

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6382958号
(P6382958)

(45) 発行日 平成30年8月29日 (2018. 8. 29)

(24) 登録日 平成30年8月10日 (2018. 8. 10)

(51) Int. Cl.

F I

G 0 6 F 3/06 (2006. 01)

G 0 6 F 3/06 3 0 4 H

G 0 6 F 13/10 (2006. 01)

G 0 6 F 13/10 3 4 0 A

G 0 6 F 13/14 (2006. 01)

G 0 6 F 13/14 3 1 0 Z

G 0 6 F 13/14 3 1 0 B

請求項の数 19 (全 23 頁)

(21) 出願番号 特願2016-516710 (P2016-516710)
 (86) (22) 出願日 平成26年5月27日 (2014. 5. 27)
 (65) 公表番号 特表2016-526229 (P2016-526229A)
 (43) 公表日 平成28年9月1日 (2016. 9. 1)
 (86) 国際出願番号 PCT/US2014/039480
 (87) 国際公開番号 W02014/193770
 (87) 国際公開日 平成26年12月4日 (2014. 12. 4)
 審査請求日 平成29年5月16日 (2017. 5. 16)
 (31) 優先権主張番号 13/904, 989
 (32) 優先日 平成25年5月29日 (2013. 5. 29)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 314015767
 マイクロソフト テクノロジー ライセン
 シング, エルエルシー
 アメリカ合衆国 ワシントン州 9805
 2 レッドモンド ワン マイクロソフト
 ウェイ
 (74) 代理人 100140109
 弁理士 小野 新次郎
 (74) 代理人 100075270
 弁理士 小林 泰
 (74) 代理人 100101373
 弁理士 竹内 茂雄
 (74) 代理人 100118902
 弁理士 山本 修

最終頁に続く

(54) 【発明の名称】 クラスタにおけるストレージ防御の分散

(57) 【特許請求の範囲】

【請求項 1】

ストレージデバイスへのアクセスを可能にするための方法であって、
 特定のノードにより、前記ストレージデバイスへ書込み要求を送信するステップと、
 前記特定のノードにより、前記ストレージデバイスから登録テーブルを受信するステッ
 プであって、前記登録テーブルはノードのクラスタ内の1以上の他のノードを識別する、
 ステップと、

前記特定のノードにより、前記特定のノードの第1の登録キーが前記登録テーブル内に
 存在しないかを決定するステップと、

前記特定のノードにより、第2の登録キーを、前記クラスタへの加入再承認及び前記ス
 トレージデバイスへの書込みアクセスのために前記ストレージデバイスへ送信するステッ
 プと、

前記第2の登録キーを送信するステップに基づいて、前記特定のノードにより、第1の
 登録タイマを設定するステップであって、前記ストレージデバイスにより維持される前記
 登録テーブルを変更するために、前記第1の登録タイマは、前記ノードのクラスタ内の1
 以上のノードに、時間期間を提供する、ステップと、

前記第1の登録タイマの期限切れと同時に、前記特定のノードにより前記ストレージデ
 バイスから前記登録テーブルを受信するステップであって、前記登録テーブルは、前記時
 間期間の間のノードのクラスタ内の前記1以上の他のノードからの変更を含む、ステッ
 プと、

10

20

前記特定のノードにより、前記第2の登録キーが前記登録テーブル内に記憶されているかを決定するステップと、

前記第2の登録キーが前記登録テーブルに記憶されているときに、前記クラスタ内の前記1つ以上の他のノードを通じて前記ストレージデバイスへのリモート接続を確立し、前記ストレージデバイスへ書込みするステップとを含む方法。

【請求項2】

前記クラスタに加入するステップ後、前記登録テーブルをスクラブするステップをさらに含む、請求項1に記載の方法。

【請求項3】

前記登録テーブルをスクラブするステップは、
前記登録テーブルを前記ストレージデバイスから受信するステップと、
前記登録テーブルを読み取るステップと、
前記登録テーブルから1つまたは複数の登録キーを抹消するステップであって、前記1つまたは複数の登録キーのそれぞれは、前記クラスタ内の他のノードと関連している、ステップとを含む、請求項2に記載の方法。

【請求項4】

請求項3に記載の方法であって、さらに第2の登録タイマを設定するステップを備え、前記第1の登録タイマの長さは、前記第2の登録タイマの長さの少なくとも2倍である、方法。

【請求項5】

前記第1の登録キーおよび第2の登録キーのそれぞれは、クラスタ識別子、およびノード識別子を含む、請求項1に記載の方法。

【請求項6】

コンピュータ実行可能な命令を符号化するコンピュータ可読ストレージデバイスであって、前記コンピュータ実行可能な命令は、1つまたは複数のプロセッサによって実行されると、ストレージデバイスへのアクセスを可能にするための方法を行い、前記方法は、

特定のノードにより、前記ストレージデバイスへ書込み要求を送信するステップと、
前記特定のノードにより、前記ストレージデバイスから登録テーブルを受信するステップであって、前記登録テーブルはノードのクラスタ内の1以上の他のノードを識別する、ステップと、

前記特定のノードにより、前記特定のノードの第1の登録キーが前記登録テーブル内に存在しないことを決定するステップと、

前記特定のノードにより、第2の登録キーを、前記ストレージデバイスへ送信するステップと、

前記第2の登録キーを送信するステップに基づいて、前記特定のノードにより、第1の登録タイマを設定するステップであって、前記ストレージデバイスにより維持される前記登録テーブルを変更するために、前記第1の登録タイマは、ノードのクラスタ内の1以上の他のノードに、時間期間を提供する、ステップと、

前記第1の登録タイマの期限切れと同時に、前記特定のノードにより前記ストレージデバイスから前記登録テーブルを受信するステップであって、前記登録テーブルは、前記時間期間の間のノードのクラスタ内の前記1以上のノードからの変更を含む、ステップと、

前記特定のノードにより、前記第2の登録キーが前記登録テーブル内に記憶されているかを決定するステップと、

前記第2の登録キーが前記登録テーブルに記憶されているときに、前記クラスタ内の前記1つ以上のノードを通じて前記ストレージデバイスへのリモート接続を確立し、前記ストレージデバイスへ書込みするステップとを含む方法。

【請求項7】

前記クラスタに加入するステップ後、前記登録テーブルをスクラブするための命令をさらに含む、請求項 6 に記載のコンピュータ可読ストレージデバイス。

【請求項 8】

前記登録テーブルをスクラブするステップは、
前記ストレージデバイスから前記登録テーブルを受信するステップと、
前記登録テーブルを読み取るステップと、
前記登録テーブルから 1 つまたは複数の登録キーを抹消するステップであって、前記 1 つまたは複数の登録キーのそれぞれは、前記クラスタ内の他のノードと関連している、ステップと
を含む、請求項 7 に記載のコンピュータ可読ストレージデバイス。

10

【請求項 9】

請求項 8 に記載のコンピュータ可読ストレージデバイスであって、第 2 の登録タイマを設定するための命令をさらに備え、前記第 1 の登録タイマの長さは、前記第 2 の登録タイマの長さの少なくとも 2 倍である、コンピュータ可読ストレージデバイス。

【請求項 10】

請求項 9 に記載のコンピュータ可読ストレージデバイスであって、前記ストレージデバイスへの接続を生成するための命令をさらに備える、コンピュータ可読ストレージデバイス。

【請求項 11】

前記接続は物理的接続である、請求項 10 に記載のコンピュータ可読ストレージデバイス。

20

【請求項 12】

前記第 1 の登録キー及び第 2 の登録キーのそれぞれは、クラスタ識別子、およびノード識別子を含む、請求項 6 に記載のコンピュータ可読ストレージデバイス。

【請求項 13】

ストレージデバイスへのアクセスを可能にするための方法であって、
クラスタ内の特定のノードから、前記ストレージデバイスに書込みコマンドを送信するステップであって、前記特定のノードは、関連登録キーを有する、ステップと、
前記書込みコマンドが拒否されたという通知を受信すると同時に、

前記特定のノードにより、前記ストレージデバイスに登録テーブルを要求するステップと、

30

前記特定のノードにより、前記特定のノードと関連する前記関連登録キーが前記登録テーブルに存在するかどうかを決定するステップと、

前記特定のノードと関連する前記関連登録キーが前記登録テーブルに存在しないとき、

前記特定のノードにより、前記ストレージデバイスに新規の登録キーを送信するステップと、

前記特定のノードにより、第 1 の登録タイマを設定するステップであって、前記第 1 の登録タイマは前記ストレージデバイスにより維持された前記登録テーブルを修正するために、ノードのクラスタ内の 1 以上の他のノードに時間期間を供給する、ステップと、

40

前記第 1 の登録タイマの期限切れと同時に、前記特定のノードにより、前記ストレージデバイスから前記登録テーブルを受信するステップであって、前記登録テーブルは、前記時間期間の間にノードのクラスタ内の前記 1 以上の他のノードからの修正を含む、ステップと、

前記特定のノードにより、前記新規の登録キーが前記登録テーブルに記憶されているかどうかを決定するステップと、

前記新規の登録キーが前記登録テーブルに記憶されているとき、前記クラスタに加入するステップであって、前記クラスタに加入するステップは、前記クラスタ内の 1 以上のノードを通じて前記ストレージデバイスへのリモート接続を確立することにより、前記特定のノードが前記ストレージデバイスに書き込むことを可能にする、ステップと

50

を含む方法。

【請求項 1 4】

前記クラスタに加入するステップ後、前記登録テーブルをスクラブするステップをさらに含む、請求項 1 3 に記載の方法。

【請求項 1 5】

請求項 1 4 に記載の方法であって、前記登録テーブルをスクラブするステップは、
前記ストレージデバイスから前記登録テーブルを受信するステップと、
前記登録テーブルを読み取るステップと、
前記登録テーブルから 1 以上の登録キーを抹消するステップであって、1 以上の登録キーのそれぞれは、前記クラスタ内の他のノードと関連している、ステップと、
を含む方法。

10

【請求項 1 6】

請求項 1 5 に記載の方法であって、第 2 の登録タイマを設定するステップをさらに含み、前記第 1 の登録タイマの長さは、前記第 2 の登録タイマの長さの少なくとも 2 倍である、方法。

