

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2005-526298

(P2005-526298A)

(43) 公表日 平成17年9月2日(2005.9.2)

(51) Int.Cl.⁷

G06F 12/00

F I

G06F 12/00 533J

G06F 12/00 514M

G06F 12/00 517

テーマコード (参考)

5B082

審査請求 未請求 予備審査請求 有 (全 31 頁)

(21) 出願番号 特願2003-514370 (P2003-514370)
 (86) (22) 出願日 平成14年7月15日 (2002.7.15)
 (85) 翻訳文提出日 平成16年3月15日 (2004.3.15)
 (86) 国際出願番号 PCT/US2002/022366
 (87) 国際公開番号 W02003/009092
 (87) 国際公開日 平成15年1月30日 (2003.1.30)
 (31) 優先権主張番号 60/305,986
 (32) 優先日 平成13年7月16日 (2001.7.16)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 607305,978
 (32) 優先日 平成13年7月16日 (2001.7.16)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 09/975,590
 (32) 優先日 平成13年10月11日 (2001.10.11)
 (33) 優先権主張国 米国 (US)

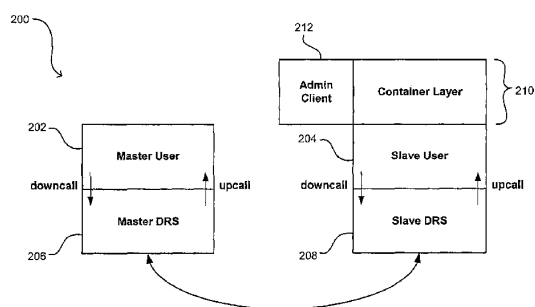
(71) 出願人 500105160
 ビーイーエイ システムズ, インコーポ
 レイテッド
 BEA Systems, Inc.
 アメリカ合衆国 カリフォルニア 951
 31, サン ノゼ, ノース ファース
 ト ストリート 2315
 2315 North First St
 reet, San Jose, CAL
 IFORNIA 95131 U. S. A
 .
 (74) 代理人 100082005
 弁理士 熊倉 禎男
 (74) 代理人 100067013
 弁理士 大塚 文昭

最終頁に続く

(54) 【発明の名称】 データ複製プロトコル

(57) 【要約】

一フェーズ法又は二フェーズ法を用いて、ネットワークを通じてデータを複製することができる。一フェーズ法の場合は、各々のスレーブがマスターからデルタを要求できるように、データのオリジナル・コピーを含むマスター・サーバが、該データの現在の状態についてのバージョン番号をネットワークで送信する。要求されたデルタは、スレーブを適切なバージョンのデータに更新するために必要なデータを含む。二フェーズ法の場合は、マスター・サーバは、情報のパケットを各々のスレーブに送信する。各々のスレーブがコミットを処理できる場合には、スレーブにより情報のパケットをコミットすることができる。



【特許請求の範囲】**【請求項 1】**

ネットワークを通じてマスター・サーバからスレーブ・サーバにデータを複製する方法であって、

前記マスター・サーバ上に格納された前記データの変更に關し、該データの現在の状態についてのバージョン番号を含む情報のパケットを、該マスター・サーバから前記スレーブ・サーバに送信し、

前記スレーブ・サーバ上の前記データが前記パケット内に含まれる前記バージョン番号に対応するように更新されたかどうかを、該スレーブ・サーバが判断することを可能にし

10

、
前記スレーブ・サーバ上の前記データが前記パケット内に含まれる前記バージョン番号に対応しない場合には、該スレーブ・サーバを更新するのに必要とされる情報を含むデルタが、前記マスター・サーバから該スレーブ・サーバに送信されることを要求する、
段階を含むことを特徴とする方法。

【請求項 2】

前記データのオリジナル・コピーを前記マスター・サーバ上に格納する段階をさらに含むことを特徴とする請求項 1 に記載の方法。

【請求項 3】

各々のスレーブ・サーバについて、ローカル・ディスク上の前記データを永続的にキャッシングする段階をさらに含むことを特徴とする請求項 1 に記載の方法。

20

【請求項 4】

前記データが変更された場合には、前記マスター・サーバ上の前記データの現在の状態についての固有のバージョン番号を判断する段階をさらに含むことを特徴とする請求項 1 に記載の方法。

【請求項 5】

ネットワークを通じてマスター・サーバからスレーブ・サーバにデータを複製する方法であって、

前記マスター・サーバ上に格納された前記データの現在の状態に関するバージョン番号を、該マスター・サーバから前記スレーブ・サーバに送信し、

前記スレーブ・サーバが前記マスター・サーバから送信された前記バージョン番号に対応する前記データの現在の状態を反映するように更新されたかどうかを、該スレーブ・サーバが判断することを可能にし、

30

前記スレーブ・サーバが前記マスターにより送信された前記バージョン番号に対応しない場合には、該スレーブ・サーバを更新するのに必要とされる情報を含むデルタが、該マスター・サーバから該スレーブ・サーバに送信されることを要求する、
段階を含むことを特徴とする方法。

【請求項 6】

前記マスター・サーバから前記スレーブ・サーバに前記デルタを送信する段階をさらに含むことを特徴とする請求項 5 に記載の方法。

【請求項 7】

前記スレーブ・サーバへの前記デルタをコミットする段階をさらに含むことを特徴とする請求項 5 に記載の方法。

40

【請求項 8】

前記デルタをコミットした後、前記スレーブ・サーバの前記バージョン番号を更新する段階をさらに含むことを特徴とする請求項 5 に記載の方法。

【請求項 9】

前記バージョン番号を前記マスター・サーバからスレーブ・サーバに定期的に送信する段階をさらに含むことを特徴とする請求項 5 に記載の方法。

【請求項 10】

前記スレーブ・サーバが前記バージョン番号の受信を確認するまで、前記バージョン番

50

号をスレーブ・サーバに送信する段階をさらに含むことを特徴とする請求項 5 に記載の方法。

【請求項 1 1】

スレーブ・サーバを更新するのに必要な前記バージョン番号を有するデータを含ませる段階をさらに含むことを特徴とする請求項 5 に記載の方法。

【請求項 1 2】

データが受け取るとすぐに、前記スレーブを更新するのに必要な該データをコミットする段階をさらに含むことを特徴とする請求項 1 1 に記載の方法。

【請求項 1 3】

前記マスター・サーバから前記デルタを送信する前に、該デルタの範囲を判断する段階をさらに含むことを特徴とする請求項 5 に記載の方法。 10

【請求項 1 4】

マスター・サーバ及び少なくとも 1 つのスレーブ・サーバを含むネットワークを通じてデータを複製する方法であって、

前記マスター・サーバ上に格納された前記データの変更に、該データの現在の状態についてのバージョン番号を含み、さらに、該データの前の変更及び前の変更の各々についてのバージョン番号に関する情報のパケットを、ネットワークで該マスター・サーバからスレーブ・サーバの各々に送信し、

前記スレーブ・サーバが前記バージョン番号に対応するように更新されたかどうかを、各々のスレーブ・サーバが判断することを可能にし、 20

前記スレーブ・サーバが前の変更を逃さなかった場合には、各々のスレーブ・サーバが前記情報をコミットすることを可能にし、

前記スレーブ・サーバが前記情報のパケットをコミットする前に、前の変更が、前記マスター・サーバから該スレーブ・サーバに送信されることを、前の変更を逃した各々のスレーブ・サーバが要求することを可能にする、
段階を含むことを特徴とする方法。

【請求項 1 5】

スレーブ・サーバへの前記情報のパケットをコミットする段階をさらに含むことを特徴とする請求項 1 4 に記載の方法。

【請求項 1 6】

スレーブ・サーバが前記更新をコミットできない場合には、前記情報のパケットの前記コミットを中止する段階をさらに含むことを特徴とする請求項 1 4 に記載の方法。 30

【請求項 1 7】

前記マスター・サーバから前記デルタを送信する前に、該デルタの範囲を判断する段階をさらに含むことを特徴とする請求項 1 4 に記載の方法。

【請求項 1 8】

前記デルタの前の変更の各々の範囲を含ませる段階をさらに含むことを特徴とする請求項 1 4 に記載の方法。

【請求項 1 9】

マスター・サーバ及び少なくとも 1 つのスレーブ・サーバを含むネットワークを通じてデータを複製する方法であって、 40

マスター・サーバ上に格納された前記データの変更に、前の状態についての前のバージョン番号及び該データの新しい状態についての新しいバージョン番号を含み、さらに、該データの前の変更及び各々の前の変更についての前のバージョン番号に関連する情報のパケットを、ネットワークで前記マスター・サーバから各々のスレーブ・サーバに送信し、

前記スレーブ・サーバ上の前記データが前記パケット内に含まれる前記前のバージョン番号に対応するかどうかを、各々のスレーブ・サーバが判断することを可能にし、

前記スレーブ・サーバ上の前記データが前記パケット内に含まれる前記前のバージョン番号に対応する場合には、各々のスレーブ・サーバが前記情報のパケットをコミットする 50

ことを可能にする、

段階を含み、前記コミットが、前記スレーブ・サーバの前記バージョンを前記新しいバージョン番号にも更新し、

前記スレーブ・サーバが前記情報のパケットをコミットする前に、前記スレーブを前記前のバージョン番号に更新するのに必要な前記情報を含むデルタが、前記マスター・サーバから送信されることを、該前のバージョン番号に対応しないスレーブ・サーバの各々が要求することを可能にする、

段階が設けられたことを特徴とする方法。

【請求項 20】

マスター・サーバ及び少なくとも 1 つのスレーブ・サーバを含むネットワークを通じてデータを複製する方法であって、 10

マスター・サーバ上に格納された前記データの変更に關し、前の状態についてのバージョン番号及び該データの新しい状態についてのバージョン番号を含み、さらに、該データの前の変更及び前の変更の各々についてのバージョン番号に関する情報のパケットを、ネットワークで前記マスター・サーバからスレーブ・サーバの各々に送信し、

前記スレーブ・サーバ上の前記データが前記パケット内に含まれる前記前のバージョン番号に対応するかどうかを、各々のスレーブ・サーバが判断することを可能にし、

前記スレーブ・サーバ上の前記データが前記パケットに含まれる前記前のバージョン番号に対応する場合には、各々のスレーブ・サーバが前記情報のパケットをコミットすることを可能にする、 20

段階を含み、前記コミットが、前記スレーブ・サーバの前記バージョンを前記新しいバージョン番号にも更新し、

前記スレーブを前記新しいバージョン番号に更新するのに必要な前記情報を含むデルタが、前記マスター・サーバから送信されることを、前記前のバージョン番号に対応しないスレーブ・サーバの各々が要求することを可能にする、

段階が設けられたことを特徴とする方法。

【請求項 21】

ネットワークを通じて、マスター・サーバかた、少なくとも 1 つのスレーブ・サーバにデータを複製する方法であって、

前記マスター・サーバ上に格納された前記データの変更に關し、前記データの現在の状態についてのバージョン番号を含む情報のパケットを、前記マスター・サーバからスレーブ・サーバに送信し、 30

スレーブ・サーバへの前記情報のパケットを受信し、

前記スレーブ・サーバが前記パケット内に含まれる前記バージョン番号に対応するように更新されたかどうか、さらに、該パケット内に含まれる該バージョン番号に対応するように更新する必要がある場合には、該スレーブ・サーバが該情報のパケットを処理することができるかどうかを、該スレーブ・サーバが判断することを可能にし、

前記スレーブ・サーバが更新されることを必要とするかどうか、及び該スレーブ・サーバが前記更新を処理することができるかどうかを示す信号を、該スレーブ・サーバから前記マスター・サーバに送信し、 40

前記スレーブ・サーバが前記パケット内に含まれる前記情報にコミットすべきであるかどうかを示す応答信号を、前記マスター・サーバから該スレーブ・サーバに送信し、

前記応答信号により前記スレーブ・サーバが前記パケット内に含まれる前記情報にコミットすべきであることが示された場合には、前記スレーブ・サーバへの該情報のパケットをコミットする、

