

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5684110号
(P5684110)

(45) 発行日 平成27年3月11日 (2015. 3. 11)

(24) 登録日 平成27年1月23日 (2015. 1. 23)

(51) Int. Cl. F I
 HO 4 L 12/701 (2013. 01) HO 4 L 12/701
 HO 4 L 12/751 (2013. 01) HO 4 L 12/751

請求項の数 12 (全 14 頁)

| | | | |
|---------------|-------------------------------|-----------|----------------------------------|
| (21) 出願番号 | 特願2011-509861 (P2011-509861) | (73) 特許権者 | 598036300 |
| (86) (22) 出願日 | 平成20年5月23日 (2008. 5. 23) | | テレフオンアクチーボラゲット エル エム エリクソン (パブル) |
| (65) 公表番号 | 特表2011-521573 (P2011-521573A) | | スウェーデン国 ストックホルム エスー |
| (43) 公表日 | 平成23年7月21日 (2011. 7. 21) | | 1 6 4 8 3 |
| (86) 国際出願番号 | PCT/EP2008/056376 | (74) 代理人 | 100076428 |
| (87) 国際公開番号 | W02009/141013 | | 弁理士 大塚 康德 |
| (87) 国際公開日 | 平成21年11月26日 (2009. 11. 26) | (74) 代理人 | 100112508 |
| 審査請求日 | 平成23年4月22日 (2011. 4. 22) | | 弁理士 高柳 司郎 |
| | | (74) 代理人 | 100115071 |
| | | | 弁理士 大塚 康弘 |
| | | (74) 代理人 | 100116894 |
| | | | 弁理士 木村 秀二 |
| | | (74) 代理人 | 100130409 |
| | | | 弁理士 下山 治 |

最終頁に続く

(54) 【発明の名称】 ルーティングテーブルを維持する方法およびオーバーレイネットワーク内で使用するためのノード

(57) 【特許請求の範囲】

【請求項 1】

C h o r d 分散型ハッシュテーブル (D H T) ベースのオーバーレイネットワークのノードにおけるルーティングテーブルを維持する方法であって、所与のノードが有するルーティングテーブルは、近隣のサクセッサノード及びプレディセッサノードのセットの各々に対して、ノードのオーバーレイネットワークアドレスと該ノードの物理的口ケータとの間のマッピングを含み、

前記ノードが前記オーバーレイネットワークに新たに参加した他のノードを知ることができるように、前記オーバーレイネットワークのノード間で D H T 保守メッセージを周期的に交換し、

前記オーバーレイネットワークからのノードの撤退に際して、前記撤退するノードの近隣ノードのうちの前記撤退を検知した一つのノードから、前記撤退するノードの他の近隣ノードの各々に、前記撤退を示すとともに受信ノードのルーティングテーブルに含まれていないノードに対する 1 つ以上のマッピングを含む離脱要求を送信し、

前記他の近隣ノードの各々において離脱要求を受信し、前記受信ノードにおいて前記ルーティングテーブルを更新するために前記マッピングを利用することを特徴とする方法。

【請求項 2】

前記離脱要求を送信する前記ノードは、周期的に送信されるキープアライブメッセージに対する応答が前記撤退するノードから送信されないことの結果として前記撤退するノードの撤退を検出し、前記離脱要求を送信することにより前記離脱要求に対して反応するこ

とを特徴とする請求項 1 に記載の方法。

【請求項 3】

前記受信ノードの前記ルーティングテーブル内に含まれていないノードに対する前記 1 つ以上のマッピングは、前記離脱要求を送信する前記ノードの前記ルーティングテーブル内に含まれているマッピングであることを特徴とする請求項 1 または 2 に記載の方法。

【請求項 4】

離脱要求を受信するノードにおいて、前記受信ノードが前記撤退するノードの近隣ノードのうち前記離脱要求を送信したノードにとって未知の近隣ノードを知っているかどうかを判断し、知っているとは判断された場合には離脱要求を 1 つ以上の前記近隣ノードに送信し、前記離脱要求は、受信ノードにおけるルーティングテーブル内に含まれていない少なくとも 1 つのノードに対する物理的ロケータマッピングへの 1 つ以上のオーバーレイネットワークアドレスを含み、前記受信ノード或いはその他の受信ノードにおいて前記離脱要求を受信すると、前記ノードの前記ルーティングテーブルを更新することを特徴とする請求項 1 乃至 3 のいずれか 1 項に記載の方法。

10

【請求項 5】

前記オーバーレイネットワークのノードにおいて、ルーティングテーブルに含まれていないノードに対する 1 つ以上のマッピングをキャッシュし、前記オーバーレイネットワークからのノードの撤退が発生した場合に、1 つ以上の前記キャッシュされたマッピングを含めるように、前記撤退するノードの少なくとも一つの近隣ノードのルーティングテーブルを更新することを特徴とする請求項 1 乃至 4 のいずれか 1 項に記載の方法。

20

【請求項 6】

前記離脱要求は、前記受信ノードに対する新たなサクセッサおよびプレディセッサルーティングテーブルの一つを含むことを特徴とする請求項 1 乃至 5 のいずれか 1 項に記載の方法。

【請求項 7】

前記オーバーレイネットワークは、分散型ハッシュテーブルネットワークであり、前記オーバーレイネットワークアドレスはハッシュ値であることを特徴とする請求項 1 乃至 6 のいずれか 1 項に記載の方法。

【請求項 8】

前記受信ノードにおいて前記ルーティングテーブルを更新するステップは、前記撤退するノードに対応するマッピングを削除することと、前記受信された離脱要求に含まれている新たなマッピングを前記ルーティングテーブルに加えることとを含むことを特徴とする請求項 1 乃至 7 のいずれか 1 項に記載の方法。