【請求項 1 7】

請求項 1 5 に記載の方法であって、1 以上の登録キーを抹消するステップは、関連する前記ノードが前記ストレージデバイスへのアクセスを拒否されるようにする、方法。

【請求項 1 8】

請求項 1 3 に記載の方法であって、前記新規の登録キーは、クラスタ識別子、およびノード識別子を含む、方法。

20

【請求項 1 9】

前記第 1 の登録タイマは、前記ノードのクラスタ内の全てのノードが前記登録テーブルをスクラブするのに十分な時間期間で構成されている、請求項 1 に記載の方法。

【発明の詳細な説明】

【背景技術】

【0 0 0 1】

[0001] ノードのクラスタがストレージデバイスにアクセスできるという通常の共有ストレージ状況では、クラスタ内の少なくとも 1 つのノードが、ストレージデバイスに接続される。結果として、ストレージデバイスに接続されているノードは、ストレージデバイスの防御の対処を担う。しかしながら、クラスタが複数のストレージデバイスにアクセスできる状況では、クラスタ内の単一のノードを各ストレージデバイスに接続することはできない。結果として、ストレージデバイスの一部は、保護されない場合がある。

30

【発明の概要】

【発明が解決しようとする課題】

【0 0 0 2】

[0002] これらおよび他の概括的考慮事項に関して、諸実施形態は作成されている。また、比較的具体的な問題を論じているが、諸実施形態は、この背景技術で特定された具体的な問題を解決すること限定されるべきでないことを理解されたい。

【課題を解決するための手段】

40

【0 0 0 3】

[0003] この概要は、「発明を実施するための形態」の項でさらに後述する概念選択を簡略化された形態で紹介するために提供される。この概要は、特許請求される主題の主要な特徴または本質的な特徴を特定するように意図するものでも、特許請求される主題の範囲を決定する際の助けとして使用すべきように意図するものでもない。

【0 0 0 4】

[0004] 本開示の実施形態は、ストレージデバイスへのアクセスを可能にし、クラスタ内の様々なノードによってアクセス可能な 1 つまたは複数のストレージデバイスを保護するための方法およびシステムを提供する。具体的には、1 つまたは複数の実施形態は、ノードがいかにクラスタに加入承認され得、それによって、ノードのクラスタ内の少なくとも

50

1つのノードに接続されているストレージデバイスへの読取り/書込みアクセスを取得することができるかについて記載する。加えて、1つまたは複数の実施形態は、ノードがクラスタと関連する登録テーブルを監視し、認識されていないノードからのエントリを除去することが可能であることを提供する。除去されたそれらのノードに関しては、ノードは、クラスタに対する再承認を求めるために、登録テーブルに再登録しようと試みることができる。

【0005】

[0005]以下に説明されるように、ストレージデバイスへのアクセスを求めるノードは、クラスタ通信プロトコルを使用して、クラスタ内に入る。ノードは、一旦、クラスタに加入承認されると、クラスタによって利用される1つまたは複数のストレージデバイスへのアクセスを得る資格を有することができる。ストレージデバイスへのアクセスを得るために、ストレージデバイスへのアクセスを求めるノードは、ストレージデバイスに登録キーを送信する。ストレージデバイスに登録後、ノードは、登録タイマを設定する。諸実施形態では、登録タイマは、要求されたアクセスがストレージデバイスへのアクセスを求めるノードに付与されるべきであるかどうかをクラスタ内の各ノードが決定する機会を得るその間の時間期間に相当する。登録タイマの期限切れと同時に、ストレージデバイスへのアクセスを求めるノードは、ストレージデバイスから登録テーブルを受信する。一旦、登録テーブルを受信されると、ノードは、その登録キーが登録テーブルに記憶されているかどうかを決定する。登録キーが登録テーブルに記憶されている場合、ノードは、ストレージデバイスへのアクセスが許可される。より具体的には、ノードには、ストレージデバイスへの書込みアクセスが付与される。

【0006】

[0006]諸実施形態は、コンピュータ処理、コンピューティングシステムとして、またはコンピュータプログラム製品もしくはコンピュータ可読媒体などの製造品として実装可能である。コンピュータプログラム製品は、コンピュータシステムによって読取り可能で、コンピュータ処理を実行するためのコンピュータプログラム命令を符号化した、コンピュータストレージ媒体であってよい。コンピュータプログラム製品はまた、コンピューティングシステムによって読取り可能で、コンピュータ処理を実行するためのコンピュータプログラム命令を符号化した、搬送波上の伝播信号であってもよい。

【0007】

[0007]限定的でなく、包括的でない実施形態は、以下の図を参照して説明される。

【図面の簡単な説明】

【0008】

【図1】[0008]本開示の1つまたは複数の実施形態による、クラスタ内の複数のノードが、それぞれのストレージデバイスに接続されているシステムを示す図である。

【図2】[0009]本開示の1つまたは複数の実施形態による、クラスタにおけるメンバシップを要求するための方法を示す図である。

【図3】[0010]本開示の1つまたは複数の実施形態による、ノードのクラスタと関連するストレージデバイスへのアクセスを決定するための方法を示す図である。

【図4】[0011]本開示の1つまたは複数の実施形態による、クラスタにおける加入再承認を要求するための方法を示す図である。

【図5】[0012]本開示の1つまたは複数の実施形態による、クラスタ内の様々なノードがいかにして物理的ストレージデバイスに接続可能であることを示すブロック図である。

【図6】[0013]本開示の1つまたは複数の実施形態とともに使用可能なコンピューティングデバイスの例示的な物理的構成要素を示すブロック図である。

【図7A】[0014]本開示の1つまたは複数の実施形態とともに使用可能なモバイルコンピューティングデバイスの簡略化されたブロック図である。

【図7B】本開示の1つまたは複数の実施形態とともに使用可能なモバイルコンピューティングデバイスの簡略化されたブロック図である。

【図8】[0015]本開示の1つまたは複数の実施形態とともに使用可能な分散されたコンピ

10

20

30

40

50

ューティングシステムの簡略化されたブロック図である。

【発明を実施するための形態】

【0009】

[0016]様々な実施形態が、その一部を成し、具体的な例示的实施形態を示す添付の図面を参照して、より十分に後述される。しかしながら、諸実施形態は、多数の種々の形態で実装可能であり、本明細書に示される実施形態に限定されるものと見なすべきでなく、そうではなくて、これらの実施形態は、本開示が、徹底的で完全になるように、および当業者に対して実施形態の範囲を十分に伝えることになるように提供される。諸実施形態は、方法、システム、またはデバイスとして実践可能である。したがって、実施形態は、ハードウェアの実装形態、全面的にソフトウェアの実装形態、またはソフトウェアの態様とハードウェアの態様とを組み合わせた実装形態の形態を取ることができる。そのため、以下の詳細な説明は、限定する意味で解釈すべきではない。

【0010】

[0017]図1は、本開示の1つまたは複数の実施形態による、クラスタ102内の複数のノードが、それぞれのストレージデバイスに接続されているシステム100を示している。図1に示すように、クラスタ102は、複数のノード102A~102Dを含むことができる。4つのノードが示されているが、クラスタ102が、4つより多いノード、または4つより少ないノードを有する場合があることが企図される。特定の实施形態では、ノードが、例えば、パーソナルコンピュータ、タブレット、ラップトップ、スマートフォン、およびパーソナルデジタルアシスタントなど、コンピューティングデバイスであってよい。他の実施形態では、ノードが、サーバコンピューティングデバイスであってもよい。

【0011】

[0018]図1はまた、クラスタ102内の各ノードが、1つまたは複数のストレージデバイスに接続されていることを示している。特定の实施形態では、ストレージデバイスは、ダイレクトアタッチドストレージデバイス(すなわち、ホストシステムまたはホストデバイスに直接、接続されているストレージデバイス)とすることができる。また、ストレージデバイスが、1つまたは複数のバスを使用して、クラスタ内のいくつかのノードによってアクセス可能であることも企図される。例えば、1つまたは複数のノードは、ストレージデバイスに物理的に接続可能である一方、クラスタ内の他のノードは、リモートバスを使用して、ストレージデバイスに接続することができる。加えて、単一のノードが、様々なストレージデバイスに対する複数の物理的接続と、様々なストレージデバイスに対する1つまたは複数のリモート接続とを有することができる。また、クラスタ内の各ノードは、クラスタ内の他のノードのそれぞれの活動および接続をビューでき得ることが企図される。つまり、システム100は、一部のストレージデバイスが、一部のノードにとって利用可能である一方、他のストレージデバイスは、それらのノードにとって利用できないという点で非対称とすることができる。

【0012】

[0019]例えば、図1に示すように、ノード102Aとノード102Bは、ストレージデバイス104に接続され、ノード102Cは、ストレージデバイス104とストレージデバイス106に接続され、ノード102Dは、ストレージデバイス106とストレージデバイス108に接続されている。特定の实施形態では、ストレージデバイス104~108は、ストレージプールを備える。ストレージプール内の各ストレージデバイスにアクセスできる単一のノードがクラスタ102内に存在しないとき、クラスタ102内の各ノードは、防御アルゴリズムを動作させて、確実に、クラスタの一部であるノードだけがストレージデバイスに読取り/書込みアクセスできるようにする役割を担う。したがって、クラスタ102内の各ノードは、それらが、ストレージプールにおいて接続されているストレージデバイスを同時に保護する。

【0013】

[0020]図1に戻って参照すると、ノード102A、102B、および102Cのそれぞれは、ストレージデバイス104に接続されている。論じたように、ノードのそれぞれは

10

20

30

40

50

、ストレージデバイス 104 への物理的接続、またはリモート接続（すなわち、ストレージデバイス 104 への物理的接続を有するノードを通じたストレージデバイス 104 への接続）を有することができる。ノード 102 A、102 B、および 102 C は、ストレージデバイス 104 に接続されているので、各ノードは、ストレージデバイス 104 に読み取り/書き込みアクセスできる。さらには、クラスタ 102 内のノード 102 A、102 B、および 102 C のそれぞれは、クラスタ 102 内の他のノードの存在を検出し、クラスタ内の他のノードのそれぞれの活動を決定することができる。