段階を含むことを特徴とする方法。

【請求項 22】

少なくとも 1 つのスレーブ・サーバの各々が前記データをコミットできるかどうかを判断する段階をさらに含むことを特徴とする請求項 21 に記載の方法。

【請求項 23】

少なくとも1つのスレーブ・サーバの各々が、応答を前記マスター・サーバに送り返したかどうかを判断する段階をさらに含むことを特徴とする請求項21に記載の方法。

【請求項24】

少なくとも1つのスレーブ・サーバのいずれかが前記データをコミットできるかどうかを判断する段階をさらに含むことを特徴とする請求項21に記載の方法。

【請求項25】

少なくとも1つのスレーブ・サーバの各々が前記コミットを処理できる場合に限り、前記データをコミットする段階をさらに含むことを特徴とする請求項21に記載の方法。

【請求項26】

少なくとも1つのスレーブ・サーバの各々が前記コミットを処理できない場合に限り、前記データを中止する段階をさらに含むことを特徴とする請求項21に記載の方法。 10

【請求項27】

前記コミットを処理することができるそれらのスレーブへの前記データをコミットする段階をさらに含むことを特徴とする請求項21に記載の方法。

【請求項28】

前記コミットを処理することができなかった前記少なくとも1つのスレーブ・サーバのいずれかに前記更新をマルチキャストする段階をさらに含むことを特徴とする請求項21に記載の方法。

【請求項29】

前記新しいバージョン番号を、前記コミットを処理することができなかった前記少なくとも1つのスレーブ・サーバのいずれかに一定間隔で送信する段階をさらに含むことを特徴とする請求項21に記載の方法。 20

【請求項30】

デルタが、前記コミットを処理することができなかったスレーブ・サーバに送信されることを要求する段階をさらに含むことを特徴とする請求項21に記載の方法。

【請求項31】

ネットワークを通じてデータを複製する方法であって、

(a) 前記複製を、一フェーズ法で達成すべきか、又は二フェーズ法で達成すべきかを判断し、

(b) 前記マスター・サーバ上に格納された前記データの変更に關し、該データの現在の状態についてのバージョン番号を含む情報のパケットを、該マスター・サーバから前記スレーブ・サーバに送信し、 30

スレーブ・サーバへの前記情報のパケットを受信し、

前記スレーブ・サーバ上の前記データが前記バージョン番号に対応するように更新されたかどうかを、前記スレーブ・サーバが判断することを可能にし、

前記スレーブ・サーバが前記バージョン番号に対応しない場合には、該スレーブ・サーバを更新するのに必要とされる情報を含むデルタが、前記マスター・サーバから該スレーブ・サーバに送信されることを要求する、

ことによって、一フェーズ法で達成されるように定められた複製情報を送信し、

(c) 前記マスター・サーバ上に格納された前記データの変更に關し、該データの現在の状態についてのバージョン番号を含む情報のパケットを、該マスター・サーバから前記スレーブ・サーバに送信し、 40

前記スレーブ・サーバが前記バージョン番号に対応するように更新されたかどうか、さらに、該スレーブ・サーバが前記情報のパケットを処理できるかどうかを、該スレーブ・サーバが判断することを可能にし、

前記スレーブ・サーバが更新されることを必要とするかどうか、及び該スレーブ・サーバが前記情報のパケットを処理できるかどうかを示す信号を、該スレーブ・サーバから前記マスター・サーバに送信し、

前記スレーブ・サーバが前記情報のパケットをコミットすべきであるかどうかを示す応答信号を、前記マスター・サーバから該スレーブ・サーバに送信し、 50

前記応答信号により前記スレーブ・サーバが前記情報のパケットをコミットすべきであることが示された場合には、該スレーブ・サーバへの該情報のパケットをコミットすることによって、二フェーズ法で達成されるように定められた複製情報を送信する、段階を含むことを特徴とする方法。

【請求項 3 2】

ネットワークを通じてデータを複製する方法であって、

(a) 前記複製を、一フェーズ法で達成すべきか、又は二フェーズ法で達成すべきかを判断し、

(b) 前記データの現在の状態についてのバージョン番号を、マスター・サーバからスレーブ・サーバに送信し、

前記スレーブ・サーバ上のデータが前記バージョン番号に対応しない場合には、デルタが前記マスター・サーバから該スレーブ・サーバに送信されることを要求する、ことによって、一フェーズ法で複製されるべきデータを送信し、

(c) 情報のパケットを前記マスター・サーバからスレーブ・サーバに送信し、

前記スレーブ・サーバが前記情報のパケットを処理できるかどうかを判断し、

前記スレーブ・サーバが前記情報のパケットを処理できる場合には、該スレーブ・サーバへの前記情報のパケットをコミットする、ことによって、二フェーズ法で複製されるべきデータを送信する、段階を含むことを特徴とする方法。

【請求項 3 3】

ネットワーク上でマスターから複数のスレーブにデータを複製する方法であって、

(a) 前記複製を、一フェーズ法で達成すべきか、又は二フェーズ法で達成すべきかを判断し、

(b) 前記データの現在の状態についてのバージョン番号を前記マスターから各々のスレーブに送信し、

前記バージョン番号に対応しないデータを含むデルタが、前記マスターから各々のスレーブに送信されることを要求する、

ことによって、一フェーズ法で複製されるべきデータを送信し、

(c) 情報のパケットを前記マスターから各々のスレーブに送信し、

前記複数のスレーブの各々が前記情報のパケットを処理できる場合には、該スレーブへの該情報のパケットをコミットする、

ことによって、二フェーズ法で複製されるべきデータを送信する、段階を含むことを特徴とする方法。

【請求項 3 4】

一フェーズ法又は二フェーズ法を用いて、ネットワーク上でマスターから複数のスレーブにデータを複製する方法であって、

(a) 前記スレーブ上の前記データを更新するために、デルタが前記マスターから該スレーブに送信されることを、各々のスレーブが要求できるように、該データの現在の状態についてのバージョン番号を該マスターから各々のスレーブに送信することによって、一フェーズ法で複製されるべきデータを送信し、

(b) あらゆるスレーブが情報のパケットをコミットできる場合、各々のスレーブによりコミットされるべき情報のパケットを前記マスターから該スレーブに送信することによって、二フェーズ法で複製されるべきデータを送信する、段階を含むことを特徴とする方法。

【請求項 3 5】

一フェーズ法又は二フェーズ法を用いて、各々がクラスター・マスター及び少なくとも1つのクラスター・スレーブを含む、クラスター化されたネットワーク上でデータを複製する方法であって、

(a) 他のクラスター・マスターがそれぞれデルタを要求できるように、前記データの

10

20

30

40

50

現在の状態についてのバージョン番号を第1のクラスター・マスターから他の全てのクラスター・マスターに送信することによって、一フェーズ法で複製されるべきデータを送信し、

(b) 前記他のクラスター・マスターが情報のパケットをコミットできる場合には、該他のクラスター・マスターによりコミットされるべき情報のパケットを、前記第1のクラスター・マスターから他のクラスター・マスターの各々に送信することによって、二フェーズ法で複製されるべきデータを送信する、

段階を含むことを特徴とする方法。

【請求項36】

一フェーズ法により、そのクラスター・マスターを有する前記クラスターにおいて、前記データをクラスター・マスターの各々からクラスター・スレーブの各々に送信する段階をさらに含むことを特徴とする請求項35に記載の方法。

【請求項37】

二フェーズ法により、そのクラスター・マスターを有する前記クラスターにおいて、前記データをクラスター・マスターの各々からクラスター・スレーブの各々に送信する段階をさらに含むことを特徴とする請求項10に記載の方法。

【請求項38】

(a) 前記マスター・サーバ上に格納されたデータの変更に關し、前記データの現在の状態についての現在のバージョン番号を含み、さらに、該データの前の変更及び前の変更の各々についてのバージョン番号に関する情報のパケットを、ネットワーク上でマスター・サーバから各々のスレーブ・サーバに送信する手段と、

(b) 前記スレーブ・サーバが前記現在のバージョン番号に対応するように更新されたかどうかを、各々のスレーブ・サーバが判断するのを可能にする手段と、

(c) 前記スレーブ・サーバが前の変更を逃した場合には、各々のスレーブが前記情報をコミットすることを可能にする手段と、

(d) 前記スレーブ・サーバが前記情報のパケットをコミットする前に、前の変更が前記マスター・サーバから該スレーブ・サーバに送信されることを、前記前の変更を逃したスレーブ・サーバの各々が要求することを可能にする手段と、

を備えることを特徴とするコンピュータ可読媒体。

【請求項39】

ネットワークを通じてデータを複製するための、サーバ・コンピュータにより実行するためのコンピュータ・プログラム製品であって、

(a) 前記マスター・サーバ上に格納されたデータの変更に關し、前記データの現在の状態についての現在のバージョン番号を含み、さらに、該データの前の変更及び前の変更の各々についてのバージョン番号に関する情報のパケットを、ネットワーク上でマスター・サーバから各々のスレーブ・サーバに送信するコンピュータ・コードと、

(b) 前記スレーブ・サーバが前記現在のバージョン番号に対応するように更新されたかどうかを、各々のスレーブ・サーバが判断するのを可能にするコンピュータ・コードと、

(c) 前記スレーブ・サーバが前の変更を逃した場合に、各々のスレーブ・サーバが前記情報をコミットするのを可能にするコンピュータ・コードと、

(d) 前記スレーブ・サーバが前記情報のパケットをコミットする前に、前の変更が前記マスターから該スレーブ・サーバに送信されるように、前記前の変更を逃したスレーブ・サーバの各々が要求することを可能にするコンピュータ・コードと、

を備えることを特徴とするコンピュータ・プログラム製品。

【請求項40】

ネットワークを通じてデータを複製するためのシステムであって、

(a) 前記マスター・サーバ上に格納されたデータの変更に關し、前記データの現在の状態についての現在のバージョン番号を含み、さらに、該データの前の変更及び前の変更の各々についてのバージョン番号に関する情報のパケットを、ネットワーク上でマスター

10

20

30

40

50

・サーバから各々のスレーブ・サーバに送信する手段と、

(b) 前記スレーブ・サーバが前記現在のバージョン番号に対応するように更新されたかどうかを、各々のスレーブ・サーバが判断するのを可能にする手段と、

(c) 前記スレーブ・サーバが前の変更を逃した場合に、各々のスレーブが前記情報をコミットすることを可能にする手段と、

(d) 前記スレーブ・サーバが前記情報のパケットをコミットする前に、前の変更が前記マスター・サーバから該スレーブ・サーバに送信されることを、前記前の変更を逃したスレーブ・サーバの各々が要求することを可能にする手段と、
を備えることを特徴とするシステム。

【請求項 4 1】

10

プロセッサと、

(a) 前記マスター・サーバ上に格納されたデータの変更に、前記データの現在の状態についての現在のバージョン番号を含み、さらに、該データの前の変更及び前の変更の各々についてのバージョン番号に関する情報のパケットを、ネットワーク上でマスター・サーバから各々のスレーブ・サーバに送信し、

(b) 前記スレーブ・サーバが前記現在のバージョン番号に対応するように更新されたかどうかを、各々のスレーブ・サーバが判断するのを可能にし、

(c) 前記スレーブ・サーバが前の変更を逃した場合には、各々のスレーブが前記情報をコミットすることを可能にし、

(d) 前記スレーブ・サーバが前記情報のパケットをコミットする前に、前の変更が前記マスター・サーバから該スレーブ・サーバに送信されることを、前記前の変更を逃したスレーブ・サーバの各々が要求することを可能にする、
ように構成された、前記プロセッサにより実行されるオブジェクト・コードと、
からなることを特徴とするコンピュータ・システム。

20

【請求項 4 2】

ネットワークを通じてデータを複製するためのシステムであって、

a . i . スタート方法と呼び出すことによりデータ複製プロセスを起動するようになり、さらに、前記データのオリジナル・コピーに関する情報を送信するようになったマスター・ユーザ層と、

i i . 前記スタート方法を含み、前記マスター・ユーザ層からの前記コール及び前記データの前記オリジナル・コピーに関する前記情報を受信するようになり、さらに、該データの該オリジナル・コピーに関する該情報の少なくとも幾らかを含むデータ複製パケットを生成し送信するようになったマスター・サービス層と、
を備える、前記データのオリジナル・コピーを含むマスター・サーバと、