30

【請求項 9】

Chord 分散型ハッシュテーブル (DHT) ベースのオーバーレイネットワーク内で使用するためのノードであって、

近隣のサクセッサノード及びプレディセッサノードのセットの各々に対する、前記ノードのオーバーレイネットワークアドレスと前記ノードの物理的ロケータとの間のマッピングを含む、ルーティングテーブルを格納するメモリと、

前記ノードが前記オーバーレイネットワークに新たに参加した他のノードを知ることができるように、前記オーバーレイネットワークの他のノードとの間で DHT 保守メッセージを周期的に交換する処理ユニットを備え、

40

前記処理ユニットは、前記オーバーレイネットワークからの近隣ノードの撤退に際して、前記撤退するノードを特定するとともに受信ノードのルーティングテーブル内に含まれていない少なくとも一つのノードに対するマッピングを含む離脱要求を前記ノードの 1 つ以上の近隣ノードに送信する処理ユニットと、を備えることを特徴とするノード。

【請求項 10】

近隣ノードから離脱要求が受信された場合に、前記メモリに格納されている前記ルーティングテーブルから前記離脱要求において特定されている撤退するノードに対応するマッピングを削除し、前記離脱要求に含まれている 1 つ以上の新しいマッピングを前記ルーテ

50

ィングテーブルに追加する、さらなる処理ユニットを備えることを特徴とする請求項9に記載のノード。

【請求項11】

離脱要求が受信された場合に、前記離脱要求が前記撤退するノードから発生したものかどうかを判定し、そうでない場合に、前記ルーティングテーブルを検査して前記撤退するノードの近隣ノードのうち前記離脱要求を送信しているノードが知らないノードを特定し、該特定されたノードに対して、前記離脱するノードを特定するとともに受信ノードのルーティングテーブル内に含まれていない少なくとも1つのノードに対する1つ以上のマッピングを含む離脱要求を送信する、さらなる処理ユニットを備えることを特徴とする請求項9または10に記載のノード。

10

【請求項12】

前記ルーティングテーブル内に含まれていないノードに対する1つ以上のマッピングをキャッシュするさらなるメモリと、前記オーバーレイネットワークから隣接ノードの撤退が発生した場合に、1つ以上の前記キャッシュされたマッピングを含むように前記ルーティングテーブルを更新するさらなるプロセッサと、を備えることを特徴とする請求項9乃至11のいずれか1項に記載のノード。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、オーバーレイネットワークにおける分散ハッシュテーブルを維持するためのメカニズムに関する。本発明は特に、オーバーレイネットワークからのノードの離脱 (leaving) を取り扱うための最適化手順に適用できる。

20

【背景技術】

【0002】

ピアツーピアすなわちP2Pのネットワークは、ファイル共有およびVoIP電話を含む多種多様なサービスを容易にするために、処理能力および通信帯域を含めて、加入ノードの共同利用リソースを用いる。セントラルサーバがない場合、特定のP2Pサービスは、リソース位置を最適化するために“オーバーレイネットワーク”を用いる場合がある。オーバーレイネットワークは、下位層のネットワーク(たとえば、インターネット)における可能な限り多くの物理リンクにわたるパスを表わす仮想リンクにより接続されるノードを備える。オーバーレイネットワークにおける各ノードは、オーバーレイネットワーク内の他の特定ノードへの一連のリンクが入っているルーティングテーブルを維持する。リソース要求は、リソース要求が、そのリソースに係るノードに到達するまでノード間を通過する。

30

【発明の概要】

【発明が解決しようとする課題】

【0003】

分散ハッシュテーブル(DHT)は、リソース名(“キー”)をオーバーレイネットワーク内の位置に対応付けるための効率的な手段を提供する。DHTは、キー、たとえば曲名、SIP URIs等を有限の値空間、たとえば128ビットに対応付けるためにハッシングアルゴリズムを用いる。ハッシングアルゴリズムは、値空間にわたってハッシュ値が比較的一様に広がることを保証するように選ばれる。このようにして、たとえば、100個の曲名のハッシングは同様に、値空間にわたって比較的等間隔である100個のハッシュ値をもたらす。オーバーレイネットワーク内のノードは、ユーザ名で識別され、ユーザ名はそれ自身がそれぞれのハッシュ値に変換される。各ノードはそれから、自身の値に隣接する、値空間内の一連のハッシュ値に係るようになる。実際、ノードは、自身が“所有する”リソース名に一致するリソースが得られることのできる位置(たとえば、IPアドレス)を記憶するであろう。オーバーレイネットワークにおけるノードがリソース要求を受信する場合、ノードは対応するハッシュ値を所有するかどうかを決定する。ハッシュ値を所有すると、ノードはリソースの位置を要求元に(オーバーレイネットワークを経て)返す。ハッシュ値を所有していないと、ノードは、リソース要求のハッシュ値に最も近いハッシ

40

50

ハッシュ値を有するハッシュテーブル内の該当ノードを識別するために、自身のルーティングテーブルを調べ、そしてリソース要求を当該ノードに送る。受信ノードは、リソース要求に対応するハッシュ値を所有しており、したがってリソース位置を知っているノードにリソース要求が達するまで、同手順を繰り返す。