【0014】

[0021] 諸実施形態では、特定のストレージデバイスに対する権利は、永続的な予約（persistent reservation）によって決定される。つまり、例えば、ストレージデバイス 104 などのストレージデバイスは、オフラインであるときでも、または再起動したときでも、特定のノードの予約を維持する。議論を進めるために、特定のノードの予約は、あるノードがある特定のストレージデバイスを予約したときに行われ、別の認可されていないノードがそのストレージデバイスにアクセスするのを防ぐ。

【0015】

[0022] 上記の例に戻って参照すると、ノード 102 A、102 B、および 102 C のそれぞれは、クラスタ 102 の一部であるという理由から、ストレージデバイス 104 に読み取り/書き込みアクセスできる。下記に詳細に説明されることになるが、クラスタ 102 内の各ノードは、時間 t において防御アルゴリズムを動作させて、クラスタ 102 内の他のいずれのノードも、(i) クラスタ内の他のノードに対する接続性か、または (ii) ストレージデバイス 104 に対する接続性を失っているかどうかを決定する。

【0016】

[0023] 例えば、ノード 102 A が、ノード 102 B と 102 C に対する接続性、またはストレージデバイス 104 に対する接続性を失う場合、ノード 102 A はもはや、ストレージデバイス 104 に（少なくとも）書き込みアクセスできない、したがって、ストレージデバイスへのアクセスが禁止されるべきであることを、ノード 102 B またはノード 102 C は独立に決定する。接続が失われている場合、ノード 102 B および 102 C は、ノード 102 A のワークロードを取り、また、ストレージデバイスへの書き込みをノード 102 A に許可することにより、ストレージデバイス 104 上のデータを破損させる可能性がある、ノード 102 A がもはや、ストレージデバイス 104 への書き込みを確実にできないようにするステップを取る。ノード 102 A は、ノード 102 B と 102 C に対する接続性を失っている場合はあるものの、ストレージデバイス 104 に対する接続性をなおも有することは可能であることが企図される。同様に、ノード 102 A は、ストレージデバイス 104 に対する接続性を失った場合であっても、なおもノード 102 B および/またはノード 102 C に接続可能であることが企図される。

【0017】

[0024] 上記の例に戻ると、ノード 102 A がストレージデバイス 104 に書き込むことを禁止するためには、ノード 102 B または 102 C は、ノード 102 A をノード登録テーブルから抹消する要求をストレージデバイス 104 に送信する。より具体的には、要求は、ノード 102 A と関連する登録キーをノード登録テーブルから抹消するように、ストレージデバイスに送信される。結果として、ストレージデバイス 104 はもはや、ノード 102 A と関連する物理パス、またはリモートパスから書き込みコマンドを受け付けなくなる。特定の実施形態では、ノード 102 A は、ストレージデバイス 104 に書き込みアクセスできなくなるものの、なおもストレージデバイス 104 への読み取りアクセスはできることになる。

【0018】

[0025] 特定の実施形態では、ノード 102 B または 102 C はいずれも、互いに独立に抹消要求を送信することができる。別の実施形態では、クラスタのノードは、特定のノードが抹消されるべきであるかどうかに関して合意に達するように求められる場合がある。さらなる別の実施形態では、ノードがクラスタから除去されるべきであることをノード自

10

20

30

40

50

体が決定することができる。例えば、ノード 102A は、ノード 102A が、他のノードのうちの 1 つに対する接続、またはストレージデバイス 104 に対する接続を失っていると決定する場合、それ自体からストレージデバイスへの 1 つまたは複数のパスを除去し、あるいはノード登録テーブルからその登録キーを除去するようにストレージデバイス 104 に命令することができる。

【0019】

[0026] ノード登録テーブルを参照すると、特定の実施形態では、ノード登録テーブルは、ストレージデバイス 104 によって維持され、ストレージデバイス 104 に対して書込みアクセスできるノードをリストする。諸実施形態では、ノード登録テーブルは、ストレージデバイスに対して書込みアクセスできるノードごとに登録キーを含む。特定の実施形態では、登録キーは、以下の形式、すなわち、(i) (クラスタ内のすべてのノードについて同じである) クラスタグローバル意識別子の 32 ビットハッシュ、(ii) 8 ビットキーリビジョン(key revision)、(iii) 8 ビットノード番号、および(iv) 16 ビットシグネチャを有する 64 ビット整数を含む。登録キーの具体的なサイズおよび構成が示されているが、登録キーは、登録キーが各ノードに固有である限り、任意の数のビットを有しても、また様々な構成を有してもよいことが企図される。

10

【0020】

[0027] 以下により詳細に説明されることになるように、ノードは、一旦、抹消されると、クラスタへの加入再承認を要求することができる。加入再承認を要求するために、抹消されたノードは、更新済みの登録キーをストレージデバイスに送信することができる。一旦、ノードがストレージデバイスに再登録されると、クラスタ内の他のノードのそれぞれは、ノードがクラスタに加入再承認されるべきであるかどうかについての決定を行う。クラスタ内のノードの決定に基づいて、加入再承認を求めるノードには、加入再承認が付与され得る、または加入再承認が拒否され得る。諸実施形態では、クラスタ内のノードは、加入再承認を求めるノードの接続速度、加入再承認を求めるノードの信頼度、および加入再承認を求めるノードが、ストレージプール内の他のストレージデバイスに対してできるアクセスなどを含む任意の数の要因に基づいて、その決定を行うことができる。

20

【0021】

[0028] 図 2 は、本開示の 1 つまたは複数の実施形態による、クラスタにおけるメンバシップを要求するための方法 200 を示している。特定の実施形態では、方法 200 は、図 1 のクラスタ 102 などのクラスタにおけるメンバシップを要求するノードによって使用され得る。上記に論じたように、ノードは、一旦、クラスタに加入承認されると、その特定のストレージデバイスに対して読取り/書込みアクセスでき得る。

30

【0022】

[0029] 具体的には、加入するノードは、アクティブなクラスタへの加入承認を得るために、クラスタ通信プロトコルを使用して他のノードと通信しようと試みることができる。この場合、加入するノードは、一旦、アクティブ状態に入ると、方法 200 を実行して、ストレージへのアクセスを得ることになる。加入するノードが、クラスタプロトコルを介して、他の加入するノード、またはアクティブなノードと通信するのに失敗し、アクティブなクラスタは存在しない場合もあり得ると考える場合、加入するノードは、方法 200 を実行して、ストレージへのアクセスを得て、したがって、第 1 のアクティブなノードになることができる。議論を進めるために、ノードは、クラスタへのアクセスを要求し、加入承認を得るとき、アクティブなノードとして、またはアクティブ状態に入っていると見られる。例えば、クラスタ通信プロトコルを動作中であり、現在、クラスタメンバシップに参加中であるすべてのノードは、アクティブなノードと見なされる。加えて、クラスタと関連する 1 つまたは複数のストレージデバイスにアクセスできるノードは、ストレージノードと見なされる。諸実施形態では、ストレージノードは、アクティブなノードセットのサブセットである。

40

【0023】

[0030] 図 2 に示すように、方法 200 は、ノードが、クラスタ通信プロトコルを使用し

50

て、クラスタに加入承認された後に開始する。一旦、クラスタに加入承認されると、1つまたは複数の実施形態は、登録キーを使用してストレージデバイスに登録すること(210)によって、ノードが、例えば、ストレージデバイス104(図1)など、クラスタと関連する1つまたは複数のストレージデバイスにアクセスするように求めることを実現する。論じたように、登録キーは、上記に論じた様々な構成要素を有する64ビット整数を含むことができる。

【0024】

[0031]一旦、登録キーがストレージデバイスに送信されると、フローは動作220に進み、ここで、登録タイマが設定される。特定の実施形態では、この登録タイマは、ストレージデバイスへのアクセス、またはストレージプールへのアクセスを要求しているノードによって維持され得る。しかしながら、ストレージデバイスまたはクラスタ内の別のノードもまた、登録タイマを維持することが可能であり得ることが企図される。諸実施形態では、登録タイマの長さは、時間期間 t に相当する。特定の実施形態では、時間期間 t は、3秒である。より具体的には、時間期間 t は、クラスタ内の他のいずれのノードも、(CPUロード、およびI/O待ち時間などに起因して生じ得るいずれの遅延も考慮して)3秒毎に行われるべきスクラブを行うための時間に相当する。

【0025】

[0032]タイマの期限切れと同時に、フローは動作230に進み、ここで、ノード登録テーブルは、ストレージデバイスから読み取られる。上記に論じたように、ノード登録テーブルは、ストレージデバイスによって(またはストレージプール内の少なくとも1つのストレージデバイスによって)維持され、クラスタ内のあらゆるノードと関連する各登録キーのリストを含んでいる。

【0026】

[0033]一旦、登録テーブルが受信され、読み取られると、フローは動作240に進み、ここで、ノードの登録キーがノード登録テーブルに含まれているかどうか決定される。ノードの登録キーがノード登録テーブルに含まれている場合、クラスタ内のノードのそれぞれは、防御アルゴリズムを動作させており、ストレージデバイスは、ストレージデバイスへのアクセスを求めるノードの要求を受け付けている。結果として、フローは動作250に進み、ここで、ノードは、ストレージデバイスにアクセスするように、より具体的には、ストレージデバイスに書込みアクセスできるように許可される。

【0027】

[0034]しかしながら、要求するノードの登録キーが、ノード登録テーブル内に存在しないことが動作240で決定される場合、フローは動作210に戻るようになり、ノードは、ストレージデバイスに登録しようと2回目の試みを行う。方法は、繰り返し、ストレージデバイスへのアクセスを要求するノードは、再度、登録テーブルを要求し、読み取って、その登録キーがノード登録テーブルに記憶されているかどうかを決定する。

【0028】

[0035]図3は、本開示の1つまたは複数の実施形態による、ノードのクラスタと関連するストレージデバイスへのアクセスを決定するための方法300を示している。特定の実施形態では、方法300は、アクティブなノードとして見られるクラスタ内の各ノード(すなわち、クラスタ通信プロトコルを動作中であり、現在、クラスタメンバシップに加入中であるすべてのノード)によって行われる。方法300はまた、ストレージノードと見なされるノード(すなわち、クラスタと関連する1つまたは複数のストレージデバイスにアクセスできるいずれのノード)によっても行われ得る。