30

b . i . 前記マスター・サービス層から前記データ複製パケットを受信し、該データ複製パケットを処理するようになり、さらに、該データ複製パケットに関する情報を送信するようになったスレーブ・サービス層と、

i i . 前記データ複製パケットに関する前記情報を前記スレーブ・サービス層から受信するようになり、該情報を該データ複製パケット内に格納するようになったスレーブ・ユーザ層と、

40

を備える、前記マスター・サーバからの前記データのコピーを格納するようになったスレーブ・サーバと、
からなることを特徴とするシステム。

【請求項 4 3】

前記マスター・ユーザ層が、マスター・ユーザ及びマスター・ユーザ装置のうちの少なくとも1つと通信していることを特徴とする請求項 4 2 に記載のシステム。

【請求項 4 4】

前記マスター・ユーザ層が、デルタの形態の前記データのオリジナル・コピーに関する情報を送信するようになり、前記デルタが、該データのオリジナル・コピーの前の状態と現在の状態の間の変更に関する情報を含むようになったことを特徴とする請求項 4 2 に記

50

載のシステム。

【請求項 4 5】

前記マスター・ユーザ層が、前記データのオリジナル・コピーを更新するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 4 6】

前記マスター・ユーザ層が、前記データのオリジナル・コピーへの変更をスレーブ・サーバ上に複製すべきでないことを示すロールバック・メッセージを送信するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 4 7】

前記マスター・ユーザ層が、前記複製についてのタイムアウト値を設定するようになったことを特徴とする請求項 4 2 に記載のシステム。 10

【請求項 4 8】

前記マスター・ユーザ層が、前記データのオリジナル・コピーの現在の状態と該データのオリジナル・コピーの前の状態との間のデルタを生成するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 4 9】

前記マスター・ユーザ層が、前記データのオリジナル・コピーの現在の状態と該データの該オリジナル・コピーの前の状態との間のデルタを生成するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 5 0】

前記マスター・ユーザ層が、前記データのオリジナル・コピーの状態の各々について固有のバージョン番号を生成するようになったことを特徴とする請求項 4 2 に記載のシステム。 20

【請求項 5 1】

前記マスター・サービス層が、前記データ複製パケットをマルチキャストするようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 5 2】

前記マスター・サービス層が、前記データ複製パケットを一定間隔で送信するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 5 3】

前記マスター・サービス層が、前記データ複製パケットのバージョン番号を含むようになったことを特徴とする請求項 4 2 に記載のシステム。 30

【請求項 5 4】

前記マスター・サービス層が、前記スレーブ・サーバ上の該データのコピーを更新するのに必要な情報を前記データのオリジナル・コピーの現在の状態に含ませるようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 5 5】

さらに、前記マスター・サービス層が、デルタを含むデータ複製パケットを生成し送信するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 5 6】

さらに、前記マスター・サービス層が、前記データのオリジナル・コピーの連続する状態の間のデルタを含むデータ複製パケットを生成し送信するようになったことを特徴とする請求項 4 2 に記載のシステム。 40

【請求項 5 7】

さらに、前記マスター・サービス層が、前記データのオリジナル・コピーの任意の状態の間のデルタを含むデータ複製パケットを生成し送信するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 5 8】

前記マスター・サービス層が、前記マスター・ユーザ層からデルタを要求するようになったことを特徴とする請求項 4 2 に記載のシステム。 50

【請求項 59】

前記マスター・サービス層が、コミット・メッセージをスレーブ・サービス層に送信するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 60】

前記マスター・サービス層が、コミット・メッセージをスレーブ・サービス層に一定間隔で送信するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 61】

前記マスター・サービス層が、スレーブ・サービス層へのコミット・メッセージをマルチキャストするようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 62】

前記マスター・サービス層が、中止メッセージをスレーブ・サービス層に送信するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 63】

前記マスター・サービス層が、中止メッセージをスレーブ・サービス層に一定間隔で送信するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 64】

前記マスター・サービス層が、中止メッセージをスレーブ・サービス層にマルチキャストするようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 65】

前記スレーブ・ユーザ層が、スレーブ・ユーザ及びスレーブ・ユーザ装置のうちの少なくとも 1 つと通信していることを特徴とする請求項 42 に記載のシステム。 20

【請求項 66】

前記スレーブ・ユーザ層が、前記スレーブ・サーバ上に格納されたデータの現在のバージョン番号を確認するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 67】

前記スレーブ・ユーザ層が、前記スレーブ・サーバ上に格納されたデータへの前記データ複製パケットに関する情報をコミットするようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 68】

前記スレーブ・ユーザ層が、前記スレーブ・サーバ上に格納された前記データへの更新を中止するようになったことを特徴とする請求項 42 に記載のシステム。 30

【請求項 69】

前記スレーブ・ユーザ層が、前記データ複製パケット内に含まれる準備要求を処理するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 70】

前記スレーブ・ユーザ層が、前記データ複製パケット内に含まれる準備要求に関する応答を前記スレーブ・サービス層に送信するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 71】

前記スレーブ・ユーザ層が、データをローカル・ディスク上に永続的にキャッシュするようになったことを特徴とする請求項 42 に記載のシステム。 40

【請求項 72】

前記スレーブ・ユーザ層が、前記スレーブ・サーバ上の前記データのコピーの前記バージョン番号を更新するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 73】

前記スレーブ・サービス層が、前記マスター・サービス層からデルタを要求するようになったことを特徴とする請求項 42 に記載のシステム。

【請求項 74】

前記スレーブ・サービス層が、前記スレーブ・サービス層から、前記スレーブ・サーバ上に格納された前記データの前記現在のバージョン番号を要求するようになったことを特 50

徴とする請求項 4 2 に記載のシステム。

【請求項 7 5】

前記スレーブ・サービス層が、コミット・メッセージを前記スレーブ・ユーザ層に送信するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 7 6】

前記スレーブ・サービス層が、中止メッセージを前記スレーブ・ユーザ層に送信するようになったことを特徴とする請求項 4 2 に記載のシステム。

【請求項 7 7】

マスター・サーバからスレーブ・サーバにデータを複製する方法であって、

マスター・サーバ上のマスター・データの現在の状態に関する情報を含むスタートコールを、該マスター・サーバ上でマスター・ユーザ・レベルからマスター・サービス・レベルに送信し、

a . 前記スレーブ・サーバ上のスレーブ・データが現在の状態を有するかどうかを判断するために、該スレーブ・サーバ上のスレーブ・ユーザ層を確認するようになった、該スレーブ・サーバ上のスレーブ・サービス層に前記情報を送信し、

b . 前記スレーブ・サービス層からのデルタについての要求を、前記マスター・ユーザ層からのデルタを要求し受信するようになったマスター・サービス層に送信し、

c . 前記スレーブ・データを現在の状態に持って行くのに必要な前記情報を含むデルタを、前記マスター・サービス層から、前記デルタを処理し該情報を前記スレーブ・ユーザ層に送信するようになった前記スレーブ・サービス層に送信し、

d . 前記スレーブ・ユーザ層を用いて前記スレーブ・データを更新する、
段階を含むことを特徴とする方法。

【請求項 7 8】

前記マスター・ユーザ層を用いて前記データの現在の状態についてのバージョン番号を判断する段階をさらに含むことを特徴とする請求項 7 7 に記載の方法。

【請求項 7 9】

マルチキャストにより、前記情報を前記スレーブ・サービス層に送信する段階をさらに含むことを特徴とする請求項 7 7 に記載の方法。

【請求項 8 0】

前記マスター・データの現在の状態についてのバージョン番号を含む情報を、前記スレーブ・サービス層に送信する段階をさらに含むことを特徴とする請求項 7 7 に記載の方法。

【請求項 8 1】

マスター・サーバからスレーブ・サーバにデータを複製する方法であって、

a . 前記マスター・サーバ上に格納されたマスター・データにおける前の状態から現在の状態への変更に関する情報を含む新しいデルタを、マスター・サーバ上でマスター・ユーザ・レベルからマスター・サービス・レベルに送信し、

b . 前記スレーブ・サーバ上の前記スレーブ・データが現在の状態を有するかどうかを判断するために、該スレーブ・サーバ上のスレーブ・ユーザ層を確認するようになった該スレーブ・サーバ上のスレーブ・サービス層に、前記マスター・サービス層から前記新しいデータを送信し、

c . 前記スレーブ・データを前記マスター・データの前の状態に更新するのに必要な情報を含む同期デルタへの要求を、前記スレーブ・サービス層から、前記マスター・ユーザ層から同期デルタを要求し受信するようになった該マスター・サービス層に送信し、

d . 前記同期デルタを、前記マスター・サービス層から、前記デルタを処理し、前記情報を前記スレーブ・データにコミットされるべき前記スレーブ・ユーザ層に送信するようになった前記スレーブ・サービス層に送信し、

e . 前記スレーブ・ユーザ層を用いて、前記スレーブ・データへの前記新しいデルタの前記情報をコミットする、
段階を含むことを特徴とする方法。

10

20

30

40

50

【請求項 8 2】

ネットワークを通じてデータをマスター・サーバからスレーブ・サーバに複製する方法であって、

a．前記マスター・サーバ上の前記データのオリジナル・コピーの現在の状態に関するバージョン番号を、マスター・サービス層から、スレーブ・サービス層に送信し、

b．前記スレーブ・サーバ上のデータが前記バージョン番号に対応するように更新されたかどうかを、スレーブ・ユーザ層が判断することを可能にし、

c．前記スレーブ・サーバ上のデータが前記バージョン番号に対応しない場合には、デルタが前記マスター・サービス層から前記スレーブ・サービス層に送信されることを要求する、

段階を含むことを特徴とする方法。

10

【請求項 8 3】

前記スレーブ・ユーザ層が、各々のスレーブ・サーバについてローカル・ディスク上の前記データを永続的にキャッシュすることを可能にする段階をさらに含むことを特徴とする請求項 7 7 に記載の方法。

【請求項 8 4】

前記マスター・ユーザ層が、前記マスター・サーバ上のデータの現在の状態についての固有のバージョン番号を判断することを可能にする段階をさらに含むことを特徴とする請求項 7 7 に記載の方法。

【請求項 8 5】

スレーブ・ユーザ層がスレーブ・サーバ上の前記データを更新するのに必要な前記バージョン番号を有するデータを含ませる段階をさらに含むことを特徴とする請求項 7 7 に記載の方法。

20

【請求項 8 6】

前記スレーブ・ユーザ層により受信されるとすぐに前記スレーブ・サーバを更新するのに必要な前記データをコミットする段階をさらに含むことを特徴とする請求項 7 7 に記載の方法。

【請求項 8 7】

マスター・サーバ及び少なくとも 1 つのスレーブ・サーバを含む、ネットワークを通じてデータを複製する方法であって、

30

a．前記マスター・サーバ上に格納された前記データの変更に關し、前の状態についての前のバージョン番号及び該データの新しい状態についての新しいバージョン番号を含み、さらに、該データの前の変更及び前の変更の各々について前のバージョン番号に関する情報のパケットを、ネットワークでマスター・サービス層から各々のスレーブ・サーバのスレーブ・サービス層に送信し、

b．前記スレーブ・サーバ上のデータが前記パケット内に含まれる前記前のバージョン番号に対応するかどうかを、各々のスレーブ・サーバのスレーブ・ユーザ層が判断することを可能にし、

c．前記スレーブ・サーバ上のデータが前記パケット内に含まれる前記前のバージョン番号に対応する場合には、各々のスレーブ・ユーザ層が前記情報のパケットをコミットすることを可能にし、前記コミットはまた、前記スレーブ・サーバの前記バージョンを前記新しいバージョン番号に更新し、

40

d．前記スレーブ・サービス層が前記情報のパケットをコミットする前に、前記スレーブを前記前のバージョン番号に更新するのに必要な前記情報を含むデルタが、前記マスター・サービス層から、そのスレーブ・ユーザ層に対応する前記スレーブ・サービス層に送信されることを、該前のバージョン番号に対応しないスレーブ・ユーザ層の各々が要求することを可能にする、