【0004】

図1は、リング（リング内に少数のノードのみが図示されている）状に編成されているオーバーレイネットワークを説明している。この例では、各ノードはリング中に少数の後続ノードおよび先行ノードばかりでなく、少数のより離れたノードの位置およびハッシュ値が入っているルーティングテーブルを維持する。図示しているネットワークでは、ノードXは、ルーティングテーブル内に2つのサクセッサ（successor）ノードおよび2つのプレディセッサ（predecessor）ノードに対する位置だけでなく、3つの離れたノードに対する位置を維持する。ルーティングテーブル内の多数のエントリが、ネットワークをルーティングに関してより効率的に、そしてノードの脱退（withdrawal）に対してよりロバストにすることができるが、大きなテーブルは維持するのが難しく、そしてしたがって、ネットワークの信頼性を増大させる。

10

【0005】

オーバーレイネットワーク内のノードは、近隣と周期的に連絡をとるように努めることにより、ルーティングテーブルにおける情報が最新であることを保証する。いくつかの異なるメカニズムがこの目的のために用いられる場合がある。

【0006】

1) ノードは、ルーティングテーブルに挙げられている他のノードがオーバーレイネットワークを離脱していないことを調べるために、キープアライブメッセージを周期的に送ることができる。このメカニズムは、DHT手法、たとえばPastry [A. RowstronおよびP. Druschel: Pastry: 大規模ピアツーピアシステムのための、拡張可能な分散オブジェクト位置およびルーティング、ミドルウェア、2001年]、Chord [I. Stoica, R. Morris, D. Karger, M. F. KaashoekおよびH. Balakrishnan: Chord: インターネットアプリケーションのための拡張可能なピアツーピア検索サービス、ACM SIGCOMM '01会議の会議録、2001年8月、サンディエゴ、カリフォルニア、USA] および内容でアドレスを指定できるネットワーク (CAN) [S. Ratsanamy, P. Francis, M. Handly, R. KarpおよびS. Shenker: 拡張可能な内容でアドレスを指定できるネットワーク、ACM SIGCOMM 2001、2001年8月] により用いられる。

20

30

【0007】

2) ノードは、旧エントリに置き換えて、ルーティングテーブルに挿入でき得る新しいノードについて学習するためにクエリを周期的に送ることができる（たとえば、Chord）。

【0008】

3) ノードは、近隣のルーティングテーブルにおけるエントリについての情報を要求するクエリを、直通の近隣（direct neighbours）に周期的に送ることができる。この情報は、ノード自身のルーティングテーブルを更新するのに用いられる（たとえば、Chord）。

40

【0009】

4) ノードは、自身のルーティングテーブルを近隣に周期的に送ることができる（たとえば、CAN）。

【0010】

ルーティングテーブルを維持するための別の（付加的な）手法は、リソース要求の発信元がルーティングテーブルに挿入されうるかどうかをノードが調べる工程を伴う（たとえば、Kademlia [P. MaymounkovおよびD. Mazieres: Kademlia: 排他的論理和尺度に基づくピアツーピア情報システム、IPTPS 02の会

50

議録、ケンブリッジ、USA, 2002年3月)。

【0011】

DHTにおける近隣関係の例を示す図2を考える。図では、リングトポロジが仮定されている。ノードXは、ルーティングテーブルに3つのサクセッサポイントおよび3つのプレディセッサポイントを維持する。複数のサクセッサポイントおよびプレディセッサポイントを維持する理由がロバスト性を増大させることであることは、既に明瞭であるだろう。単一のサクセッサが停止するようになる確率が p であると、その場合、すべての3つのサクセッサが同時に停止するようになる確率は p^3 である。しかしながら、極めて大きな実世界のDHTに基づくオーバーレイネットワークでは、このことはネットワークにおける接続性を維持するには十分でなく、所与のノードのすべての3つのサクセッサ(またはその代わりに、すべての3つのプレディセッサ)が十分短い期間内にネットワークを離脱すると、ネットワークは寸断する。

10

【0012】

ノードは、手順を踏んで(*gracefully*)または手順を踏まずに(*ungracefully*)ネットワークを離脱できる。ノードは、手順を踏んで離脱する場合、実際に離脱するに先立って、ネットワークを離脱しようとする意思について近隣に知らせる。ノードは、(アプリケーション層で解釈される)離脱メッセージを送ることによりこのことを行う。これは、近隣が自身のルーティングテーブルから離脱するノードを直ちに除去できるようにする。ノードが手順を踏まずにネットワークを離脱する場合、最初に近隣に知らせることなく、ネットワークを去る。したがって、近隣はノードが離脱していることを自ら検出しなければならない。手順を踏まずに離脱する理由には以下がある。(i)ノードが機能停止になった、(ii)P2Pアプリケーションが機能停止になった、または不意に休止した、そして(iii)利己的な振舞い。(iii)の代わりに、ユーザが、手順を踏んだ撤退(*departure*)につきものの遅延を回避するために手順を踏まずに離脱することを選ぶ場合がある。

20

【0013】

オーバーレイネットワークから手順を踏まずに撤退する場合、ノードは近隣が2つの異なる方法で離脱していることを知ることができる。

【0014】

1) 下位層のトランスポートプロトコルが信頼できる場合(たとえば、TCP)、近隣の撤退は、トランスポート層の接続がダウンするという事実から迅速に検出される。

30

【0015】

2) トランスポートプロトコルが信頼できない場合(たとえば、UDP)、ノードは、次周期のDHT保守メッセージを近隣に送ろうと試みるまで、近隣が離脱していることを学習しない。次周期の保守メッセージ送信を待つことに加えて、ノードはまた、近隣が実際に離脱しているということが確実となりうる前に、トランザクションが時間切れするまで待たなければならない。

