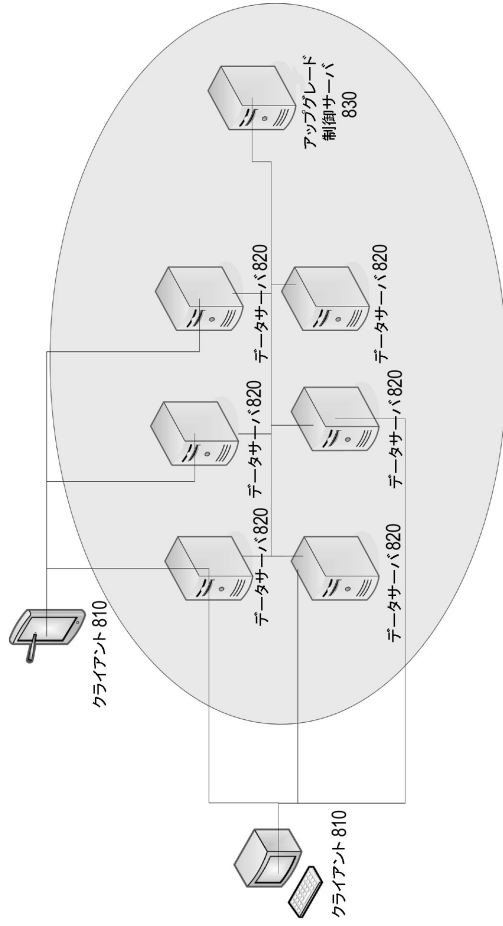
【図4】



【図8】



【 図 8 A 】



## フロントページの続き

- (72)発明者 チュー、ジャジ  
英国領、ケイマン諸島、グランド・ケイマン、ジョージ・タウン、ワン・キャピタル・プレイス、  
フォース・フロア、ピー・オー・ボックス 847 アリババ・グループ・ホールディング・リミ  
テッド内
- (72)発明者 ジャオ、シューチ  
英国領、ケイマン諸島、グランド・ケイマン、ジョージ・タウン、ワン・キャピタル・プレイス、  
フォース・フロア、ピー・オー・ボックス 847 アリババ・グループ・ホールディング・リミ  
テッド内
- (72)発明者 リン、ジャンピン  
英国領、ケイマン諸島、グランド・ケイマン、ジョージ・タウン、ワン・キャピタル・プレイス、  
フォース・フロア、ピー・オー・ボックス 847 アリババ・グループ・ホールディング・リミ  
テッド内
- (72)発明者 グ、ユエシェン  
英国領、ケイマン諸島、グランド・ケイマン、ジョージ・タウン、ワン・キャピタル・プレイス、  
フォース・フロア、ピー・オー・ボックス 847 アリババ・グループ・ホールディング・リミ  
テッド内

審査官 久保 光宏

- (56)参考文献 奥野幹也, 「My SQL Cluster構築・運用バイブル~仕組みからわかる基礎と実践のノウハウ」,  
日本, 株式会社技術評論社, 2012年 4月25日, 初版, 第404~411頁, ISBN: 978-  
4-7741-5053-6  
西村豪生(外6名), 「リソースプールを活用して無停止ソフトウェアアップデートを実現する  
セッション制御サーバミドルウェア」, 電子情報通信学会論文誌B, 日本, [online], 電子情報  
通信学会, 2014年 5月 1日, Vol.J97-B, No.5, 第393~406頁, ISSN: 1881-0209  
, [2014年5月1日検索], インターネット, <URL: [http://search.ieice.org/bin/pdf.php?lang=J  
&year=2014&fname=j97-b\\_5\\_393&abst=>](http://search.ieice.org/bin/pdf.php?lang=J&year=2014&fname=j97-b_5_393&abst=>)  
飯沢篤志(外1名), 「データベースおもしろ講座 -10- 分散データベースシステム」, bit  
, 日本, 共立出版株式会社, 1992年 2月 1日, 1992年2月号 (Vol.24, No.2), 第18  
1~192頁, ISSN: 0385-6984

## (58)調査した分野(Int.Cl., DB名)

G06F8/65-8/658  
G06F16/00-16/958  
CSDB(日本国特許庁)  
IEEEExplore(IEEE)