段階を含むことを特徴とする方法。

【請求項 8 8】

ネットワークを通じて、マスター・サーバから少なくとも 1 つのスレーブ・サーバにデ

50

ータを複製する方法であって、

a. 前記マスター・サーバ上に格納された前記データの変更に關し、該データの現在の状態についてのバージョン番号を含む情報のパケットを、前記マスター・サーバのマスター・サービス層から、該スレーブ・サーバの前記ユーザ・サービス層に送信し、

b. 前記スレーブ・サーバが前記パケット内に含まれる前記バージョン番号に対応するように更新されたかどうか、さらに、該パケット内に含まれる該バージョン番号に対応するように更新することが必要な場合には、前記スレーブ・ユーザ層が前記情報のパケットを処理できるかどうかを、前記サーバの前記スレーブ・ユーザ層が判断することを可能にし、

c. 前記スレーブ・サーバが更新されることを必要とするかどうか、及び該スレーブ・サーバが前記更新を処理できるかどうかを示す信号を、前記スレーブ・サービス層から前記マスター・サービス層に送信し、

d. 前記スレーブ・ユーザ層が前記パケット内に含まれる前記情報にコミットすべきかどうかを示す応答信号を、前記マスター・サービス層から前記スレーブ・サービス層に送信し、

e. 前記応答信号により前記スレーブ・ユーザ層が前記パケット内に含まれる前記情報にコミットすべきであることが示された場合には、前記スレーブ・サーバへの前記情報のパケットをコミットする、

段階を含むことを特徴とする方法。

10

20

【請求項 89】

a. 前記マスター・サーバ上の前記データのオリジナル・コピーの現在状態に関するバージョン番号を、マスター・サービス層からスレーブ・サービス層に送信する手段と、

b. 前記スレーブ・サーバ上のデータが前記バージョン番号に対応するように更新されたかどうかを、スレーブ・ユーザ層が判断することを可能にする手段と、

c. 前記スレーブ・サーバ上の前記データが前記バージョン番号に対応しない場合には、デルタが前記マスター・サービス層から前記スレーブ・サービス層に送信されることを要求する手段と、

を備えることを特徴とするコンピュータ可読媒体。

30

40

【請求項 90】

ネットワークを通じてマスター・サーバからスレーブ・サーバにデータを複製する、サーバ・コンピュータにより実行するためのコンピュータ・プログラム製品であって、

a. 前記マスター・サーバ上の前記データのオリジナル・コピーの現在の状態に関するバージョン番号を、マスター・サービス層からスレーブ・サービス層に送信するコンピュータ・コードと、

b. 前記スレーブ・サーバ上の前記データが前記バージョン番号に対応するように更新されたかどうかを、スレーブ・ユーザ層が判断することを可能にするコンピュータ・コードと、

c. 前記スレーブ・サーバ上の前記データが前記バージョン番号に対応しない場合には、デルタが前記マスター・サービス層から前記スレーブ・サービス層に送信されることを要求するコンピュータ・コードと、

を備えることを特徴とするコンピュータ・プログラム製品。

50

【請求項 91】

ネットワークを通じてデータを複製するためのシステムであって、

a. 前記マスター・サーバ上の前記データのオリジナル・コピーの現在の状態に関するバージョン番号を、マスター・サービス層からスレーブ・サービス層に送信する手段と、

b. 前記スレーブ・サーバ上の前記データが前記バージョン番号に対応するように更新されたかどうかを、スレーブ・ユーザ層が判断することを可能にする手段と、

c. 前記スレーブ・サーバ上の前記データが前記バージョン番号に対応しない場合には、デルタが前記マスター・サービス層から前記スレーブ・サービス層に送信されることを要求する手段と、

50

を備えることを特徴とするシステム。

【請求項 9 2】

プロセッサと、

(a) 前記マスター・サーバ上に格納された前記データのオリジナル・コピーの現在の状態に関するバージョン番号を、マスター・サービス層からスレーブ・サービス層に送信し、

(b) 前記スレーブ・サーバ上の前記データが前記バージョン番号に対応するように更新されたかどうかを、スレーブ・ユーザ層が判断することを可能にし、

(c) 前記スレーブ・サーバ上の前記データが前記バージョン番号に対応しない場合には、デルタが前記マスター・サービス層から前記スレーブ・サービス層に送信されることを要求する、

ように構成された、前記プロセッサにより実行されるオブジェクト・コードと、
からなることを特徴とするコンピュータ・システム。

【発明の詳細な説明】

【発明の詳細な説明】

【0001】

(優先権主張)

本出願は、2001年7月16日に出願された名称「データ複製プロトコル」の米国特許仮出願第60/305,986号、2001年7月16日に出願された名称「データ複製のための階層化アーキテクチャ」の米国特許仮出願第60/305,978号、2001年10月11日に出願された名称「データ複製プロトコル」の米国特許出願第09/975,590号、及び2001年10月11日に出願された名称「データ複製のための階層化アーキテクチャ」の米国特許出願第09/975,587号に基づく優先権を主張するものであり、これらは、出典の明示によりここで開示されたものとする。

【0002】

(著作権の告知)

本特許文献の開示の一部は、著作権の保護を受ける内容を含んでいる。著作権者は、本特許文献が米国特許商標庁の特許ファイル又は記録において公開されたときに、何者かが特許文献又は特許の開示を複製することに異議はないが、他の場合には、如何なる場合にも全ての著作権を留保する。

【0003】

(技術分野)

本発明は、一般に、データを転送するためのシステムに関する。より具体的には、ネットワークを通じてデータを複製するためのシステム及び方法に関する。

【0004】

(背景技術)

分散処理システムには幾つかのタイプがある。分散処理システムは、一般に、通信媒体を通じて連結された2つのコンピュータのような複数の処理デバイスを含んでいる。あるタイプの分散処理システムは、クライアント/サーバ型ネットワークである。クライアント/サーバ型ネットワークは、少なくとも2つのデバイス、典型的には中央サーバ及びクライアントを含む。付加的なクライアントを中央サーバに連結することができ、多数のサーバがあってもよく、或いはネットワークが通信媒体を通じて連結された唯一のサーバを含んでもよい。

【0005】

このようなネットワーク環境においては、アプリケーション又は情報を中央サーバから多数のワークステーション及び/又は他のサーバに送信することがしばしば望ましい。多くの場合、各々のワークステーションへの別々にインストールを行い、或いは、中央サーバから個々のワークステーションの各々及び/又はサーバに情報の新しいライブラリを別々に押し出す。これらの手法は時間がかかり、リソースの使用は非効率となる。各々のワークステーション又はサーバについてアプリケーションを別個にインストールすることが

10

20

30

40

50

、付加的かつ潜在的なエラーの源となりうる可能性もある。

【0006】

理想的には、情報の送信は、障害に対する信頼性が高く、規模拡大が可能であるべきであり、それによってプロセスがネットワークを有効に使えるようになる。従来のソリューションでは、概して、これらの目的の一方又は両方とも達成することができない。1つの簡単な手法は、各々のスレーブに個別に連絡を取り、TCP/IP接続のような2地点間のリンクを通じてデータを転送するマスター・サーバをもつことである。1つ又はそれ以上のスレーブが一時的に到達不能である場合、或いはスレーブが更新を処理する際にエラーが生じる場合、この手法によると、一貫性のないデータのコピーが行われることになる。これと対照的なのは複雑な分散同意プロトコルであり、これは、全てのデータのコピーが一貫していることを保証するように、スレーブの間に多くのクロストークを必要とする。

10

【0007】

(発明の開示)

本発明は、ネットワーク上で達成することができるように、マスター・サーバから少なくとも1つのスレーブ又は管理対象サーバにデータを複製する方法を含む。この方法においては、複製を一フェーズ法で達成すべきか又は二フェーズ法で達成すべきかを判断することができる。複製が一フェーズ法で達成すべき場合には、マスター・サーバ上のデータの現在の状態に対応するバージョン番号を送信することができる。このバージョン番号は、ネットワーク上のあらゆるスレーブ・サーバ、又はスレーブ・サーバのサブセットにだけ送信することができる。次に、バージョン番号を受信するスレーブ・サーバは、デルタがマスターから送信されることを要求することができる。デルタは、現在のバージョン番号に対応するようにそのスレーブ上のデータを更新することが必要とされるデータを含むことができる。

20

【0008】

複製が、二フェーズ法で達成されるべき場合には、マスターから各々のスレーブ、又はスレーブのサブセットに、情報のパケットを送信することができる。次に、これらのスレーブは、それらが情報のパケットをコミットできるかどうかを、マスター・サーバに回答することができる。スレーブの少なくとも幾つかがデータをコミットできる場合には、マスターは、コミットを処理すべきであるという信号をこれらのスレーブに送信することができる。コミットを処理した後、これらのスレーブは、現在のバージョン番号に更新することができる。いずれかのスレーブがコミットを処理できない場合には、コミットを中止することができる。

30

【0009】

(発明を実施するための最良の形態)

本発明は、例えば、マスター・サーバ又は「管理」サーバ(「Adminサーバ」)からスレーブ・サーバの集合又は「管理対象」サーバに、データ又は他の情報の複製を提供するものである。この複製は、従来のローカル・エリア・ネットワーク又はイーサネット(登録商標)のような何らかの適切なネットワークにわたって生じることができる。1つの実施形態において、マスター・サーバは、ネットワークの全データのオリジナルの記録を所有し、このオリジナルの記録に何らかの更新が加えられることになる。更新が生じたとき、該更新と共にデータのコピーを、各々のスレーブ・サーバに伝送することができる。1つの適用例は、Adminサーバから管理対象サーバの集合に設定情報を分散することを含む。

40

本発明による1つのシステムにおいては、データ複製サービス(DRS)のようなサービスが、設定及び配置情報をAdminサーバから適当なドメイン内の管理対象サーバに分散させることが必要である。ユーザ・データグラム・プロトコル(「UDP」)のようなマルチキャスト・プロトコルはフロー制御をもたず、システムを制圧することができるので、伝送制御プロトコル(「TCP」)のような2地点間接続によって大きなデータ項目を分散させることができる。2地点間接続には、リモート・メソッド呼び出し(RMI

50

）、ハイパーテキスト転送プロトコル（ＨＴＴＰ）、又は同様のプロトコルを用いることができる。

【００１０】

管理対象サーバは、ローカル・ディスク上のデータを永続的にキャッシュすることでもできる。そのようなキャッシュを用いない場合には、必要なデータの転送に、許容できないほどの量の時間を必要とすることがある。転送されるべき起動データの量を減少させることによって、起動速度が増すので、管理対象サーバのキャッシュする能力は重要である。Adminサーバが到達不能の場合には、キャッシングは、起動及び／又は再起動を可能にもする。再起動はより魅力的な選択肢であり、Adminサーバがサーバに起動するよう命令する場合がそうである。しかしながら、キャッシングは、利用可能なAdminサーバを用いずに、ドメインを開始する能力を提供することができる。

図１のドメイン構造１００に示されるように、１つのAdminサーバ１０２と少なくとも１つの管理対象サーバ１０４が、ドメイン１０６を含むことができる。このドメイン１０６は、起動及びシャットダウンのための管理ユニットとすることができる。１つの実施形態において、ブラウザ１０８又は他のユーザ・アプリケーション又はデバイスが、Adminサーバ１０２に起動するよう伝える。次に、Adminサーバ１０２は、ドメイン１０６内の全ての管理対象サーバ１０４に起動するよう伝え、適切な設定情報を渡す。管理対象サーバ１０４が起動した後でサーバが停止した場合には、Adminサーバ１０２が利用可能であろうとなかろうと、そのサーバが自動的に再起動することが望ましい。キャッシュされたデータは、この目的のために有用なものである。

【００１１】

Adminサーバ上のデータへの更新は、バージョン間の段階的デルタとしてパッケージすることができる。デルタは、変更されるべき設定及び／又は他の情報を含むことができる。システムをオフラインにするのは望ましくないので、ドメインが動作している間に設定を更新することが好ましい。１つの実施形態において、設定の変更は、それらがAdminサーバにより押し出された動的に生じる。その都度全部の設定を送信するのは不必要であり、煩わしすぎるので、設定の変更だけがデルタで送信される。

本発明によるプロトコルは、これに応じて他の適切な方法を用いることもできるが、更新の分散のための２つの方法を統合する。これらの分散方法は、一フェーズ法及び二フェーズ法と呼ぶことができ、一貫性と規模拡大可能性の間のトレードオフを提供することができる。規模拡大可能性を好む一フェーズ法において、各々のスレーブは、独自のペースで更新を獲得し処理することができる。スレーブは、異なる時間にマスターから更新を得ることができるが、更新が受信されるとすぐにデータをコミットすることができる。スレーブは、更新を処理する際にエラーに遭遇することがあるが、一フェーズ法においては、このことが、他のスレーブの更新処理を防止することはない。