【0016】

手順を踏んだ撤退と手順を踏まない撤退との双方の場合、最終結果は、離脱ノードの直通の近隣の各々がルーティングテーブルで1つ少ないポイントを有することである。たとえば、図2におけるノードS1がネットワークを離脱すると、ノードXは、2つのサクセッサポイント、すなわちS2およびS3を残しているだけである。また、ノードXが追加のサクセッサを見つける機会を得る前に、S2およびS3がまたネットワークを離脱すると、ノードXはもはやサクセッサを知らないの、オーバーレイネットワークは分割されるようになる。同じ状況がノードXのプレディセッサノードについて生じうる。

40

【0017】

このようにして、現存する解決策に関連する問題は、所与のノードのサクセッサまたはプレディセッサのすべてが短期間のうちに停止すると、ネットワークは分割されるようになりうるし、そしてリソース要求はギャップを埋めることができないということである。この“短期間”は、2つの連続するDHT保守メッセージ間の時間に当てはまる。そのような保守メッセージがたとえば、60秒ごとに送られるとした場合、単一のノードのすべ

50

てのサクセッサまたはプレディセッサが、この60秒の期間内にオーバーレイネットワークを離脱すると、オーバーレイネットワークは崩壊する。このことは、ネットワークが高い“解約(churn)”率に直面している場合、起こり得そうな事象である。この問題に対する直感的な解決策は、DHT保守メッセージをより頻繁に送ることであろうが、周期的な保守メッセージの送信間隔は、もたらされる信号負荷がネットワークを過負荷にするようになるので、任意には小さくできない。この問題は、S. Rhea、D. Geels、T. RoscoeおよびJ. Kubiatowicz: DHTにおける解約の取り扱い、USENIXの会議録、年次技術会議、2004年6月、により確認されている。

【課題を解決するための手段】

【0018】

本発明の目的は、ノードが停止する、またはそうでなければネットワークから脱退する場合、オーバーレイネットワークに対して混乱を最小にすることである。少なくとも、本発明の特定の実施形態は、脱退するノード、または脱退するノードの近隣ノードが他の近隣のルーティングテーブルを更新できるようにすることにより、この目的を達成する。

【0019】

本発明の第1の態様によれば、オーバーレイネットワークのノードでルーティングテーブルを維持する方法が提供され、所与ノードのルーティングテーブルには、一連の近隣サクセッサノードおよび近隣プレディセッサノードのうちの各々に対して、ノードのオーバーレイネットワークアドレスとノードの物理的ロケータとの間のマッピングが入っている。

【0020】

本方法は、ノードがオーバーレイネットワークから撤退する時点でまたはその直前に、撤退ノード(または撤退を承知している、撤退ノードの近隣ノードのうちの1つ)から、各近隣ノード(または撤退ノードの他の各近隣ノード)に離脱要求を送る工程、撤退を表示する工程および受信ノードのルーティングテーブルに入っていないノードに対する1つ以上のマッピングを入れる工程を備える。各近隣ノード(または他の各近隣ノード)が離脱要求を受信し、そして前記(複数の)マッピングを用いてルーティングテーブルを更新する。

【0021】

本発明の実施形態では、脱退ノードの近隣が自身のルーティングテーブルを置換プレディセッサノードまたは置換サクセッサノードで迅速に更新できるようにする。ネットワークにおいて高い解約である場合において、ネットワーク鎖における破断の危険性が大いに減少する。

【0022】

本発明の第2の態様によれば、オーバーレイネットワークのノードでルーティングテーブルを維持する方法が提供され、所与ノードのルーティングテーブルには、一連の近隣サクセッサノードおよび近隣プレディセッサノードのうちの各々に対して、ノードのオーバーレイネットワークアドレスとノードの物理的ロケータとの間のマッピングが入っている。

【0023】

本方法では、ノードがオーバーレイネットワークを撤退する直前に、撤退ノードから撤退ノードの各近隣ノードに離脱要求を送信し、撤退を示し、そして、受信ノードのルーティングテーブル内に入っていないノードに対する物理的ロケータマッピングへの1つ以上のオーバーレイネットワークアドレスを入れることを備える。各近隣ノードで離脱要求を受信する時点で、ノードは(複数の)マッピングを用いて自身のルーティングテーブルを更新する。

【0024】

本発明のこの態様の実施形態では、離脱要求内に入っている少なくとも1つのマッピングは、受信ノードには未知の、撤退ノードの近隣ノードに対応する。

【0025】

本発明の第3の態様によれば、オーバーレイネットワークのノードでルーティングテーブルを維持する方法が提供され、所与ノードのルーティングテーブルには、一連の近隣サク

10

20

30

40

50

セッサノードおよび近隣プレディセッサノードのうちの各々に対して、ノードのオーバーレイネットワークアドレスとノードの物理的ロケータとの間のマッピングが入っている。

【0026】

本方法では、ノードがオーバーレイネットワークから撤退する時点で、撤退を承知している、撤退ノードの近隣ノードのうちの1つから、撤退ノードの他の近隣ノードに離脱要求を送信し、ここで離脱要求は撤退を示すとともに受信ノードのルーティングテーブルに入っていないノードに対する1つ以上のマッピングを含んでいる。前記他の各近隣ノードで離脱要求を受信すると、そのノードはその(複数の)マッピングを用いて自身のルーティングテーブルを更新する。

【0027】

本発明のこの態様の実施形態では、(複数の)離脱要求を送るノードが、周期的に送信されるキープアライブメッセージに対して撤退ノードが応答しなかったという結果から撤退ノードの撤退を検出する。離脱要求を送るノードは、(複数の)離脱要求を送ることによりその検出への反応を示す。受信ノードのルーティングテーブルに入っていないノードに対する前記1つ以上のマッピングは、離脱要求を送るノードのルーティングテーブルに入っているマッピングであってもよい。