【0029】

[0036]方法300は、ノードが、ノード登録テーブルからの登録キーを「スクラブする(scrub)」(310)ときに開始する。具体的には、ノードは、アクティブなクラスタの一部ではない他のノードを探す。登録キーがディスク登録テーブルからスクラブされるとき、スクラブされた登録キーと関連するノードはもはや、クラスタと関連する特定のストレージデバイスまたはストレージプールに、(少なくとも)書込みアクセスできない。

特定の実施形態では、クラスタ内の1つのノードが、クラスタ内の別のノードとはもはや、特定のストレージデバイスまたはストレージプールに書き込みアクセスできないはずであると考えるとき、登録キーがスクラップされる。これは、ノードがクラスタ内の別のノードへの接続を失うこと、ノードがクラスタからそれ自体を除去すること、またはノードからストレージデバイスへの接続を失うことのうちの1つの結果とすることができる。ノードが、それ自体をクラスタから除去している状況では、そのノードは、その登録キーがノード登録テーブルから除去されるべきであるということを示す要求をストレージデバイスに送信することができる。別の実施形態では、クラスタ内の他のノードのうちの1つは、スクラッピング処理中、ノードがクラスタから抹消される（すなわち、除去される）ことを要求することができる。特定の実施形態では、ノードがクラスタから抹消される場合、他のノードは、除去されたノードからのコマンドがストレージデバイスに達するのを防ぐように構成可能である。

10

【0030】

[0037]図3に示すように、スクラッピング処理は、様々なサブ動作を有する。スクラッピング処理は、サブ動作311で開始し、ここで、ノードは、ストレージデバイスによって維持されるノード登録テーブルを読み取る。上記に論じたように、ノード登録テーブルは、クラスタ内の各ノードと関連する登録キーすべてのリストを含んでいる。

【0031】

[0038]次いで、フローはサブ動作312に進み、ここで、クラスタ内のアクティブなメンバシップを有していない1つまたは複数のノードが、クラスタから抹消される。諸実施形態では、クラスタ内の各ノードは、クラスタ内の他のあらゆるノードによってビューでき、リモート接続、または物理接続のいずれかによって、1つまたは複数のストレージデバイスに接続され得る。クラスタ内の各ノードが、クラスタ内の他のあらゆるノードのビューを有すると、ノード登録テーブルを読み取り中であるノードは、クラスタ内のどのノードが、ノード登録テーブル内の関連登録キー有しているのかを決定することができる。したがって、ノードは、ストレージデバイスから受信した登録テーブルをスクラップする。登録キーはテーブル内にあるが、ノードがアクティブでない場合、ノードは抹消される。

20

【0032】

[0039]特定の実施形態では、ノードは、クラスタ内の複数のノードが類似の決定（すなわち、抹消されようとしているノードはノード登録テーブルに登録キーを有していない）に達するまで、抹消され得ない。他の実施形態では、抹消されるべきノードはノード登録テーブルに関連登録キーを有していないという決定に単一のノードが達するとき、ノードが、抹消され得る。

30

【0033】

[0040]一旦、ノードが抹消されるように、要求が送信されると、アルゴリズムを動作させるノードは、ノードが成功裏に抹消されたかどうかを決定する（313）。ノードが成功裏に抹消されなかった場合、アルゴリズムを動作させるノードは、自己チェックを行って、それ自体の登録キーがノード登録テーブル内に存在するかどうかを決定する。それ自体の登録キーがノード登録テーブル内に存在しない場合、フローはサブ動作314に進み、ノードは、「ゲスト状態（guest state）」に入り、1つまたは複数のクラスタプロトコルを使用して、クラスタへの加入再承認を求める。ノードが成功裏に抹消された場合、フローは動作315に進み、ストレージデバイスがまだ予約されていない場合、ノードは、ストレージデバイスを予約する。つまり、ノードは、ストレージデバイスの所有権を取り、次いで、このストレージデバイスへのアクセスを、同じクラスタ内にある他のノードと共有する。

40

【0034】

[0041]ノード登録テーブル内の登録キーがスクラップ済みである場合、フローは動作32に進み、ここで、第2のタイマが設定される。諸実施形態では、第2のタイマは、図2で論じた登録タイマに比例する。例えば、スクラッピングアルゴリズムは、すべてのノードにおいて並行して動作する（例えば、各ノードは、すべての他のノードと独立に3秒毎にキ

50

ーをスクラブする。特定の実施形態では、スクラビングアルゴリズムは、別のノードもやはり、登録テーブルをスクラブしているかどうか、またはノードのうちの1つがストレージへのその接続を失っており、そのため、登録テーブルをスクラブすることができないかどうかを1つのノードが信頼可能に見分けることができない可能性があるという理由から、並行して動作する。

【0035】

[0042] ノードは、クラスタに加入するとき、クラスタに登録し、加入承認されるようになるのを待つ。ノードが加入承認される場合、1つの実施形態は、できる限りすぐにノードが、ディスクをスクラブし、明らかにする(surface)ことを実現する。3秒スクラビングタイマが設定され、ノードは、ストレージ上のそれ自体の登録を認証することができる。

10

【0036】

[0043] タイマの期限切れと同時に、フローは動作330に進み、ここで、ノードは、ストレージデバイスへの1つまたは複数のパスを明らかにする。つまり、ノードは、物理パス、またはリモートパスのいずれのパスが、クラスタ内の他のノードおよび/またはストレージデバイスに接続される必要があるかを決定することができる。例えば、図1を参照すると、ノード102Bが、クラスタ102に加入し、ノード102Aが、ストレージデバイス104への物理接続を確立した場合、または有する場合、物理パスは、クラスタ内の他のノードにアダプタイズされ、クラスタ内の他のノード、例えば、102Dは、ノード102Aに対するリモートパスを確立し、ノード102Aとストレージデバイス104との間の物理接続を利用することができる。

20

【0037】

[0044] 図4は、本開示の1つまたは複数の実施形態による、クラスタへの加入再承認を要求するための方法400を示している。特定の実施形態では、クラスタへの加入再承認を求めるノードは、それ自体、クラスタから除去されている可能性があり、または図3に関して上記に論じたように、クラスタ内の別のノードによって除去されている可能性もある。諸実施形態では、方法400は、クラスタ内のノードが、ストレージデバイスへの書込みコマンドを送信し、書込みが成功裏でなかったことを通知されるときに開始する。書込みが成功裏でなかった場合、ノードは、コマンドが、異なるパスの下方に送信されることを要求することができる。加えて、または代替として、ノードは、進行中である他のすべての保留になっているコマンドがキャンセルされることも要求することができる。通知を受信すると同時に、ノードは、ストレージデバイスにノード登録テーブルを要求する(410)。上記に論じたように、ノード登録テーブルは、ストレージデバイスによって維持され、クラスタ内のそれぞれのノードと関連している様々な登録キーを含む。

30

【0038】

[0045] ノード登録テーブルが、要求しているノードによって受信されたとき、ノードは、登録テーブルを読み取って(420)、それ自体の登録キーが登録テーブルに含まれているかどうかを決定する。ノードの登録キーが、ノード登録テーブルに含まれていないことが決定されるとき、ノードは、別の登録キーを使用して、ストレージデバイスに登録する(430)。上記に論じたように、登録キーは、クラスタグローバル意識別子の32ビットハッシュ、8ビットキーリビジョン、8ビットノード番号、および16ビットシグネチャを有する64ビット整数とすることができる。特定の実施形態では、ノードは、再登録しなくてはならないとき、新規の登録キーを生成することができ、その場合、登録キーの少なくとも一部分がインクリメントされ、または変更される。したがって、ノード、ストレージデバイス、またはクラスタ内の他のノードは、新規のパスがノードのために設定されなくてはならない回数を追跡することができる。いくつかの実施形態では、ノードが、加入再承認を要求した回数は、クラスタに対する加入再承認を取得するノードに影響をもたらし得る。加えて、登録キーに対する変更は、(ノードがクラスタに加入再承認されるときに)ノードからおよび/またはノードと関連するパスから生じる書込みコマンドが、ストレージデバイスに書き込まれるのをなおも待っている可能性がある、ノードから

40

50

の古くなった書込みコマンドと確実に区別できるようにすることの助けにもなる。

【 0 0 3 9 】

[0046]一旦、ノードが、その新規に生成された登録キーをストレージデバイスに登録すると、登録タイマが設定される(440)。先に論じたように、登録タイマの長さは、時間期間 t の約2.5倍に相当する。つまり、時間期間 t は、クラスタ内の各ノードが、本明細書に開示される防御アルゴリズムを動作させるのにかかる時間に相当する。

【 0 0 4 0 】

[0047]タイマの期限切れと同時に、フローは動作450に進み、ストレージデバイスへの書込みアクセスを求めるノードは、ストレージデバイスによって維持される登録テーブルを読み取る。新規に生成された登録キーが、ノード登録テーブルに存在することが決定される場合(460)、ノードには、ストレージデバイスへの書込みアクセスが付与される。

10

【 0 0 4 1 】

[0048]しかしながら、動作460で、ノードの新規に生成された登録キーが、ノード登録テーブルに存在しないことが決定される場合、フローは動作430に戻って進み、ノードは再登録し、登録タイマはリセットされる。特定の実施形態では、ノードが、決められた回数、拒否された場合、ノードはもはや、ストレージデバイスへの書込みアクセスを求めなくなる。特定の実施形態では、ノードは、決められた時間期間が経過した後に、再度、クラスタに対する加入承認、またはストレージデバイスへのアクセスを求めることができる。

20

【 0 0 4 2 】

[0049]図5は、本開示の1つまたは複数の実施形態による、クラスタ内の2つのノードが、リモートパスおよび/または物理パスを使用して、物理ディスクにそれによってアクセスすることが可能なシステム500を示すブロック図である。図1から図4に関して上に論じたように、クラスタ内の様々なノードは、ストレージプール内の1つまたは複数のストレージデバイスに接続可能である。それらの接続(またはスピンドル)は、物理接続またはリモート接続とすることができる。以下に説明されるように、ノードは、様々なパスを利用して、1つまたは複数のストレージデバイスに接続することができる。