【００１２】

一貫性を好む、本発明による二フェーズ法において、全てのスレーブがデータをうまく処理できるか、いずれのスレーブもデータをうまく処理できないかという点で、分散を「原子」にすることができる。中止の可能性を許容する準備段階及びコミット段階といった別個の段階がある。準備段階において、マスターは、各々のスレーブが更新を取得できるかどうかを判断することができる。スレーブが更新を受容できることを全てのスレーブが示す場合には、コミット段階において、コミットされるべきスレーブに新しいデータを送信することができる。スレーブ・サーバのうちの少なくとも１つが更新を取得できない場合には、更新を中止することができ、コミットはなされない。この場合には、管理対象サーバは、該管理対象サーバが準備をロールバックすべきであり、何も変更されないことを通知することができる。いずれの方法においても、更新がコミットされた時に到達不能なスレーブが最終的に更新を得るので、本発明によるこうしたプロトコルは信頼性があるものである。

【００１３】

本発明によるシステムはまた、一時的に利用できないサーバが最終的に全ての更新を受

10

20

30

40

50

信することを保証することもできる。例えば、サーバがネットワークから一時的に隔離され、その後再起動しないでネットワークに復帰する。サーバが再起動しないので、該サーバは、通常、更新を確認しない。サーバに新しい更新を定期的を確認させることによって、又はサーバが更新を受信したかどうかを調べるためにマスター・サーバに定期的を確認させることによって、ネットワークに復帰するサーバについて説明することができる。

1つの実施形態においては、マスター・サーバが、マルチキャスト「ハートビート」をスレーブ・サーバに規則的に送信する。マルチキャスト手法は信頼できないので、スレーブが任意のシーケンスのハートビートを逃すこともある。例えば、ネットワーク分割化のために、スレーブ・サーバがネットワークから一時的に切断されるか、或いはハートビートを逃させることで、スレーブ・サーバ自体が一時的にネットワークに利用できなくなることがある。したがって、ハートビートは、最近の更新についての情報のウィンドウを含むことができる。以下に説明されるように、前の更新についてのこうした情報を用いて、ネットワーク・トラフィックの量を減少させることができる。

10

【0014】

各々のマスター及び各々のスレーブ内には、少なくとも2つの層、すなわち、ユーザ層及びシステム層（すなわちDRS層）がある。ユーザ層は、データ複製システムのユーザに対応することができる。DRS層は、データ複製システム自体の実施に対応することができる。これらの参加者及び層の対話が図2に示される。

図2の起動図200に示されるように、この実施形態におけるマスター・ユーザ層202及びスレーブ・ユーザ層204は、それぞれ、マスターDRS層206及びスレーブDRS層208内にダウンコールを行う。こうしたダウンコールは、例えば、
`registerMaster(DID, verNum, listener)`
`registerSlave(DID, verNum, listener)`
 の形態をとることができ、ここで、DIDは、公知のDIDの知識から得られる識別子であり、関心あるオブジェクトを指し、verNumは、ユーザの現在のバージョン番号としてローカルの永続的記憶装置から得られ、listenerは、DRS層からのアップコールを取り扱うオブジェクトである。アップコールは、リスナ・オブジェクトについての方法を呼び出すことができる。その後、マスターは、現在のバージョン番号と共に、ハートビート又は定期的なデルタを送り始めることができる。スレーブ・ユーザ204から情報を得るようになったコンテナを含むことができるコンテナ層210が示される。可能なコンテナの例は、エンタープライズ・ジャバ・ビーンズ、ウェブ・インターフェース、及びJ2EE（ジャバ2プラットフォーム、エンタープライズ・エディション）アプリケーションを含む。他のアプリケーション及び/又はコンポーネントは、管理クライアント212のようなコンテナ層210にプラグインすることができる。ユーザ層とDRS層の間を通信するメッセージの例は、一フェーズ法については図4に、二フェーズ法については図5に示されている。

20

30

【0015】

図4は、本発明による階層化アーキテクチャにおける一フェーズ分散手法のために用いることのできる1つの基本プロセス400を示す。このプロセスにおいては、マスター・ユーザ層402は、マスターDRS層406内にダウンコール404を行い、一フェーズ分散を開始する。このコールは、システム内の全てのスレーブに対するものとしてもよく、スレーブ・サーバのサブセットだけに対するものでもよい。このコールがサブセットに対するものである場合、マスター・ユーザ層402は、更新の範囲、或いはどのスレーブが更新を受信すべきであるかを判断することができる。

40

【0016】

マスターDRS層は、マスター上のデータの現在のバージョン番号を含むハートビート408をスレーブDRS層にマルチキャストし始める410。スレーブDRS層410は、スレーブ・ユーザ層414からそのスレーブについての現在のバージョン番号を要求する412。次に、スレーブ・ユーザ層414は、そのスレーブのバージョン番号を有するスレーブDRS層410に応答する416。スレーブが同期しているか、又は既に現在の

50

バージョン番号をもっている場合には、次の更新までそれ以上の要求はなされない。スレーブが同期しておらず、該スレーブが更新の範囲内にある場合には、スレーブ D R S 層 4 1 0 は、該スレーブをマスター上のデータの現在のバージョン番号に更新するために、マスター D R S 層 4 0 6 からデルタを要求することができる 4 2 0。マスター D R S 層 4 0 6 は、マスター・ユーザ層 4 0 2 が、スレーブを更新するためにデルタを生成することを要求する 4 2 2。次に、マスター・ユーザ層 4 0 2 は、デルタをマスター D R S 層 4 0 6 に送信 4 2 4 し、該マスター D R S 層は、該デルタ及びマスターの現在のバージョン番号をスレーブ D R S 層 4 1 0 に転送し 4 2 6、該スレーブ D R S 層は、デルタをコミットされるべきスレーブ・ユーザにデルタを送信する 4 2 8。ハートビート 4 0 8 がスレーブにより受信されてからマスターが更新された場合は、現在のバージョン番号がデルタと共に送信される。 10

【 0 0 1 7 】

マスター D R S 層 4 0 6 は、バージョン番号を含むマルチキャスト・ハートビートを、スレーブ・サーバに定期的送信し続けることができる 4 0 6。このことは、利用できなかった、又はデルタを受信し処理することができなかった如何なるサーバも、該サーバがデータの現在のバージョン番号を有していないことを判断し、後にスレーブがシステムに復帰した時などにデルタを要求することを可能にする 4 2 0。

図 5 は、本発明による階層化アーキテクチャにおける二フェーズ分散手法のために用いることのできる 1 つの基本プロセス 5 0 0 を示す。このプロセスにおいては、マスター・ユーザ層 5 0 4 は、マスター D R S 層 5 0 6 内にダウンコールを行い、二フェーズ分散を開始する 5 0 4。マスター・ユーザ層 5 0 2 は、再び更新の範囲を判断することを必要とし、更新処理のための「タイムアウト」値を設定することができる。 20

【 0 0 1 8 】

マスター D R S 層 5 0 6 は、新しいデルタをスレーブ D R S 層 5 1 0 に送信する 5 0 8。スレーブ D R S 層 5 1 0 は、新しいデルタについて、準備要求をスレーブ・ユーザ層 5 1 4 に送信する 5 1 2。次に、スレーブ・ユーザ層 5 1 4 は、スレーブが新しいデルタを処理できようとできまいとスレーブ D R S 層 5 1 0 に応答する 5 1 6。スレーブ D R S 層は、応答をマスター D R S 層 5 0 6 に転送する 5 1 8。スレーブが同期していないために要求を処理できない場合には、マスター D R S 層 5 0 6 は、マスター・ユーザ層 5 0 2 にアップコールを行い、スレーブを同期させるデルタを生成し、該デルタをコミットする 5 2 0。マスター・ユーザ層 5 0 2 は、同期デルタをマスター D R S 層に送信し 5 2 2、該マスター D R S 層は、同期デルタをスレーブ D R S 層 5 1 0 に転送する 5 2 4。スレーブが同期デルタを処理できる場合には、スレーブ D R S 層 5 1 0 は、同期応答をマスター D R S 層 5 0 6 に送信し 5 2 6、ここで該スレーブが新しいデルタを処理できるようになる。スレーブが同期デルタを処理できない場合には、スレーブ D R S 層 5 1 0 は、適切な同期応答をマスター D R S 層 5 0 6 に送信する 5 2 6。次に、マスター D R S 層 5 0 6 は、スレーブが新しいデルタを処理できると該スレーブ・サーバが応答したかどうかによって、コミット・メッセージ又は中止メッセージをスレーブ D R S 層 5 1 0 に一定間隔で送信する 5 2 8。例えば、全てのスレーブがデルタを準備できる場合には、マスターは、コミット信号を一定間隔で送ることができる。その他の場合には、マスターは、コミット信号を一定間隔で送信することができる。ハートビートは、更新の範囲も含むので、スレーブは、該ハートビート内に含まれる情報を処理すべきかどうかを知るようになる。 30 40

【 0 0 1 9 】

スレーブ D R S 層は、このコマンドをスレーブ・ユーザ層 5 1 4 に転送し 5 3 0、次いで、該スレーブ・ユーザ層は、新しいデルタの更新をコミットするか又は中止する。準備段階がマスター・ユーザ層 5 0 2 によって設定されたタイムアウト値内に完了しなかった場合には、マスター D R S 層 5 0 6 は、中止を全てのスレーブに自動的に一定間隔で送信することができる 5 2 8。このことは、例えば、マスター D R S 層 5 0 6 がスレーブのうちの少なくとも 1 つに連絡し、該スレーブがコミットを処理できるかどうか判断することかできない時に生じる。タイムアウト値は、マスター D R S 層 5 0 6 が、更新を中止する 50

前に指定の時間スレーブに連絡をとろうとすることができるように、設定することができる。

【0020】

一フェーズ法における更新の場合、これらのハートビートは、各々のスレーブに、該スレーブのデータの現在のバージョンから始まるデルタを要求させることができる。こうしたプロセスは、図6のフローチャートに示されている。本発明による階層化アーキテクチャを利用してもしなくともよいこの基本プロセス600において、マスター・サーバ上の現在のデータのバージョン番号は、マスター・サーバからスレーブ・サーバに送信される602。スレーブ・サーバは、該スレーブ・サーバが現在のバージョン番号に更新されたかどうかを判断する604。スレーブが現在のバージョンでない場合には、スレーブ・サーバを更新するために、デルタがマスター・サーバから送信されることを要求する606。デルタがスレーブ・サーバに送信された時には、スレーブ・サーバは、スレーブ・データを現在のバージョンに更新する608。次に、スレーブ・サーバは、バージョン番号を現在のバージョン番号に更新する610。

10

【0021】

二フェーズ法における更新の場合、マスターは、該マスターが直前のバージョンからのデルタを積極的に各々のスレーブに送信する準備段階で始めることができる。こうしたプロセスが、図7のフローチャートに示されている。本発明による階層化アーキテクチャを利用してもしなくともよいこの基本プロセス700においては、情報のパケットが、マスターからスレーブ・サーバに送信される702。パケットを受信する各々のスレーブ・サーバは、該スレーブ・サーバがそのパケットを処理し、現在のバージョンに更新できるかどうか判断する704。スレーブ・サーバがパケットを処理できるかどうかを示す、パケット応答を受信するスレーブ・サーバの各々は、マスター・サーバに応答する706。全てのスレーブ（これにデルタが送信される）が、あるタイムアウト時間内にデルタの処理に成功したことを確認した場合には、マスターは、更新をコミットするように決定することができる。他の場合には、マスター・サーバは、更新を中止するように決定することができる。一旦決定が下されると、マスター・サーバは、更新をコミットすべきか中止すべきかを示すメッセージをスレーブ・サーバに送信する708。決定がコミットすべきである場合には、各々のサーバは、コミットを処理する710。スレーブの中の1つによりコマンドを逃した場合には、ハートビートを更に用いて、コミットが生じたか、或いは中止が生じたかを知らせることができる。