【0028】

離脱要求を受信するノードは、その受信ノードが撤退ノードの近隣ノードのうちの、該離脱要求を送っているノードが承知していない近隣ノードを承知しているかどうかを判断するようにしてもよい。その受信ノードが、離脱要求を送っているノードが承知していない近隣ノードを承知している場合、離脱要求を送っているノードに対して、離脱要求を送信する。ここで送信される離脱要求は、その受信ノードのルーティングテーブルに入っていないノードに対する物理的ロケータマッピングへの1つ以上のオーバーレイネットワークアドレスを含んでいる。受信ノードまたはさらなる各受信ノードにおいて、離脱要求を受信する時点で、ノードは自身のルーティングテーブルが更新される。

【0029】

当然のことながら、受信ノードでルーティングテーブルを更新するステップは、撤退ノードに対応するマッピングを削除すること、および受信した離脱要求に入っている新しいマッピングをルーティングテーブルに追加することを含んでもよい。

【0030】

本発明の実施形態に都合よく組み込み可能な機能は、それぞれのルーティングテーブルに含まれていないノードに対する1つ以上のマッピングをキャッシュすることである。ネットワークからノードが撤退する場合に、撤退するノードの少なくとも1つの近隣ノードのルーティングテーブルは1つ以上のキャッシュされたマッピングを含むように更新できる。

【0031】

本発明の第4の態様によれば、オーバーレイネットワーク内で用い、そして一連の近隣サクセッサノードおよび近隣プレディセッサノードの各々に対して、ノードのオーバーレイネットワークアドレスとノードの物理的ロケータとの間のマッピングが入っているルーティングテーブルを記憶するための記憶装置を備えるノードが提供される。ノードはまた、ノードまたは近隣ノードがネットワークから撤退する時点で、ノードの1つのまたは複数の近隣ノードに離脱要求を送るように構成された処理部を備え、離脱要求は撤退するノードを識別し、そして受信ノードのルーティングテーブルに入っていないノードに対する物理的ロケータマッピングへの1つ以上のオーバーレイネットワークアドレスが入っている。

【0032】

ノードは、近隣ノードから離脱要求を受信し、上記記憶装置に入っているルーティングテーブルから離脱要求で識別されている撤退ノードに対応するマッピングを削除し、そして離脱要求に入っている1つ以上の新しいマッピングをルーティングテーブルに追加するように構成されるさらなる処理ユニットを備えるようにしてもよい。また、さらなる処理ユニットが、離脱要求が撤退するノードから発しているかどうかを決定し、撤退ノードか

10

20

30

40

50

ら発していない場合には、撤退ノードのすべての近隣ノードのうちの上記離脱要求を送ったノードが承知していないノードを識別するためにルーティングテーブルを調べ、そしてそのような識別されたあらゆるノードに対して離脱要求を送るように構成されてもよい。この離脱要求には、撤退ノードを識別するとともに、受信ノードのルーティングテーブルに入っていないノードに対する1つ以上のマッピングが入っている。

【0033】

ノードは、近隣ノードがネットワークから撤退する場合に、ルーティングテーブルに含まれていないノードに対する1つ以上のマッピングをキャッシュするためのさらなる記憶装置、およびキャッシュされた1つ以上のマッピングを含むようにルーティングテーブルを更新するためのさらなるプロセッサを備え得る。

10

【0034】

本発明の第5の態様によれば、オーバーレイネットワークのノードでルーティングテーブルを維持する方法が提供され、所与のノードのルーティングテーブルには、一連の近隣サクセッサノードおよび近隣プレディセッサノードの各々に対して、ノードのオーバーレイネットワークアドレスとノードの物理的ロケータとの間のマッピングが入っている。

【0035】

本方法は、ノードに更新されたアドレス指定情報を提供するためにそれらノード間で保守メッセージを周期的に交換することを含む工程を。アドレス指定情報がピアノードに対して所与ノードで受信され、そしてかかるピアノードが所与ノードのルーティングテーブルに含まれていない場合、情報は所与ノードでキャッシュされる。所与ノードのルーティングテーブルに入っているノードがネットワークから脱退する場合、ピアノードはキャッシュされた情報を用いてルーティングテーブルに追加される。

20

【図面の簡単な説明】

【0036】

【図1】いくつかのノードを備えているDHTに基づくリング状のオーバーレイネットワークを概略的に説明する図である。

【図2】さらに、ノード間の近隣関係を示す、DHTに基づくリング状のオーバーレイネットワークを説明する図である。

【図3】図2のオーバーレイネットワークのノードを概略的に説明する図である。

【図4】図2のオーバーレイネットワークからの、手順を踏んだノードの撤退を取り扱うためのメカニズムを説明するフロー図である。

30

【図5】図2のオーバーレイネットワークからの、手順を踏まないノードの撤退を取り扱うためのメカニズムを説明するフロー図である。

【発明を実施するための形態】

【0037】

本明細書で説明している最適化DHT離脱操作は、近隣ノードのルーティングテーブルが、あるノードがネットワークから撤退することにより影響を受ける場合に、近隣ノードが自身のルーティングテーブルを迅速に更新するのを手助けするノードに依存する。撤退するノードは、たとえば手順を踏んだ撤退シナリオでの近隣ノード、または手順を踏まない撤退の場合での別の近隣ノードである場合がある。これら2つのシナリオについてこれから詳細に考察することとする。

40

【0038】

リングトポロジを用いるChord DHTに基づくオーバーレイネットワークの例を示す図2を再度参照する。Chord DHTは、本明細書では例として用いられているが、説明する手順は他のDHTに基づくオーバーレイネットワークに同様に適用できる。説明している例では、DHTに基づくオーバーレイネットワークにおける各ノードは、6つの近隣、すなわち、3つのプレディセッサノードおよび3つのサクセッサノードへのポイントを維持するという仮定をしている。勿論、提案メカニズムは任意の数のサクセッサポイントおよびプレディセッサポイントとともに機能する。