【 0 0 4 3 】

[0050]図5に示される例示的な実施形態では、クラスタは、2つのノード、Node(ノード)Aの510およびNodeBの520を有することができる。各ノードは、物理Disk530への物理接続を有することができる。2つのノードのみが示されているが、クラスタは、3つ以上のノードで構成されてよいことが企図される。加えて、各ノードは、1つまたは複数のストレージデバイスへの物理接続を有することも企図される。図5に示すように、ノードは、いくつかの異なるパスを通じて物理ディスクにアクセスでき得る。例えば、NodeAの510は、物理ディスク530への物理パスを有し、また、NodeBの520のターゲット524を経由して物理ディスク530へのリモートパスも有する。特定の実施形態では、単一のノードが、同じディスクへの複数の物理パスを有することができる。そのような実施形態では、ノードは、クラスタ内の他のノードのすべてに対してすべてのこれらの物理パスを曝すことになる。

30

40

【 0 0 4 4 】

[0051]また、図5に示すように、NodeAは、仮想ディスク511と、物理ディスク530への複数の物理パスおよびリモートパスを統合するマルチパスオブジェクト512と、例えば、NodeBの520など、別のノードを通じて、物理ディスク530への1つまたは複数のリモートパスをインスタンス化するリモートパスオブジェクト513と、NodeAの510と物理ディスク530との間の1つまたは複数の物理接続をアドバタイズするように働き、クラスタ内の他のノードがノードAの510を介して物理ディスク530に対する接続性を(例えば、リモートパスを通じて)得ることを可能にするターゲットオブジェクト514にアクセスでき、物理パスオブジェクト515が、物理ディスク530への1つまたは複数の物理接続もしくはパスをインスタンス化する。

50

【 0 0 4 5 】

[0052]同様に、N o d e Bの5 2 0は、仮想ディスク5 2 1と、N o d e Bの5 2 0から物理ディスク5 3 0への複数の物理パスおよびリモートパスを統合するマルチパスオブジェクト5 2 2と、例えば、N o d e Aの5 1 0など、別のノードを通じて、物理ディスク5 3 0への1つまたは複数のリモートパスをインスタンス化するリモートパスオブジェクト5 2 3その5 1 3と、クラスタ内の他のノードに物理ディスク5 3 0への物理パスをアダプタイズするターゲットオブジェクト5 2 4と、物理ディスク5 3 0への1つまたは複数の物理接続もしくはパスをインスタンス化する物理パスオブジェクト5 2 5にアクセスできる。N o d e Aの5 1 0とN o d e Bの5 2 0の両方に1つのリモートパスが示されているが、単一のノードが複数のリモートパスを有し得ることが企図される。また、各ノードが複数の物理パスを有し得ることも企図される。

10

【 0 0 4 6 】

[0053]諸実施形態では、様々なコマンドが物理ディスク5 3 0に送信されるのに経由する好ましいパスは、物理パスである。例えば、新規のディスクが検出されるとき、クラスタの1つまたは複数のノードは、ディスクを登録し、または予約することになる。上記に論じたように、この処理は、図2に関して上記に論じた防御アルゴリズムを動作させること、続いて、ノードから物理ディスクへの物理パスを創出することを含む。特定の実施形態では、各物理パス、または物理パスの新規の各インスタンスは、クラスタ識別子、ノード識別子、および再生識別子（物理パスがインスタンス化されるたびにインクリメントされる物理パスの固有の番号）を含む登録キーを有する。諸実施形態では、パスの登録キーは、関連ノードの登録キーに相当してよい。一旦、物理接続が確立され、ノードが登録キーを使用してディスクに登録されると、ノードのマルチパスオブジェクトおよびターゲットオブジェクトは、新規に確立された物理パスについて通知される。次いで、その情報が、クラスタ内の他のノードに伝送され、したがって、他のノードは、物理ディスクへの物理接続を有するノードのターゲットを介してリモートパスを確立することができる。

20

【 0 0 4 7 】

[0054]上記に論じたように、1つまたは複数のノードが、クラスタ内の1つまたは複数の他のノードへの、あるいは物理ディスクへの接続を失う場合があることが企図される。そのような事象においては、クラスタ内の接続されたノードのうちの1つは、接続解除されたノードからの1つまたは複数のパスが、除去されることを要求し、また、ストレージデバイスが、接続解除されたノードと関連する1つまたは複数のパス（例えば、物理パスもしくはリモートパス）からの書込み要求を取るのを中止することも要求することになる。同様に、接続解除されたノードへのリモート接続を有する各ノードと関連するターゲットもまた、接続解除されたノードからのコマンドを受け取るのを中止することができる。そのようなアクションにより、接続解除されたノードは、通信線上にはあるが、まだ完了していない可能性があるストレージデバイスに対する追加のおよび/または二重の書込みを送信しないようになる。つまり、登録キーをストレージデバイスから除去すること、およびターゲットを通じて書込みコマンドをブロックすることは、接続解除されたノードが、物理パスまたはリモートパスを使用して、ディスクに書き込むことが確実にできないようにすることの助けになる。

30

40

【 0 0 4 8 】

[0055]例えば、図5を参照すると、N o d e Aの5 1 0は、その物理パスオブジェクト5 1 5を介して物理ディスク5 3 0へのその物理接続を失っている場合がある。しかしながら、図示するように、N o d e Aの5 1 0はまた、N o d e Bの5 2 0のターゲットオブジェクト5 2 4を通じた物理ディスク5 3 0へのリモートパス5 1 3も有する。加えて、物理ディスク5 3 0への接続を失う前に、N o d e Aの5 1 0は、まだ完了していない物理ディスク5 3 0への書込みコマンドを送信していた場合もある。N o d e Aの5 1 0は、物理ディスク5 3 0への接続性を失ったとき、その書込みコマンドが実行されたのか、または拒否されたのかについての知識が全くない場合がある。

【 0 0 4 9 】

50

[0056]しかしながら、Node Aの510が、物理ディスク530とすぐに再接続することが許容され、実行された可能性がある、もしくは実行されていない可能性があるコマンドをいずれも再提出することが許容された場合、またはNode Aの510が、物理ディスク530に対する追加のコマンド(Node Aの510がその接続を失ったことにより、順序から外れてもよい)を送信することが許可された場合、そのようなアクションは、物理ディスク530内のデータを破損させることになる可能性がある。そのような破損を防ぐために、Node Bの520は、Node Aの510と関連する物理パスおよび/またはすべてのリモートパスをプリエンプト(preempt)する。

【0050】

[0057]一旦、Node Aの510と関連する物理パスおよび/またはリモートパスがプリエンプトされると、物理ディスク530は、Node Aの510と関連するパスからのコマンドを受け付けなくなる。各ノードについての各パスが関連識別子を有すると、物理ディスク530は、どのコマンドが、パスのそれぞれの識別子に基づいてノードと関連付けられているかを決定することができる。特定の実施形態では、物理ディスク530は、物理パス間を区別する。したがって、物理ディスクの視点から見れば、I/Oがリモートパスを経由して来た場合、リモートパスが接続されているターゲットをホストするノードからI/Oが来たかのように見えることになる。要するに、リモートパスI/Oフェンシングが、ターゲットにおいて実行される一方、物理パスI/Oフェンシングは、物理ディスク530レベルにおいて実行される。

【0051】

[0058]例を進めるために、クラスタ内の各ノードは、他のあらゆるノードの各スピンドルまたはパスを見ることができる。したがって、Node Bの520は、Node Aの510が、物理ディスク530への接続を失ったことが分かり得る。結果として、Node Bの520は、リモートパス523を破棄することになる。しかしながら、Node Aの510がクラスタ内の他のノードと通信することができない場合、Node Bの520は、物理ディスク530に、Node Aの510からの書き込みコマンドを拒否するように命令することができる。

【0052】

[0059]特定の実施形態では、一旦、物理ディスク530が、Node Aの510の物理パスからのコマンドの拒否を開始すると、Node Aの510のマルチパスオブジェクト512は、コマンドの拒否を検出する。結果として、マルチパスオブジェクト512は、すべての他の既存の物理パスに、いずれかが有効であるかどうかを決定するようにクエリすることができる。1つの物理パスがなおも、有効である場合、有効物理パスは、マルチパスオブジェクト512に追加される。しかしながら、有効な物理パスオブジェクトが全くない場合、新規のマルチパスオブジェクトが創出され、物理パスオブジェクト515は、新規の登録キーを含む新規の物理パスをインスタンス化する。新規の物理パスおよびその関連登録キーは、生成されると、古い物理パスと関連した今や機能していない(defunct)識別子とは別に設定する新規の再生識別子を有することになる。

【0053】

[0060]加えて、例えば、Node Aの510などのノードが、新規の識別子を使用して、クラスタに加入再承認を要求するとき、新規の識別子は、クラスタ内の他のノードにアドバタイズされる。したがって、他のノードのリモートパスオブジェクトは、Node Aの510の物理パスの新規の識別子を使用して、物理ディスク530に接続することができる。上記に論じたように、物理ディスク530は、古い物理パスからのコマンドを受け付けないことを知ると、新規の物理パスからのコマンドおよびその関連識別子を受け付け、そのとき、Node Aの510が、図2～図4に関して上述した方法によりクラスタに加入再承認を求める。

【0054】

[0061]図5に戻って参照すると、Node Aの510およびNode Bの520は、互いへの接続性を失う場合、仮想ディスクに書き込まれていないアプリケーションのキャッ

10

20

30

40

50

シュにデータが存在すること、または物理ディスクに書き込まれていない仮想ディスクにデータが存在することもある。したがって、諸実施形態は、接続解除されたノードからのバス上の残存するすべてのコマンドを消失させる（drain）こと、および接続解除されたノードと関連するバスからのさらなるコマンドが受け付けられないことを実現する。

【0055】

[0062]本明細書に記載の実施形態および機能性は、デスクトップコンピュータシステム、ワイヤードおよびワイヤレスのコンピューティングシステム、モバイルコンピューティングシステム（例えば、モバイル電話、ネットブック、タブレットタイプまたはスレートタイプのコンピュータ、ノートブックコンピュータ、およびラップトップコンピュータ）、ハンドヘルドデバイス、マルチプロセッサシステム、マイクロプロセッサベースのまたはプログラマブルの家庭用電化製品、ミニコンピュータ、ならびにメインフレームコンピュータを限定することなく含めた数多くのコンピューティングシステムによって動作することができる。