20

30

【0022】

スレーブは、最初にマスターから現在のバージョン番号を得ることなく、キャッシュされたデータを用いて、直ちに起動及び/又は再起動するように構成することができる。上述したように、本発明による1つのプロトコルは、スレーブがローカル・ディスク上のデータを永続的にキャッシュすることを可能にする。このキャッシングは、転送されることが必要なデータの量を減少させることによって、システムの起動に必要とされる時間を減少させ、規模拡大可能性を改善するものである。このプロトコルは、マスターが到達不能な場合に、スレーブが起動及び/又は再起動することを可能にすることにより、信頼性を改善することができる。さらに、バージョン間の段階的デルタとして更新をパッケージすることを可能にすることができる。キャッシュ・データが存在しない場合には、スレーブは、マスターを待つか、又はデータ自体を取り出すことができる。スレーブがキャッシュを有する場合には、該スレーブは、依然として同期しないで起動することを望まないであろう。スレーブが待つことを承知している場合には、起動時間を減少させることができる。

40

【0023】

このプロトコルは、状況に合わせて、マスター又はスレーブがデータ転送のイニシアチブを取ることができるという点で、双方向性とするすることができる。例えば、スレーブは、ドメインの起動中、マスターからデルタを取り出すことができる。スレーブが、デルタを更新しようと意図するものと異なるバージョンにあると判断した時には、該スレーブは、デルタがその現在のバージョンからの現在のシステム・バージョンに要求することができ

50

る。スレーブはまた、一フェーズ分散中にデルタを取り出すこともできる。ここで、システムは、ハートビートを読み取り、更新をとらえ損なったことを判断し、適切なデルタを要求することができる。

スレーブはまた、例外的な状況から回復するのに必要とされる時にデルタを取り出すこともできる。例外的な状況は、例えば、システムのコンポーネントが同期していない時に存在する。スレーブがデルタを取り出す時、該デルタは、データの任意のバージョン間にあることができる。換言すれば、どんなに多くの隔たった反復がそれらのバージョン間にあっても、デルタは、スレーブの現在のバージョンとシステム（又はドメイン）の現在のバージョンの間にあることができる。この実施形態においては、ハートビートの可用性、及びデルタを受信する能力が、システムの同期をもたらすことができる。

10

【0024】

スレーブがデルタを取り出す能力に加えて、マスターは、二フェーズ分散中、デルタをスレーブに押し出す能力を有することができる。1つの実施形態において、これらのデルタは、常に、データの連続するバージョンの間にある。この二フェーズ分散手法は、参加者の間に不一致が生じる可能性を最小にすることができる。スレーブ・ユーザは、更新をクライアントに公開することなく、又は更新をロールバック不能にすることなく、できる限り速く準備を処理することができる。このことは、競合についてサーバを確認するようなタスクを含むことができる。「ディスクが満杯である」又は「一貫性のない設定である」といったメッセージを送信することによって、スレーブのいずれかがエラーを知らせる場合には、更新を全体的にロールバックすることができる。

20

【0025】

しかしながら、非一貫性が生じる可能性が依然としてある。例えば、ソケットを開くことができないといった理由のために、コミットを処理する際にエラーが生じ得る。サーバはまた、様々な時間に更新をコミットし、公開することができる。データが、全く同時に全ての管理対象サーバに到達するはできないので、幾らかのリプル効果が生じ得る。このリプル効果を最小にしようとして、マルチキャストを使用することが小さな時間窓を提供することができる。1つの実施形態においては、コミットをとらえ損なった場合に準備されたスレーブが停止する、信号をとらえ損なった場合にマスターがクラッシュするなど。

マルチキャストに対するベストエフォート型手法は、スレーブ・サーバにコミット信号を逃させることがある。コミット段階中マスターが途中でクラッシュした場合、回復のためのロゲイン又は手段はないであろう。マスターが、残りのスレーブに、該マスターがコミットする必要があることを伝える方法はない。中止する際、バージョンが適切にロールバックされない場合には、幾つかのスレーブが、データをコミットすることで終わることができる。1つの実施形態において、残りのスレーブは、一フェーズ分散を用いて更新を得ることができる。このことは、例えば、管理対象サーバが、Adminサーバから受信したハートビートに応答してデルタを取り出す時に起こり得る。この手法はシステムの規模拡大可能性を維持し、如何なるコミット・エラー又はバージョン・エラーをも回避するためにシステムが分散を制限した場合に失われるであろう。

30

【0026】

システムにより管理される各々のデータ項目は、ドメイン全体にわたって公知である、固有の長寿命ドメイン識別子(DID)をもつように構造化することができる。データ項目は、各々がドメイン内のサーバの幾つかのサブセットに関連する多くのコンポーネントからなる、大きく複雑なオブジェクトとすることができる。これらのオブジェクトは、一貫性のあるユニットとすることができるので、幾つかの小さなオブジェクトより、少数の大きなオブジェクトをもつ方が望ましい。例として、単一のデータ項目又はオブジェクトは、config.xmlファイル又はアプリケーション-EARファイルといったコードファイルを含むシステムについての全ての設定情報を表すことができる。データ項目内の所定のコンポーネントは、スレッドの数に関して個々のサーバに関連することができ、配置されたサービスに関してクラスターに関連することができ、或いは安全証明書に関す

40

50

るドメイン全体に関連することもできる。2つのバージョンの間のデルタは、これらのコンポーネントのうちの幾つか又は全てについての新たな値から構成することができる。例えば、コンポーネントは、ドメインのメンバーに配置された全てのエンタープライズ・ジャバ・ビーンズを含むことができる。デルタは、これらのジャバ・ビーンズのサブセットだけに対する変更を含むことができる。

【0027】

デルタの「範囲」は、該デルタ内の関連するコンポーネントを有する全てのサーバの組を指すことができる。本発明によるAdminサーバは、デルタの範囲を判断するために、設定の変更を解釈することができる。マスター上のDRSシステムは、データを適切なスレーブに送信するために、その範囲を知る必要がある。マスターが、各々の更新においてサーバのサブセットに触れるだけでよい時、あらゆる設定の更新をあらゆるサーバに送信することは、時間とリソースの無駄であろう。

10

分散を制御するために、マスター・ユーザは、連続するバージョン間のデルタと共に各々の更新の範囲を提供することができる。この範囲は、ドメイン内の同じネームスペースからとることのできる、サーバ及び/又はクラスターに関する名前の組として表すことができる。1つの実施形態においては、DRSは、名前をアドレスにマッピングするために、リゾルバ・モジュールを用いる。クラスター名は、そのクラスター内の全てのサーバのアドレスの組にマッピングすることができる。これらのアドレスは、仮想マシンに関するように、相対的なものとすることができる。当該技術分野において周知であり使われているように、リゾルバは、介在するファイアウォールがあるかどうかを判断し、サーバが「ファイアウォールの内側」にあるかどうかに関する、「内側」アドレス又は「外側」アドレスのいずれかを返すことができる。Adminサーバ又は他のサーバは、設定データを用いて対応するリゾルバを初期化することができる。

20

【0028】

管理されるデータ項目の各々についての固有の長寿命ドメイン識別子(DID)と共に、データ項目の各々のバージョンは、長寿命のバージョン番号をもつこともできる。適切なバージョンに関する混乱のために、サーバが不適正に更新したり更新に失敗したりしないように、各々のバージョン番号を更新の試みに固有のものとすることができる。同様に、中止された二フェーズ分散のバージョン番号を再使用することはできない。マスターは、バージョン番号だけが与えられた任意の2つのバージョン間のデルタを生成することができる。マスターが、こうしたデルタを生成できない場合には、データ又はアプリケーションの完全なコピーを提供することができる。

30

【0029】

データ複製サービスは、できる限り一般的なものにしておくことが望ましい。したがって、このシステムのユーザに少数の仮定を課すことができる。このシステムは、例えば、

- ・このシステムは、バージョン番号を増加させる方法を含むことができる、
- ・このシステムは、バージョン番号をマスターにもスレーブにも永続的に格納することができる、

- ・このシステムは、バージョン番号を比較し、等しいかどうかを判断する方法を含むことができる、

40

という3つの主な仮定に依存することができる。これらの仮定は、インターフェース「VersionNumber」のような、DRSインターフェースのユーザレベルの実施によって提供することができる。こうしたインターフェースは、ユーザが、バージョン番号抽象化の特定の概念及び実施を提供することを可能にできる一方、システムがそのバージョン番号の属性にアクセスを有することを保証する。例えば、ジャバでは、VersionNumberインターフェースを、以下のように実行することができる。

```

package weblogic.drs;
public interface VersionNumber extends Serializable {
    VersionNumber increment();
    void persist() throws Exception;
    boolean equals (VersionNumber anotherVN);
    boolean strictlyGreaterThan(VersionNumber anotherVN);
}

```

ユーザがシステムに提供できるこの抽象化の単純化した実施は、大きな正の整数となるであろう。この実施はまた、当該技術分野において「直列化可能」と呼ばれる、システムがネットワークを介してマスターからスレーブまでデルタ情報を転送できることをも保証することができる。

【0030】

上記の抽象化を用いる場合、ユーザレベルにおけるデルタの詳細な内容の概念から抽象化することが有用である。このシステムは、デルタ情報構造の知識を必要としなくてもよく、実際に、該構造を判断することさえできなくてもかまわない。デルタの実施を、直列化可能とすることもでき、システムがネットワークを介してマスターからスレーブまでデルタのバージョン情報を転送できることを保証する。

適切なDID及びバージョン番号と共に、各々のデータ項目についての記録のコピーをマスターに永続的に格納させることが望ましい。二フェーズ分散を始める前に、マスターは、提案された新しいバージョン番号を永続的に格納し、マスターが機能しなくなる場合に該新しいバージョン番号が再使用されないことを保証することができる。スレーブは、そのDID及びバージョン番号と共に、各々の関連データ項目の最新のコピーを永続的に格納することができる。スレーブがその都度データ又はプロトコルを取得しなければならないように、必要なキャッシングを行うように該スレーブを構成することもできる。このことは、あらゆる場合について望ましいわけではないが、起こり得る特定の状況を取り扱うために許容することができる。

【0031】

さらに、本発明によるシステムは、同時制限を含むことができる。例えば、所定の範囲にわたる所定のDIDについての更新の二フェーズ分散中、特定の操作が許容されないようにすることができる。こうした操作は、空でない交差部分を有する範囲にわたって同じDID上の範囲のメンバーシップを修正するといった、一フェーズ法の更新又は二フェーズ法の更新を含むことができる。

少なくとも1つの実施形態において、マスターDRSは、ドメイン内の各々のサーバのスレーブDRSに対して、ハートビート又はパケット情報を規則的にマルチキャストする。各々のDIDについて、ハートビートは、各々の更新バージョン番号、前のバージョンに対するデルタの範囲、及びその更新がコミットされたか又は中止されたかを含む、最新の更新についての情報ウィンドウを含むことができる。現在のバージョンのについての情報をいつも含ませることができる。正確さ又は有効性のためではなく、マスターに戻るトラフィックの量を最小にするために、古いバージョンについての情報を用いることもできる。

【0032】

デルタ内に古いバージョン情報を含ませる場合、スレーブは、準備時に予期していた更新の一部分をコミットし、より新しい更新を取り扱うために新しいデルタを要求することができる。矢継ぎ早の更新がウィンドウを許容できない寸法まで増加させることがあるが、少なくとも幾つかの一定の設定可能な数のハートビートの間、所定のバージョンについての情報を含ませることができる。別の実施形態においては、一旦マスターが全てのスレーブが更新を受信したと判断すると、古いバージョンについての情報を廃棄することができる。

マルチキャスト・ハートビートは、考慮すべき幾つかの特性をもつ。これらのハートビートは、非同期式すなわち「一方通行」とすることができる。その結果、スレーブがハートビートに応答するまでには、マスターは、新たな状態に進んでいることがある。さらに、全てのスレーブが正確に同時に応答するとは限らない。したがって、マスターは、スレーブがその状態についての知識をもっていないと仮定し、デルタが更新することを意図するものを含むことができる。スレーブがハートビートの任意のシーケンスを逃すともあるので、これらのハートビートも、信頼性を欠くものである。このことは、ハートビート内に古いバージョン情報を再び含ませることにつながる。1つの実施形態においては、ハートビートは、それらが送信された順序でスレーブに受信される。例えば、スレーブは、該スレーブがバージョン6をコミットするまで、バージョン7をコミットすることができない。サーバは、該サーバが6を受信するまで待つか、或いは単純に6を廃棄し、7をコミットすることができる。この順序付けは、バージョンが逆戻りすることによりもたらされる混乱の可能性を排除することができる。