【0039】

50

図2で、ノードXは3つのサクセッサ、S1、S2、およびS3を有する。オーバーレイネットワークが重大な“解約”に遭遇していると、ノードXのサクセッサのすべてが短時間フレームのうちにオーバーレイネットワークを離脱することを選ぶ可能性がある。ノードが手順を踏んで撤退することを仮定すると、ノードは自身の近隣ノードに離脱要求を送るであろう。しかしながら、上述したように、ノードXは、すべての3つのサクセッサがネットワークを離脱してしまう前に、他のサクセッサノードのどれをも識別する時間を有しない場合がある。このシナリオで、ネットワークの分割を回避するために、オーバーレイを離脱しようとする各ノードは、実際にネットワークを離脱する前に、近隣が自身のルーティングテーブルを代替の近隣ノードで埋めるのに協力する。

【0040】

以下の本文で、用語“近隣テーブル”が、直接につながる近隣へのポイントが入っているルーティングテーブルの一部を引用するのに用いられ、用語“プレディセッサテーブル”が、プレディセッサポイントが入っている近隣テーブルの一部を引用するのに用いられ、そして用語“サクセッサテーブル”が、サクセッサポイントが入っている近隣テーブルの一部を引用するのに用いられる。

【0041】

図2のノードS1が、手順を踏んでオーバーレイネットワークを離脱することを選ぶ場合を考える。ネットワークでのすべての他のノードと同じように、S1は自身のルーティングテーブルに近隣ノードへのポイントを維持している。図2では、ノードS1が3つのプレディセッサポイントおよび3つのサクセッサポイントを維持することを仮定している。ノードS1のサクセッサはS2、S3およびAである。ノードS1のプレディセッサは、X、P1およびP2を含む。ネットワークを離脱する前に、ノードS1は近隣のノードNの各々に対して以下の手順を繰り返す。

【0042】

ノードNがノードS1のプレディセッサであると、ノードS1は、Nに対する新しいサクセッサテーブルを構築し、そしてサクセッサテーブルをノードNに送られる離脱メッセージに含める。ノードS1は、ノードNに対して創出するサクセッサテーブルに自身を含めるべきではない。サクセッサテーブルにはS1のサクセッサばかりでなく、S1とNとの間のノードを含めることができる。実際には、NがS1の直接のプレディセッサでない場合に、そのようになる。これらの介在するノードのいくつかは、それまではノードNには未知である場合がある。(Chord DHTアルゴリズムに従って、新規に加わったノードのサクセッサのみがノードを承知していて、他のノードは、次周期のDHT保守メッセージが予定される場合に新規のノードについて学習するであろう)。

【0043】

ノードNがノードS1のサクセッサであると、ノードS1はノードNに対する新しいプレディセッサテーブルを創出し、そしてプレディセッサテーブルをノードNに送られる離脱メッセージに含める。このテーブルには、ノードS1とノードNとの間に位置するS1のサクセッサのどれをも含む。ノードS1は、勿論、ノードNに送られるプレディセッサテーブルに自身を含めるべきではない。

【0044】

ノードS1からのサクセッサ/プレディセッサテーブルが入っている離脱メッセージを受信して、ノードNは最初にノードS1を自身のルーティングテーブルから取り除く。次に、Nは離脱メッセージで伝えられるノードのリストを細かく調べる。リストにある各ノードに対して、ノードNはノードを自身のサクセッサテーブルおよびプレディセッサテーブルにあるエントリと比較し、そしてあるノードが存在していない場合に、そのノードをテーブルの正しい位置に挿入する。

【0045】

この手順は、オーバーレイネットワークが安定状態に留まっており、そしてノードS1の撤退後でも完全な接続性を保持することを保証する。すなわち、ノードS1の撤退は決してネットワークの運用を妨害しない。

10

20

30

40

50

【 0 0 4 6 】

ノード（たとえば、図2のノードA）が近隣に知らせずにネットワークを離脱すると、異なる状況が生じる。これは、P2Pアプリケーションの突然のクラッシュまたは何らかの他の異常終了のために、または離脱ノード側での利己的な振舞いのために起こる可能性がある。ノードAの近隣のノードBが、オーバーレイネットワーク内で、ノードAが停止していることを（たとえば、ノードAに宛てた周期的なキープアライブメッセージ、保守メッセージまたは任意の他のメッセージが停止しているという事実から）最初に検出するノードであると仮定する。このとき、ノードBは、ノードAの撤退をノードAの他の近隣に知らせ、そして他の近隣の近隣テーブルの内容を更新させる責務を負わされる。しかしながら、ノードBはノードAの近隣テーブルを完全には再現できない（この例では、ノードBは、ノードAの近隣であるが、しかしノードBの近隣ではないS1に関する知識を有していない）ので、ノードB以外の近隣ノードは、ノードBが承知していないそれらの近隣ノードの近隣テーブルを更新することによる手順に関与しなければならない。

10

【 0 0 4 7 】

ノードAの撤退を最初に検出する近隣ノードが、ノードAの最も離れたサクセッサまたは最も離れたプレディセッサである例外的な場合では、かかる近隣は、他の近隣ですでに知らないものとなっている、利用できる情報を有しない。このようにして、最も離れた近隣ノードはより近くの近隣ノードの近隣テーブルの内容を更新できない。しかしながら、最も離れた近隣ノードはさらに、ノードAに代わって、より近くの近隣ノードが他の近隣ノードに近隣テーブルを伴った離脱要求を送らせるようになる空の離脱要求を送ることができる（より近くの近隣ノードが他のノードに恩恵をもたらす情報を有する）。