【0056】

[0063]加えて、本明細書に記載の実施形態および機能性は、分散されたシステム（例えば、クラウドベースのコンピューティングシステム）を通じて動作することができ、その場合、アプリケーション機能性、メモリ、データの記憶および検索ならびに様々な処理機能は、Internet、またはイントラネットなど、分散されたコンピューティングネットワークを通じて互いとりモートに動作可能である。様々なタイプのユーザインターフェイスおよび情報は、オンボードコンピューティングデバイス表示部を介して、あるいは1つまたは複数のコンピューティングデバイスと関連するリモート表示装置を介して表示可能である。例えば、様々なタイプのユーザインターフェイスおよび情報は、様々なタイプのユーザインターフェイスおよび情報が投影される壁面に表示することができ、その壁面と相互作用することができる。本発明の実施形態がそれにより実践可能な数多くのコンピューティングシステムとの相互作用は、キーストローク入力、タッチ画面入力、音声または他の音響入力、およびジェスチャ入力などを含み、このジェスチャ入力では、関連コンピューティングデバイスに、コンピューティングデバイスの機能性を制御するようにユーザジェスチャを取り込み、解釈するための検出（例えば、カメラ）機能性が備わっている。

【0057】

[0064]図6～図8および関連の記載では、本発明の実施形態がその中で実践可能である種々の動作環境についての議論が行われる。しかしながら、図6～図8に関して図示され、論じられるデバイスおよびシステムは、例示および図示の目的のためであり、本明細書に記載される本発明の実施形態を実践するのに使用可能な膨大な数のコンピューティングデバイス構成を限定するものではない。

【0058】

[0065]図6は、本発明の実施形態がそれにより実践可能なコンピューティングデバイス105の物理的構成要素（すなわち、ハードウェア）を示すブロック図である。後述のコンピューティングデバイスの構成要素は、上記に記載のノードまたはコンピューティングデバイスに適し得る。基本構成では、コンピューティングデバイス105は、少なくとも1つの処理装置602、およびシステムメモリ604を含むことができる。コンピューティングデバイスの構成およびタイプに応じて、システムメモリ604は、これらに限定されないが、揮発性ストレージ（例えば、ランダムアクセスメモリ）、不揮発性ストレージ（例えば、読取り専用メモリ）、フラッシュメモリ、またはそのようなメモリの任意の組合せを含むことができる。システムメモリ604は、オペレーティングシステム605と、ソフトウェアの様々なアプリケーション620を動作させるのに適している1つまたは複数のプログラムモジュール606とを含むことができる。オペレーティングシステム605は、例えば、コンピューティングデバイス105の動作を制御するのに適し得る。さらには、本発明の実施形態は、グラフィックスライブラリ、他のオペレーティングシステム、または任意の他のアプリケーションプログラムと併せて実践可能であり、いずれかの

特定のアプリケーションまたはシステムに限定されない。この基本構成は、図 6 に、破線 608 内のそれらの構成要素によって示されている。コンピューティングデバイス 105 は、追加の特徴または機能性を有することができる。例えば、コンピューティングデバイス 105 はまた、例えば、磁気ディスク、光ディスク、またはテープなどの追加のデータストレージデバイス（リムーバブルおよび/または非リムーバブル）を含むことができる。そのような追加のストレージは、図 6 に、リムーバブルストレージデバイス 609 および非リムーバブルストレージデバイス 610 によって示されている。

【0059】

[0066] 上記のように、いくつかのプログラムモジュールおよびデータファイルが、システムメモリ 604 に記憶され得る。処理装置 602 において実行する一方、プログラムモジュール 606 は、これらに限定されないが、図 1 ~ 図 4 に示す方法の段階のうちの 1 つまたは複数を含む処理を行うことができる。本発明の実施形態により使用可能な他のプログラムモジュールには、電子メール、および連絡先アプリケーション、ワードプロセッシングアプリケーション、スプレッドシートアプリケーション、データベースアプリケーション、スライドプレゼンテーションアプリケーション、描画またはコンピュータ支援アプリケーションのプログラムなどを含めることができる。

【0060】

[0067] さらに、本発明の実施形態は、個別の電子要素、論理ゲートを含んだパッケージ化もしくは集積化された電子チップ、マイクロプロセッサを利用する回路、または電子要素もしくはマイクロプロセッサを含んだ単一のチップを含む電気回路で実践可能である。例えば、本発明の実施形態は、システムオンチップ（system-on-a-chip: SOC）により実践可能であり、その場合、図 6 に示す構成要素のそれぞれ、または多くが、単一の集積回路上に集積化可能である。そのような SOC デバイスは、1 つまたは複数の処理装置、グラフィックス装置、通信装置、システム仮想化装置、および様々なアプリケーション機能性を含むことができ、それらのすべては、単一の集積回路としてチップ基板上に集積化される（または焼き込まれる）。本明細書に記載される機能性は、SOC により動作するとき、単一の集積回路（チップ）上のコンピューティングデバイス 105 の他の構成要素とともに集積化されるアプリケーションに固有の論理機構により動作可能である。本発明の実施形態はまた、これらに限定されないが、機械技術、光技術、流体技術、および量子技術を含んだ、例えば、AND、OR、およびNOTなどの論理演算を行うことができる他の技術を使用しても実践可能である。加えて、本発明の実施形態は、汎用のコンピュータ内で、または任意の他の回路もしくはシステムにおいて実践可能である。

【0061】

[0068] コンピューティングデバイス 105 はまた、キーボード、マウス、ペン、音入力デバイス、タッチ入力デバイスなどの 1 つまたは複数の入力デバイス 612 を有することができる。表示部、スピーカ、プリンタなどの出力デバイス（複数可）614 もまた、含められ得る。上述のデバイスは、例であり、他のデバイスが使用されてもよい。コンピューティングデバイス 104 は、他のコンピューティングデバイス 618 との通信を可能にする 1 つまたは複数の通信接続部 616 を含むことができる。適切な通信接続部 616 の例は、これらに限定されないが、RF 伝送器、受信器、および/または送受信器回路、ユニバーサルシリアルバス（universal serial bus: USB）、パラレルポート、および/またはシリアルポートを含む。

【0062】

[0069] 本明細書に使用される用語コンピュータ可読媒体は、コンピュータストレージ媒体を含むことができる。コンピュータストレージ媒体には、コンピュータ可読命令、データ構造、またはプログラムモジュールなど、情報を記憶するための任意の方法または技術で実装される揮発性と不揮発性の、リムーバブルと非リムーバブルの媒体を含めることができる。システムメモリ 604、リムーバブルストレージデバイス 609、および非リムーバブルストレージデバイス 610 はすべて、コンピュータストレージ媒体の例（すなわ

10

20

30

40

50

ち、メモリストレージ)である。コンピュータストレージ媒体には、RAM、ROM、電氣的に消去可能な読取り専用メモリ(EEPROM)、フラッシュメモリ、もしくは他のメモリ技術、CD-ROM、デジタル多用途ディスク(DVD)、もしくは他の光ストレージ、磁気カセット、磁気テープ、磁気ディスクストレージ、もしくは他の磁気ストレージデバイス、または情報を記憶するのに使用可能であり、コンピューティングデバイス105によってアクセス可能な任意の他の製造品を含めることができる。任意のそのようなコンピュータストレージ媒体は、コンピューティングデバイス105の一部とすることができる。コンピュータストレージ媒体は、搬送波、または他の伝播もしくは変調データ信号を含まない。

【0063】

10

[0070]通信媒体は、コンピュータ可読命令、データ構造、プログラムモジュール、または搬送波もしくは他の輸送機構などの変調データ信号における他のデータによって実現可能であり、任意の情報配信媒体を含む。用語「変調データ信号」は、信号における情報を符号化するような形で、1つまたは複数の特性を設定または変更した信号について説明することができる。限定ではなく、例として、通信媒体は、ワイヤードネットワークまたは直接ワイヤード接続などのワイヤード媒体、ならびに音響、無線周波数(RF)、赤外線、および他のワイヤレス媒体など、ワイヤレス媒体を含むことができる。

【0064】

[0071]図7Aおよび図7Bは、本発明の実施形態がそれにより実践可能なモバイルコンピューティングデバイス700、例えば、モバイル電話、スマートフォン、タブレットパーソナルコンピュータ、およびラップトップコンピュータなどを示している。図7Aに関しては、諸実施形態を実装するためのモバイルコンピューティングデバイス700の1つの実施形態が示されている。基本構成では、モバイルコンピューティングデバイス700は、入力要素と出力要素との両方を有するハンドヘルドコンピュータである。モバイルコンピューティングデバイス700は、典型的には、表示部705と、ユーザがモバイルコンピューティングデバイス700に情報を入力することを可能にする1つまたは複数の入力ボタン710とを含む。モバイルコンピューティングデバイス700の表示部705はまた、入力デバイス(例えば、タッチ画面表示部)として機能することができる。オプションのサイド入力要素715が含まれている場合、このサイド入力要素715は、さらなるユーザ入力を可能にする。サイド入力要素715は、ロータリスイッチ、ボタンまたは任意の他のタイプの手動入力要素とすることができる。代替の実施形態では、モバイルコンピューティングデバイス700は、より多くの入力要素、またはより少ない入力要素を組み込むことができる。例えば、表示部705は、いくつかの実施形態では、タッチ画面でない場合もある。さらなる別の代替の実施形態では、モバイルコンピューティングデバイス700は、セルラ電話など、携帯電話システムである。モバイルコンピューティングデバイス700はまた、オプションのキーパッド735を含むことができる。オプションのキーパッド735は、物理キーパッドであっても、またはタッチ画面表示部において作成される「ソフト」キーパッドであってもよい。様々な実施形態では、出力要素は、グラフィカルユーザインターフェイス(graphical user interface: GUI)を示すために表示部705、視覚指示器720(例えば、発光ダイオード)、および/または音響トランスデューサ(例えば、スピーカ)を含む。いくつかの実施形態では、モバイルコンピューティングデバイス700は、ユーザに触覚フィードバックを提供するための振動トランスデューサを組み込んでいる。さらなる別の実施形態では、モバイルコンピューティングデバイス700は、音響入力(例えば、マイクロフォンジャック)、音響出力(例えば、ヘッドフォンジャック)など、入力ポートおよび/または出力ポート、ならびに外部デバイスに信号を送信し、もしくは外部デバイスから信号を受信するための映像出力(例えば、HDMI(登録商標)ポート)を組み込んでいる。