10

【0033】

上述のように、ドメインは、図3に示されるようなクラスター化を利用することもできる（マルチキャスト・ハートビート・スライドの特性）。この実施形態についての一般的なネットワーク・トポロジは、マスターを含むハブ・アイランドに接続されたマルチキャスト・アイランドの集合である。マルチキャスト・トラフィックは、ハブから外へ2地点間で転送される。一フェーズ法で分散することができる小さなデルタを、マルチキャストを通じて直接伝送することができる。他の全ての場合、デルタは、2地点間リンクを通じて伝送することができる。ツリー構造化された2地点間転送スキームをハブ・スポーク式マルチキャスト構造の上にオーバーレイし、マスターにおけるボトルネックを減少させることができる。

20

【0034】

図3のドメイン図300において、管理対象サーバ302の1つ又はそれ以上を、クラスター304とも呼ばれるマルチキャスト・アイランド内にまとめることができる。ドメイン308のAdminサーバ306は、ハブ・アイランド312のマスターとして働き、ブラウザ310を介するといった、ドメインへの入力点である。Adminサーバ306は、クラスター・マスターと呼ばれる、クラスター内の管理対象サーバの1つに連絡する。この実施形態におけるAdminサーバは、デルタ又はメッセージを各々のクラスター・マスターにマルチキャストすることができ、次に、各々のクラスター・マスターは、マルチキャストによってデルタ又はメッセージをそのクラスター内の他の管理対象サーバに転送する。クラスター・マスターは、如何なる設定情報も所有せず、代わりに、Adminサーバから情報を受信する。クラスター・マスターがオフラインになるか又はクラッシュした場合には、ドメイン内の別の管理対象サーバがクラスター・マスターにとって代わることができる。この場合、オフラインのサーバが第2のクラスター・マスターとして復帰することを防止するために、機構を所定の位置に置くことができる。このことは、クラスター又はシステム・インフラストラクチャにより取り扱うことができる。

30

【0035】

2つ以上のドメインが存在することもできる。この場合は、ネスト状ドメインすなわち「シンジケート」をおくことができる。各々のドメイン・マスターが他のドメイン・マスターに情報を押し出す能力をもつことができるので、各々のドメイン・マスターに直接触れることにより、情報を該ドメイン・マスターに行き渡らせることができる。しかしながら、ドメイン・マスターにマルチキャストすることは望ましくない。

40

一フェーズ分散において、マスター・ユーザは、更新の分散をトリガするために、ダウンコールを行うことができる。こうしたダウンコールは、

`startOnePhase(DID, newVerNum, scope)`

の形態をとることができ、ここで、DIDは、更新されたデータ項目又はオブジェクトのIDであり、newVerNumは、そのオブジェクトの新しいバージョン番号であり、scopeは、その更新が適用される範囲である。マスターDRSは、新しいバージョン

50

番号を進ませ、該新しいバージョン番号をディスクに書き込み、その後のハートビート内に情報を含ませることにより、応答することができる。

【0036】

スレーブDRSがハートビートを受信した時、該スレーブDRSは、該スレーブDRSが関心ある最新の更新に関する情報のウィンドウを分析することにより、取り出しを必要とするかどうかを判断することができる。スレーブの現在のバージョン番号がウィンドウ内にあり、該スレーブが後にコミットされる更新のいずれの範囲にも入っていない場合には、如何なるデータも取り出すことなく、単純に最新のバージョン番号に進ませることができる。このプロセスは、スレーブが最新である些細な場合を含むことができる。他の場合には、スレーブDRSは、マスターDRSからのデルタに対する2地点間コール、又は

10

`createDelta(DID, curVerNum)`

の形態をとることができる別の同様の要求を行うことができ、ここで、`curVerNum`は、スレーブの現在のバージョン番号であり、ドメイン・マスター又はクラスター・マスターに送り返される。この要求を取り扱うために、マスターDRSは、`createDelta(curVerNum)`のようなアップコールを行うことができる。このアップコールは、デルタ及び新しいバージョン番号を獲得し、それらをスレーブDRSに戻すために、適切なりスナを通じて行うことができる。スレーブが最後にハートビートを受信してから新しいバージョン番号が変更されたかもしれないので、新しいバージョン番号を含ませるべきである。デルタは、最も新しくコミットされた更新までのものだけにすることができる。如何なる進行中の二フェーズ更新も、別個の機構によって取り扱うことができる。次に、スレーブDRSは、`commitOnePhase(newVerNum, delta)`のような、スレーブ・ユーザへのアップコールを行い、その後、新しいバージョン番号に進むことができる。

20

【0037】

二フェーズ更新分散をトリガするために、マスター・ユーザは、`startTwoPhase(DID, oldVerNum, newVerNum, delta, scope, timeout)`のようなダウンコールを行うことができ、ここで、`DID`は更新されるべきデータ項目又はオブジェクトのIDであり、`oldVerNum`は前のバージョン番号であり、`newVerNum`は新しいバージョン番号であり（前のバージョン番号からの1段階）、`delta`取り出されるべき連続するバージョン間のデルタであり、`scope`は更新の範囲であり、`timeout`はそのジョブについての最大有効期間である。「準備」及び「コミット」は同期式であるので、ジョブに対して特定の期限を設定することが望ましい。前のバージョン番号を含ませることができ、これに対して異なるバージョン番号のサーバはデルタを取らない。

30

【0038】

1つの実施形態におけるマスターDRSは、範囲内にある全てのサーバを検討し、各々のスレーブDRSに対して、`prepareTwoPhase(DID, oldVerNum, newVerNum, delta, timeout)`のような2地点間コールを行う。次に、スレーブは、適切なタイムアウト値を取得することができる。バイナリ・コードを含むデルタのように、デルタが大きい場合には、2地点間プロトコルを用いることができる。一フェーズ法を用いて、例えば、キャッシュ・サイズの修正のような小さな設定変更だけを含むことができる小さな更新を行うことができる。アプリケーションの追加のような大きな変更が一貫性をもってサーバに到達することがより重要なので、この手法を用いることができる。代わりに、マスターが、存在する場合にはクラスター・マスターのところにいき、該クラスター・マスターにコールを行わせることができる。クラスター・マスターの代理を務めるマスターをもつことは、システムの規模拡大可能性を改善することができる。

40

【0039】

1つの実施形態において、スレーブ又はクラスター・マスターへの各々のコールは、マ

50

スターDRSにより取り扱われる「到達不能」、「OutOfSync」、「Nak」、「Ack」といった4つの応答のうちの1つを生成する。応答が「到達不能」である場合には、当該サーバに到達することができず、再試行のためにキューに入れることができる。応答が「OutOfSync」である場合には、再試行のためにキューに入れることができる。その間、サーバは、マスターからの取り出しを用いて該サーバ自体を同期させようとするので、再試行時にデルタを受信できるようになる。応答が「NoAck」すなわち肯定応答がない場合には、ジョブは中止される。サーバがジョブを許容できない場合にこの応答が与えられる。応答が「Ack」である場合には、如何なる動作もとられない。

【0040】

スレーブを準備するため、マスターDRSは、prepareTwoPhaseのような方法呼び出すことができる。マスターDRSからの「準備」要求を受信すると、スレーブDRSは、まず、その現在のバージョン番号が更新されるべき古いバージョン番号と等しいかどうかを確認する。等しくない場合には、スレーブは、「OutOfSync」応答を返すことができる。次に、スレーブは、あたかもちょうどハートビートを受信したかのように、マスターDRSからデルタを取り出すことができる。最終的に、マスターDRSは、prepareTwoPhaseを再試行することができる。この手法は、マスターにデルタを取り出させるよりも簡単であるが、マスターの慎重な設定を必要とする。応答を長く待ちすぎることはジョブのタイムアウトを招くので、マスターの設定が必要となる。さらに、十分に応答を待たないことは、「OutOfSync」応答を得るのに余分な要求をもたらすことになる。スレーブからの取り出し要求の完了時に、再試行をトリガ

10

20

【0041】

スレーブが同期している場合には、該スレーブは、prepareTwoPhase(newVerNum, delta)のように、できる限り深くサーバ内のスレーブ側のクライアント層に、アップコールを行うことができる。次に、返される結果物としての「Ack」又は「Nak」を、マスターDRSに送信することができる。応答が「Ack」であった場合には、スレーブは、特別の準備状態になることができる。応答が「Nak」であった場合には、スレーブは、如何なる更新の記録もフラッシュすることができる。何らかの理由で、後でコミットされることになる場合には、スレーブは、これを一フェーズ分散として獲得することができ、その後機能しなくなる。

30

【0042】

マスターDRSがタイムアウト期間内にあらゆるサーバからなんとか「Ack」を収集する場合には、該マスターDRSは、twoPhaseSucceeded(newVerNum)のようなコミット・アップコールを行い、新しいバージョン番号に進むことができる。マスターDRSがいずれかのサーバから「Nak」を受信した場合、又はタイムアウト時間が終了した場合には、マスターDRSは、twoPhaseFailed(newVerNum, reason)のような中止アップコールを行い、バージョン番号を変更しないままにしておくことができる。ここで、reasonは、如何なる「Nak」応答のロールアップも含む例外である。いずれにしても、中止/コミット情報をその後のハートビートに含ませることができる。

40

どんな時でも、マスターDRSは、cancelTwoPhase(newVerNum)のような取消しダウンコールを行うことができる。次に、マスターDRSは、ジョブが進行中でない場合には例外を廃棄することによりこのコールを取り扱うことができ、或いは停止が生じる予定であるように働く。

【0043】

準備されたスレーブDRSが、新しいバージョン番号がコミットされたことを示すハートビートを取得した場合には、該スレーブDRSは、commitTwoPhase(newVerNum)のようなアップコールを行い、新しいバージョン番号に進むことができる。準備されたスレーブDRSが、代わりに、新しいバージョン番号が中止されたことを示すハートビートを取得した場合には、該スレーブはジョブを中止することができる。

50

ウィンドウが新しいバージョン以上に進んだ場合にスレーブがハートビートを取得したとき、該スレーブは、同じデータ項目について新たな `prepareTwoPhase` コールを取得するか、又はジョブをタイムアウトする。このような場合には、スレーブは、`abortTwoPhase(newVerNum)` のようなアップコールを行い、バージョン番号を変更しないままにしておくことができる。このことは、スレーブが準備された後、該スレーブがコミットする前に、マスター・サーバが機能しなくなる場合のような、状況を適切に取り扱うことを保証する 1 つの方法である。

【 0 0 4 4 】

本発明の好ましい実施形態の上記の説明は、図示及び説明目的で提供された。これは、網羅的であるか又は本発明を開示された精密な形態に制限することを意図するものではない。明らかに、多くの修正及び変形が当業者には明らかであろう。実施形態は、本発明の原理及びその実際の用途を最もよく説明するために選択されて述べられたものであり、したがって、当業者が種々の実施形態について及び考慮される特定の用途に適当な種々の修正をもって本発明を理解できるようにするものである。本発明の範囲は、特許請求の範囲及びその均等技術により定義されることが意図される。

10

【図面の簡単な説明】

【 0 0 4 5 】

【図 1】本発明の 1 つの実施形態によるドメイン構造の図である。

【図 2】本発明の 1 つの実施形態による階層化アーキテクチャの図である。

【図 3】本発明の 1 つの実施形態によるクラスター化されたドメイン構造の図である。

20

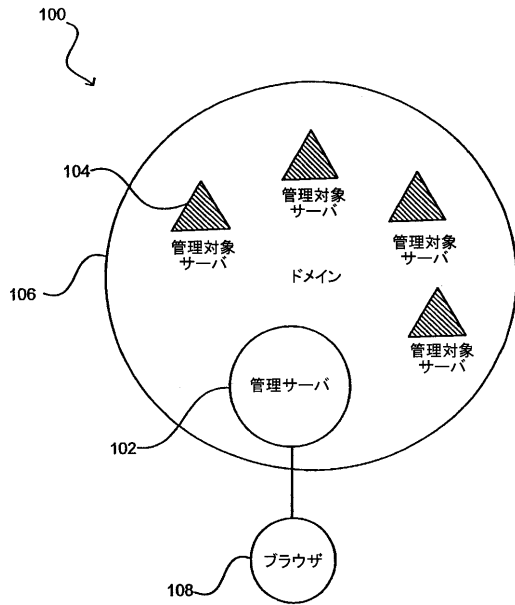
【図 4】本発明の 1 つの実施形態による階層化アーキテクチャのための一フェーズ法の図である。