20

【 0 0 4 8 】

例として、図2のノードXが自身の近隣、すなわち、ノードS1、S2、S3、P1、P2およびP3に知らせずに、ネットワークを離脱することを仮定する。さらに、たとえばノードXがノードS1からのキープアライブメッセージに回答し損なうために、ノードS1が、ノードXが離脱していることを検出する最初のノードであると仮定する。ノードXの停止を検出した直後に、ノードS1は、ノードXの近隣テーブルの内容のできる限り正確な表現を創出することになる。この場合、S1はノードXの6つの近隣のうちの5つに関する情報を再現できる。これは、S1がノードXのすべてのサクセッサおよびノードXのプレディセッサのうちの2つを知っているからである。次に、ノードS1は、S1が承知している、ノードXの各近隣Nに対して以下の手順を行う。

30

【 0 0 4 9 】

ノードS1はノードNの近隣テーブルの内容を構築する。

【 0 0 5 0 】

ノードS1は、S1がノードNに対して構築した近隣テーブルからノードXを取り除き、ノードXの回復された近隣テーブルから適切なノードを選び、そしてこれらをノードNに対する新しい近隣テーブルに挿入する。より具体的には、ノードNがノードXのプレディセッサであると、ノードS1はノードNに対する新しいサクセッサテーブルを創出する。その代わりに、ノードNがノードXのサクセッサであると、ノードS1はノードNに対する新しいプレディセッサテーブルを創出する。

40

【 0 0 5 1 】

ノードS1は、ノードXに代わって、ノードNに離脱要求を送る。離脱要求には、S1が創出している新しいサクセッサテーブルまたはプレディセッサテーブルを含める。

【 0 0 5 2 】

ノードNは、ノードS1から受信した離脱要求に基づいて自身の近隣テーブルを更新する。

【 0 0 5 3 】

ノードXのルーティングテーブルの内容を決定する場合、ノードS1はノードXの最も離れているプレディセッサ、すなわち、P3のID（identity）を知らない。しかしながら、P1とP2の双方はP3のIDを知っている。したがって、ノードXの近隣Nが、ノ

50

ードXに代わって別のノードから送られた離脱要求を受信する（離脱要求の発信元のソースアドレスがノードXのアドレスと一致しないので、近隣Nがこれを検出できる）と、近隣ノードNは以下の処理動作を実行する。

【0054】

ノードNは、ノードXの近隣テーブルの内容の表現を再現する。より具体的には、ノードNがノードXのプレディセッサであると、ノードNはXのプレディセッサテーブルを再現する。他方、ノードNがノードXのサクセッサであると、ノードNはノードXのサクセッサテーブルの表現を再現する。

【0055】

ノードNが構築した近隣テーブルに基づいて、ノードNは、ノードXの任意の近隣のうち、離脱要求の送信元（すなわち、ノードS1）が承知していないノードを承知しているかどうかを調べる。ノードNが図2のノードP2であると仮定すると、P2は離脱要求の送信元、ノードS1がノードP3（ノードXの第3のプレディセッサ）を承知していないことを検出する。

10

【0056】

ノードNは、ノードXに代わって、離脱要求の送信元が承知していない、ノードXの近隣の各々に離脱要求を送る。再び、ノードNが図2のノードP2であり、そして送信元がS1であると仮定すると、P2はノードP3に離脱要求を送ることになるであろう。

【0057】

最後に、ノードNは、離脱要求の送信元（すなわち、ノードS1）が承知していない、ノードXの近隣のリストを（離脱要求に対してノードNが生成する応答で）返送する。離脱要求の送信元（すなわち、ノードS1）はそれから、これらのノードを自身のプレディセッサテーブルに挿入できる。

20

【0058】

上記例では、ノードXがオーバーレイネットワークを離脱していることを検出する最初のノードが、サクセッサ（S1）であると仮定した。ノードXがオーバーレイネットワークを離脱していることを検出する最初のノードがプレディセッサであると、プレディセッサ（たとえば、P1）は、ノードXの最も離れたサクセッサ、すなわちこの例ではS3に関する知識を有しない。この場合、ノードXのサクセッサのうちの1つ、たとえばノードS1が上述の工程を実行する。

30

【0059】

保守運用は、DHTネットワークで実行され、そして新しいノードについて学習し、近隣ノードの状態を調べるための、参加ノード間における周期的なメッセージ交換を含む。しかしながら、DHTに基づくオーバーレイネットワークにおけるノードは慣習的に、ノードが自身のルーティングテーブルを更新するために近隣から受信する情報の一部のみを用い、そして情報の残りを廃棄する。たとえば、図2のネットワーク例では、ノードXは直通的近隣P1からノードEの存在について学習するけれども、ノードXがプレディセッサテーブルに3つのプレディセッサポイント、プレディセッサP1、P2およびP3のみを維持しているので、通常、ノードXはノードEの情報を廃棄する。しかしながら、ノードXがノードEの連絡情報をキャッシュすると、ノードXはノードP1（およびまたP2およびP3）の手順を踏まない、考えうる撤退からより迅速に回復できる。言い換えると、ノードXはノードEを自身の近隣テーブルに挿入しないけれども、ノードXは近隣キャッシュを実装する別のデータ構造にノードEの連絡情報を記憶する。ノードP1が突然ネットワークを離脱すると、ノードXは、ノードEが活性しているか、そしてルーティングテーブルに追加できるかどうかを確認するために、キャッシュを調べ、そしてたとえばノードEに連絡する。