【0065】

[0072]図7Bは、モバイルコンピューティングデバイスの1つの実施形態のアーキテクチャを示すブロック図である。つまり、モバイルコンピューティングデバイス700は、

50

いくつかの実施形態を実装するシステム（すなわち、アーキテクチャ）702を組み込んでいることができる。1つの実施形態では、システム702は、1つまたは複数のアプリケーション（例えば、ブラウザ、eメール、スケジュール管理（calendar）、連絡先マネージャ、メッセージングクライアント、ゲーム、およびメディアクライアント/プレーヤ）を動作させることができる「スマートフォン」として実装される。いくつかの実施形態では、システム702は、集積化されたパーソナルデジタルアシスタント（personal digital assistant: PDA）とワイヤレス電話など、コンピューティングデバイスとして集積化される。

【0066】

[0073] 1つまたは複数のアプリケーションプログラム766は、メモリ762にロード可能であり、オペレーティングシステム764において、またはオペレーティングシステム764と関連して動作することができる。アプリケーションプログラムの例には、電話ダイヤルプログラム、eメールプログラム、パーソナル情報管理（personal information management: PIM）プログラム、ワードプロセッシングプログラム、スプレッドシートプログラム、Internetブラウザプログラム、およびメッセージングプログラムなどが含まれる。システム702はまた、メモリ762内に不揮発性ストレージ領域768を含む。不揮発性ストレージ領域768は、システム702が電源を落とされた場合、失うべきでない永続的情報を記憶するために使用可能である。アプリケーションプログラム766は、eメールアプリケーションなどによって使用されるeメールまたは他のメッセージなど、不揮発性ストレージ領域768の情報を使用し、記憶することができる。同期化アプリケーション（図示せず）もまた、システム702において常駐し、ホストコンピュータに常駐する対応する同期化アプリケーションと相互作用して、不揮発性ストレージ領域768に記憶された情報がホストコンピュータにおいて記憶された対応する情報と継続して同期化されるようにプログラミングされる。理解すべきであるように、他のアプリケーションは、メモリ762にロード可能であり、モバイルコンピューティングデバイス700において動作することができる。

【0067】

[0074] システム702は、電源770を有し、この電源770は、1つまたは複数のバッテリーとして実装可能である。電源770は、バッテリーを補填する、もしくは再充電するACアダプタまたは電力供給されるドッキングクレードルなど、外部電源をさらに含むこともあり得る。

【0068】

[0075] システム702はまた、無線周波数通信を送受信する機能を行う無線機772も含むことができる。無線機772は、通信キャリアまたはサービスプロバイダを介して、システム702と「外部世界（outside world）」との間のワイヤレス接続性を容易にする。無線機772との間の伝送は、オペレーティングシステム764の制御の下、実施される。言い換えれば、無線機772によって受信される通信は、オペレーティングシステム764を介してアプリケーションプログラム766に伝えられていくことができ、逆も同様である。

【0069】

[0076] 視覚指示器720は、視覚的な通知を提供するために使用可能であり、および/または音響インターフェイス774は、音響トランスデューサ725を介して可聴通知をもたらすために使用可能である。図示の実施形態では、視覚指示器720は、発光ダイオード（LED）であり、音響トランスデューサ725は、スピーカである。これらのデバイスは、電源770に直接、連結可能であり、それにより、プロセッサ760および他の構成要素がバッテリー電力を保存するためにシャットダウンすることがあり得ても、それらのデバイスは、アクティブ状態であるとき、通知機構によって指図された継続時間の間は、オンのままである。LEDは、ユーザがデバイスの電源オン状態を示すアクションを取るまでは無期限に、オンのままであるようにプログラミング可能である。音響インターフェイス774は、ユーザに可聴信号を提供し、ユーザから可聴信号を受信するために使用

される。例えば、音響インターフェイス 774 は、音響トランスデューサ 725 に連結されていることに加えて、電話会話を容易にするなどのために、可聴入力を受信するマイクロフォンにも連結可能である。本発明の実施形態により、マイクロフォンはまた、後述されるように、通知の制御を容易にする音響センサとしても働くことができる。システム 702 は、静止画像、および映像ストリームなどを記録するオンボードカメラ 730 の動作を可能にする映像インターフェイス 776 をさらに含むことができる。

【0070】

[0077] システム 702 を実装するモバイルコンピューティングデバイス 700 は、追加の特徴または機能性を有することができる。例えば、モバイルコンピューティングデバイス 700 はまた、磁気ディスク、光ディスク、またはテープなどの追加のデータストレージデバイス（リムーバブルおよび/または非リムーバブル）も含むことができる。そのような追加のストレージは、不揮発性ストレージ領域 768 によって図 7B に示されている。

10

【0071】

[0078] モバイルコンピューティングデバイス 700 によって生成され、または取り込まれ、システム 702 を介して記憶されるデータ/情報は、上記に記載したように、モバイルコンピューティングデバイス 700 においてローカルに記憶可能であり、すなわち、データは、無線機 772 を介して、またはモバイルコンピューティングデバイス 700 と、モバイルコンピューティングデバイス 700 に関連する別個のコンピューティングデバイス、例えば、Internet など、分散されたコンピューティングネットワークにおけるサーバコンピュータとの間のワイヤード接続を介して、デバイスによってアクセス可能な任意の数のストレージ媒体において記憶可能である。理解すべきであるように、そのようなデータ/情報は、無線機 772 を介して、または分散されたコンピューティングネットワークを介して、モバイルコンピューティングデバイス 700 によりアクセス可能である。同様に、そのようなデータ/情報は、電子メールおよびコラボレーション型のデータ/情報共有システムを含む、よく知られているデータ/情報の転送およびストレージ手段に従って、記憶し、使用するためにコンピューティングデバイス間で容易に転送可能である。

20

【0072】

[0079] 図 8 は、上述したクラスタにおけるメンバシップを提供し、維持するためのシステムのアーキテクチャの 1 つの実施形態を示している。例えば、ノード登録テーブル、識別子、ならびにノード間およびノードと物理ディスクとの間の様々なパスは、種々の通信チャネルまたは他のストレージタイプで記憶され得る。例えば、様々な識別子は、ディレクトリサービス 822、ウェブポータル 824、メールボックスサービス 826、インスタントメッセージングストア 828、またはソーシャルネットワーキングサイト 830 を使用して記憶され得る。サーバ 820 は、データおよび/または接続タイプを、クラスタ内の 1 つまたは複数の他のサーバもしくはノードに提供することができる。1 つの例として、サーバ 820 は、ネットワーク 815 を通じてデータをクライアントへとウェブを介して提供するウェブサーバとすることができる。例として、クライアントコンピューティングデバイスは、コンピューティングデバイス 105 として実装可能であり、パーソナルコンピュータ、タブレットコンピューティングデバイス 610、および/またはモバイルコンピューティングデバイス 700（例えば、スマートフォン）において実現可能である。クライアントコンピューティングデバイス 105、610、700 のこれらの実施形態のいずれもが、ストア 816 からコンテンツを取得することができる。

30

40

【0073】

[0080] 本発明の実施形態は、例えば、本発明の実施形態による方法、システム、およびコンピュータプログラム製品のブロック図および/または動作図を参照して上記に記載されている。ブロックに注記された機能/行為は、任意の流れ図に示される順序から外れて行われてもよい。実際、例えば、連続して示された 2 つのブロックが、実質的に同時に実行されても、またはブロックが、時として、関与している機能性/行為に応じて、逆の順

50

序で実行されてもよい。

【 0 0 7 4 】

【0081】本出願において提供された１つまたは複数の実施形態の説明および例示は、特許請求される本発明の範囲をいかなる形でも限定または制限するように意図するものではない。本出願において提供される実施形態、例、および詳細は、所有権を譲渡し、当業者が、特許請求される発明の最良な様態を作成し、使用することを可能にするのに十分であると見なされる。特許請求される本発明は、本出願において提供されるいかなる実施形態、例、または詳細にも限定するものとは見なすべきではない。組み合わせで、または別個に示され、説明されているかどうかにかかわらず、（構造的にも、方法論的にも）様々な特徴は、特定の特徴の組とともに実施形態を生み出すために選択的に含まれ、または省略されるように意図される。本出願の説明および例示が提供されているので、当業者は、特許請求される本発明のより広範な範囲から逸脱しない、本出願において実現される概括的な発明的概念のより広い態様の趣旨の範囲に入る変形形態、修正形態、および代替の実施形態を想定することができる。

10

【 図 1 】

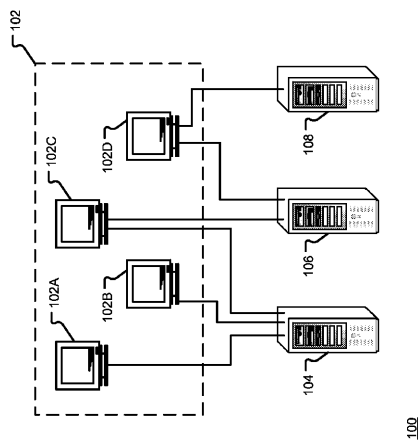
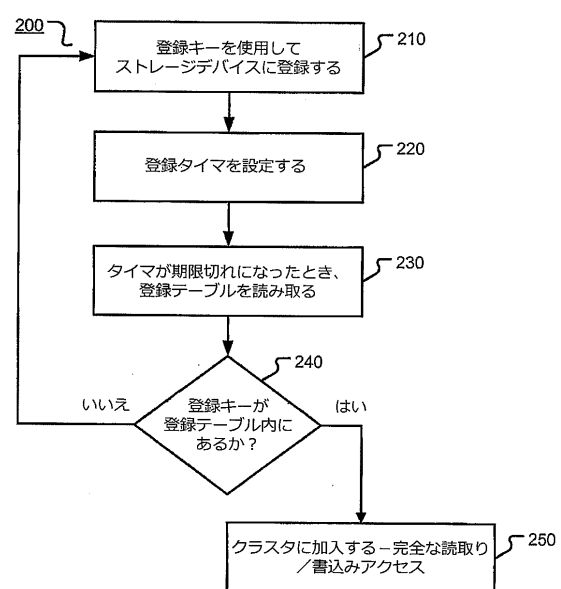
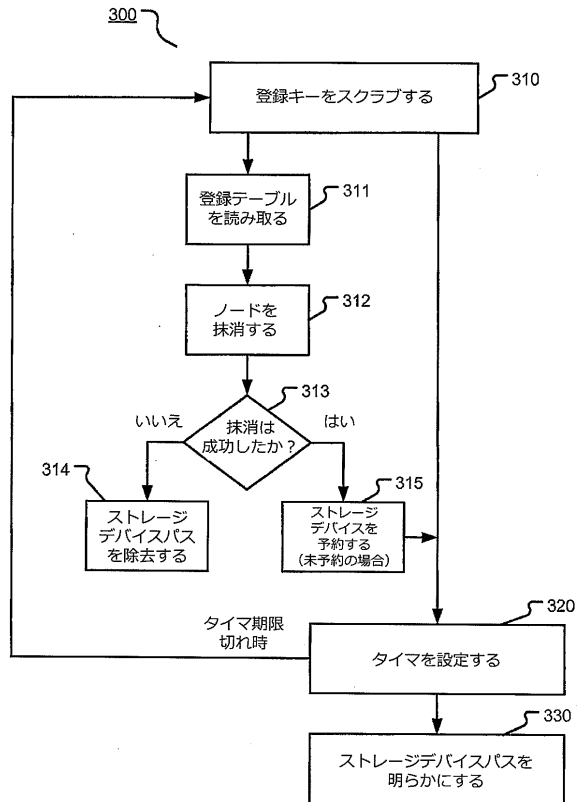