【図 5】本発明の 1 つの実施形態による階層化アーキテクチャのための二フェーズ法の図である。

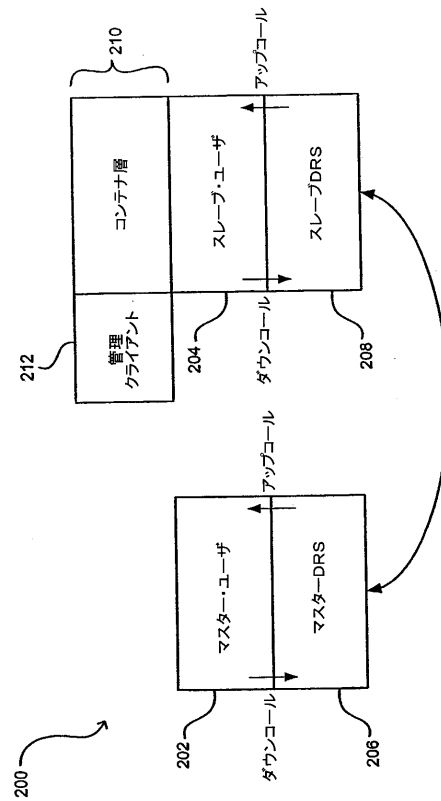
【図 6】本発明の 1 つの実施形態による一フェーズ法のフローチャートである。

【図 7】本発明の 1 つの実施形態による二フェーズ法のフローチャートである。

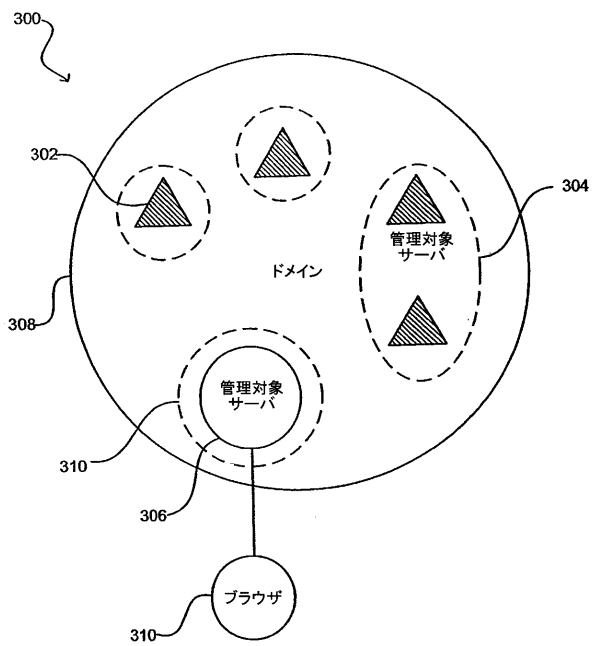
【図 1】



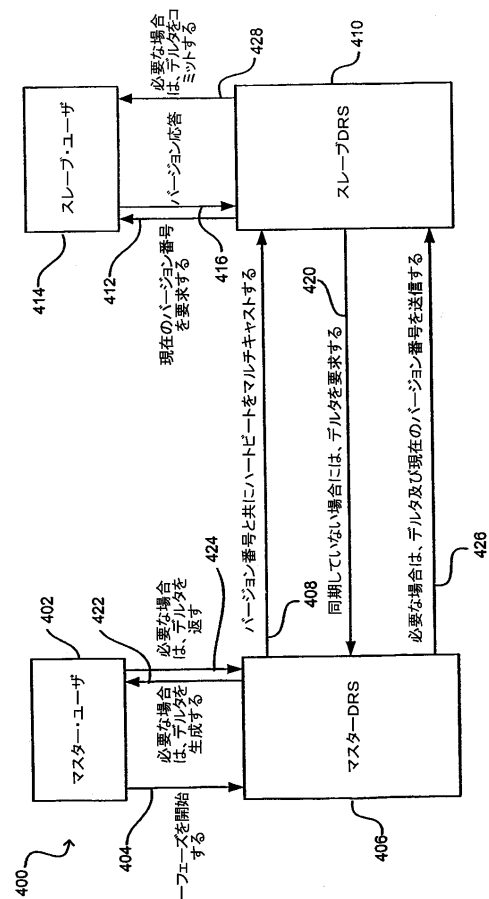
【図 2】



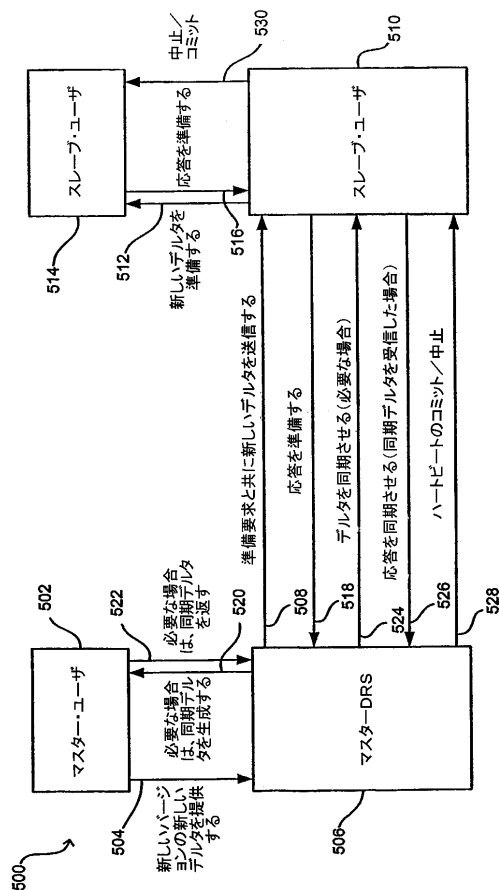
【図 3】



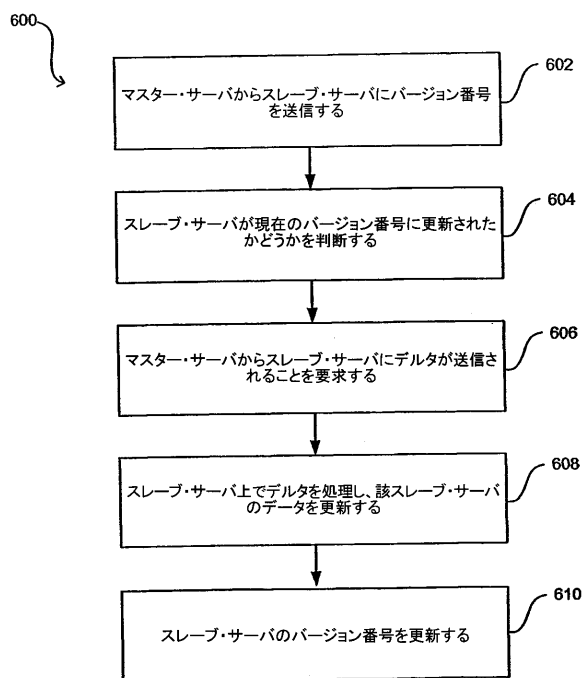
【図 4】



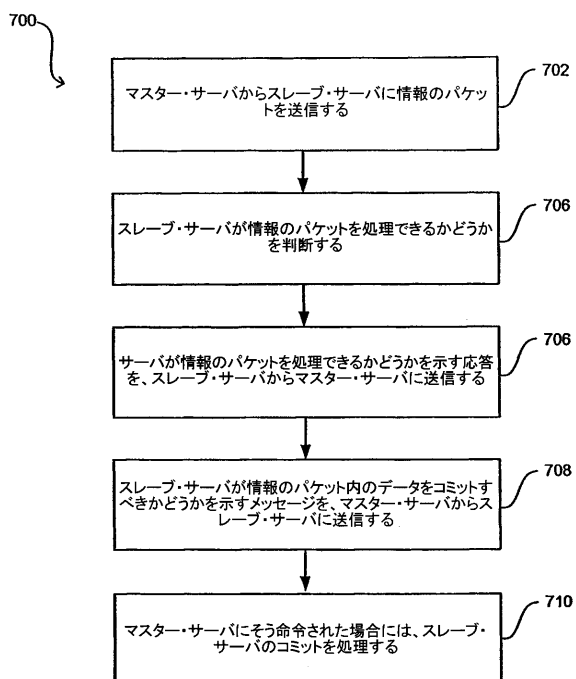
【 図 5 】



【 図 6 】



【 図 7 】



【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/US02/22366
A. CLASSIFICATION OF SUBJECT MATTER IPC(7) : G06F 15/16; G06F 7/00, 17/30 US CL : 709/203; 707/2, 8, 10, 101 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) U.S. : 709/203, 208, 209, 242, 246		
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) Please See Continuation Sheet		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	US 6,088,694 A (BURNS et al.) 11 JULY 2000 (11.07.2000), col. 4 line 30 to col. 6 line 43.	1-93
A	US 5,920,867 A (VAN HUBEN et al.) 06 JULY 1999 (06.07.1999), col. 6 line 50 to col. 8 line 3.	1-93
A, P	US 6,263,372 B1 (HOGAN et al.) 17 JULY 2001, (17.07.2001), col. 3 line 23 to col. 9 line 45.	1-93
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input 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 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 18 November 2002 (18.11.2002)		Date of mailing of the international search report 28 FEB 2003
Name and mailing address of the ISA/US Commissioner of Patents and Trademarks Box PCT Washington, D.C. 20231 Facsimile No. (703)305-3230		Authorized officer Ayaz R Sheikh <i>Ayaz R Sheikh</i> Telephone No. 703 305 3900

INTERNATIONAL SEARCH REPORT

PCT/US02/22366

Continuation of B. FIELDS SEARCHED Item 3:

West, Derwent, EPO, JPO

master, slave, client, server, replicating, duplicating, copying, updating, changing, modifying, delta, phase, version

フロントページの続き

(31)優先権主張番号 09/975,587

(32)優先日 平成13年10月11日(2001.10.11)

(33)優先権主張国 米国(US)

(81)指定国 AP(GH,GM,KE,LS,MW,MZ,SD,SL,SZ,TZ,UG,ZM,ZW),EA(AM,AZ,BY,KG,KZ,MD,RU,TJ,TM),EP(AT, BE,BG,CH,CY,CZ,DE,DK,EE,ES,FI,FR,GB,GR,IE,IT,LU,MC,NL,PT,SE,SK,TR),OA(BF,BJ,CF,CG,CI,CM,GA,GN,GQ,GW, ML,MR,NE,SN,TD,TG),AE,AG,AL,AM,AT,AU,AZ,BA,BB,BG,BR,BY,BZ,CA,CH,CN,CO,CR,CU,CZ,DE,DK,DM,DZ,EC,EE,ES, FI,GB,GD,GE,GH,GM,HR,HU,ID,IL,IN,IS,JP,KE,KG,KP,KR,KZ,LC,LK,LR,LS,LT,LU,LV,MA,MD,MG,MK,MN,MW,MX,MZ,N O,NZ,OM,PH,PL,PT,RO,RU,SD,SE,SG,SI,SK,SL,TJ,TM,TN,TR,TT,TZ,UA,UG,UZ,VN,YU,ZA,ZM,ZW

(74)代理人 100074228

弁理士 今城 俊夫

(74)代理人 100086771

弁理士 西島 孝喜

(72)発明者 ジェイコブス ディーン バーナード

アメリカ合衆国 カリフォルニア州 9 4 7 0 7 パークリー マデラ ストリート 1 7 4 7

(72)発明者 ラマー リート

アメリカ合衆国 カリフォルニア州 9 4 1 3 3 サン フランシスコ グリーン ストリート
4 1 1 # 2 エイ

(72)発明者 スリニヴァサン アナンサム パーラ

アメリカ合衆国 カリフォルニア州 9 4 1 3 1 サン フランシスコ サンチェ ストリート
1 6 1 0

F ターム(参考) 5B082 GA05 HA02 HA03