FIG. 1

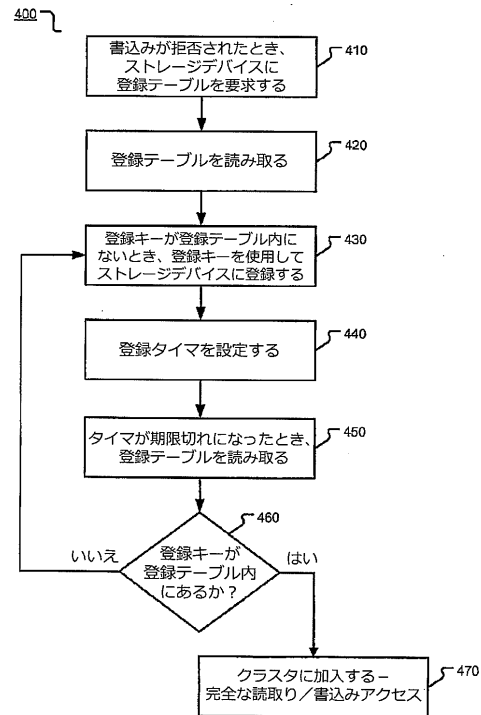
【圖 2】



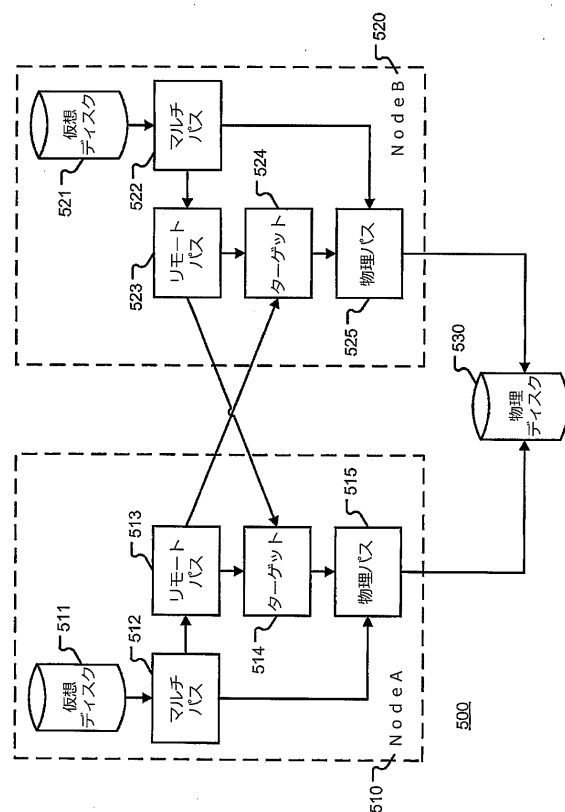
【図 3】



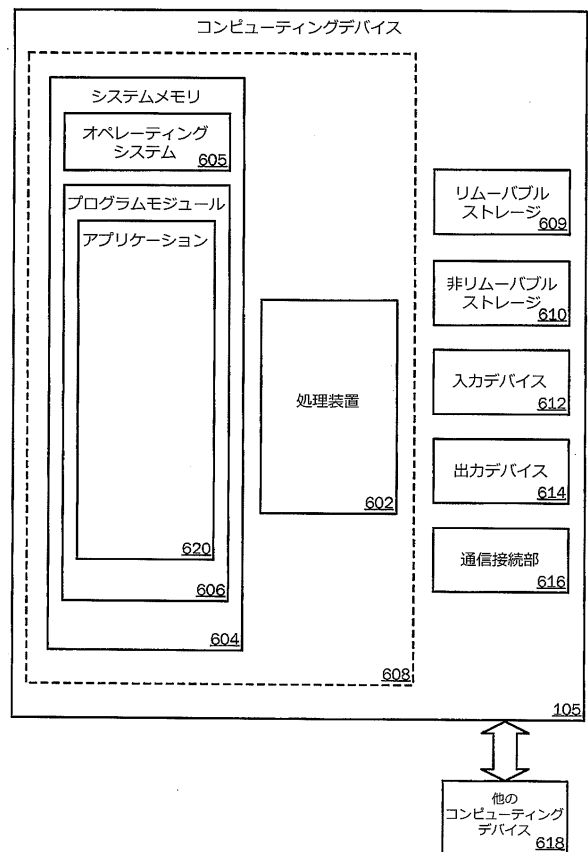
【図 4】



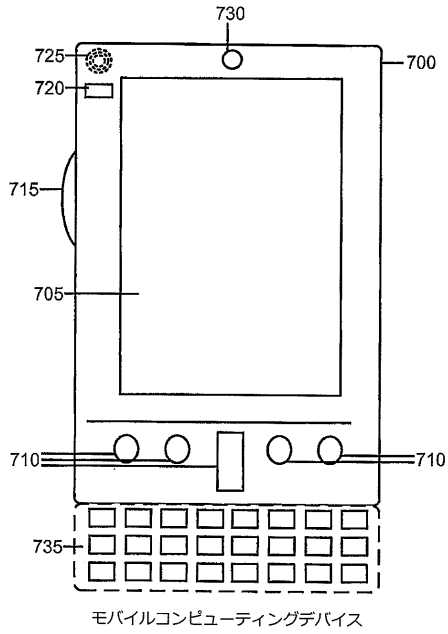
【図 5】



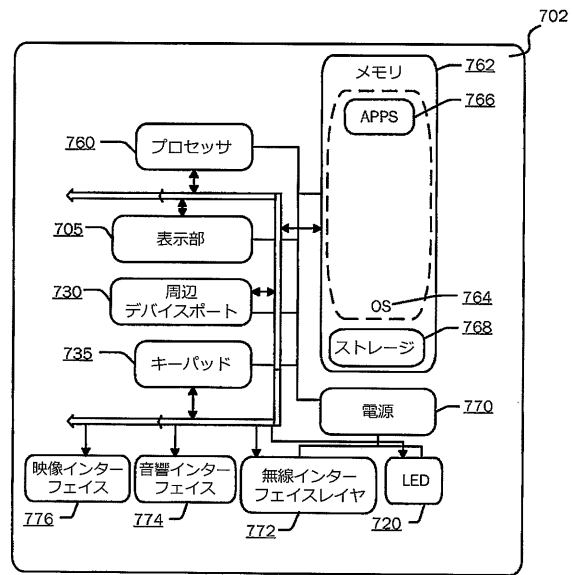
【図 6】



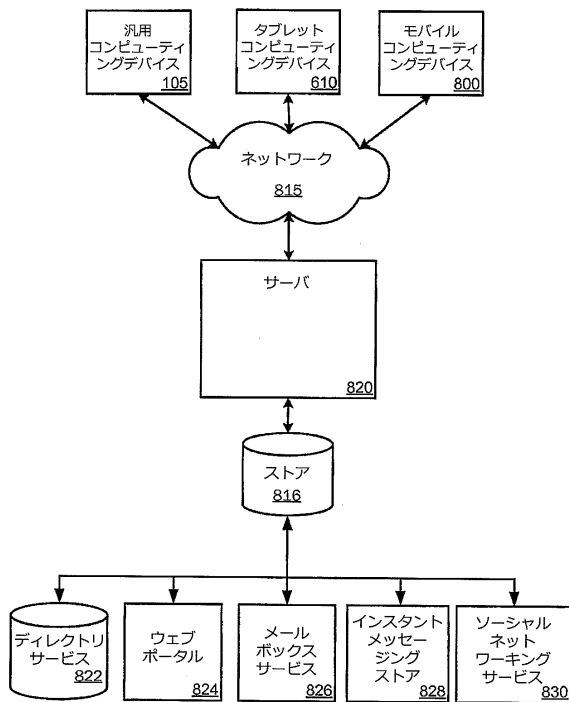
【図 7 A】



【図 7 B】



【図 8】



フロントページの続き

(74)代理人 100162846

弁理士 大牧 綾子

(72)発明者 クズネトソフ, ヴァチェスラフ

アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9, レッドモンド, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテンツ (8 / 1 1 7 2)

(72)発明者 シャンカー, ヴィノッド・アール

アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9, レッドモンド, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテンツ (8 / 1 1 7 2)

(72)発明者 ダマト, アンドレア

アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9, レッドモンド, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテンツ (8 / 1 1 7 2)

(72)発明者 ディオン, デーヴィッド・アレン

アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9, レッドモンド, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテンツ (8 / 1 1 7 2)

審査官 田上 隆一

(56)参考文献 特表 2 0 1 2 - 5 0 3 2 4 9 (J P , A)

米国特許出願公開第 2 0 0 3 / 0 0 6 5 7 8 2 (U S , A 1)

特表 2 0 1 1 - 5 2 6 0 3 8 (J P , A)

(58)調査した分野(Int.Cl., D B 名)

G 0 6 F 3 / 0 6

G 0 6 F 1 3 / 1 0

G 0 6 F 1 3 / 1 4