40

【0060】

以下、上述したメカニズムを実装するのに相応しいDHTに基づくオーバーレイネットワークのノード1が概略的に説明されている図3を参照する。ノード1は、ノードのためのルーティングテーブルを記憶するように構成された記憶装置2を備える。ノードには、オ

50

ーバレイネットワークにおける他のノードとのインタフェース3（典型的には、IPネットワーク、たとえばインターネットとのインタフェース）が備わっている。第1の処理部4は、インタフェース3を通して近隣ノードから離脱要求を受信するように構成されている。離脱要求が受信された場合、第1の処理部はルーティングテーブルを上述のように更新するように構成されている。離脱要求はまた、撤退ノードが要求を送ったノードであるかどうかを決定する第2の処理部5により受信される。撤退ノードが要求を送ったノードでない場合には、第2の処理部5が、要求を送ったノードが承知していない、撤退ノードの近隣ノードを識別するために、ルーティングテーブルを調べる。第2の処理部はそれから、1つ以上のさらなる離脱要求を構築し、そしてこれらを識別されたノードに送る。

【0061】

第3の処理部6は、ノードのネットワークからの手順を踏んだ撤退を取り扱うように構成されている。第3の処理部6は離脱要求をノードの近隣に送り、各近隣ノードに対する1つ以上の代替マッピングを識別することによりこれを行う。この第3の処理部はまた、近隣が手順を踏まずにオーバーレイネットワークを離脱する場合、近隣ノードにキープアライブメッセージを周期的に送ることに関与する、接続性検出器7により通知されるように、そして撤退ノードの近隣に適切な離脱要求を送るよう構成されている。

【0062】

第4の処理部8は、メモリ9内に、その時点のルーティングテーブルに入っていないノードに対するマッピングが入っているキャッシュを維持する。第4の処理部が近隣から離脱要求を受信する場合、第4の処理部はキャッシュメモリからマッピングを抽出し、そしてこれをルーティングテーブルに追加する場合がある。

【0063】

次に、手順を踏んだ撤退シナリオで適用されるメカニズムを説明している図4を参照する。工程S1からS3までで、撤退ノードは離脱要求を近隣ノード（ここでは、3つのみが示されている）に送る。工程S4からS6までで、近隣ノードの各々はルーティングテーブルを更新する。

【0064】

図5は、手順を踏まない撤退のシナリオの場合に適用されるメカニズムを説明している。工程T1で、撤退ノードの第1の近隣ノードが周期的なキープアライブメッセージを撤退ノードに送る。第1の近隣で、キープアライブへの応答がその間に受信されていないタイムアウトT2の後、第1の近隣は、撤退ノードに対するエントリを削除するために工程T3でルーティングテーブルを更新する。工程T4およびT5で、第1の近隣は撤退ノードに代わって、離脱要求を撤退ノードの近隣に送る。離脱要求には、前述のように新しいマッピングが入っている。離脱要求の受信時点で、工程T6およびT7で近隣は自身のルーティングテーブルを更新する。近隣のうちの1つ、この場合では第3の近隣は、第1の近隣が承知していない、撤退ノードのさらなる近隣を承知している。工程T8で、近隣は、かかるさらなる近隣にさらなる離脱要求を送る。工程T9で、さらなる近隣が自身のルーティングテーブルを更新する。

【0065】

本明細書で説明したメカニズムは、DHTに基づくオーバーレイネットワークのロバスト性を向上するので、本メカニズムは、インターネット技術特別委員会（IETF）のP2PSIP作業部会により標準化途上にある、P2PSIP電話のような重大なDHTに基づくシステムに特に有用である。

【0066】

当業者には当然のことながら、本発明の範囲を逸脱することなく、様々な変更が上述した実施形態になされる可能性がある。

10

20

30

40

【図1】

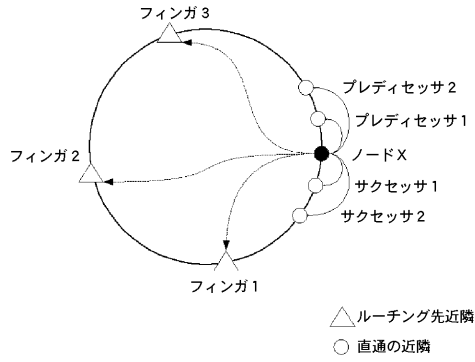


Figure 1

【図3】

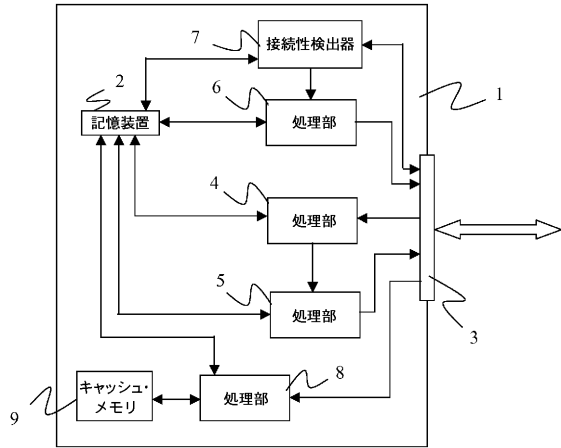


Figure 3

【図2】

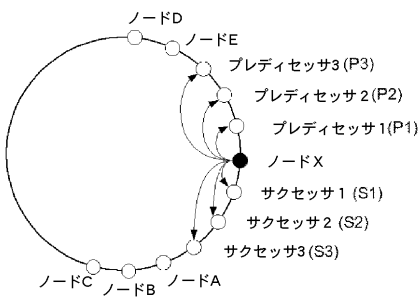


Figure 2

【図4】

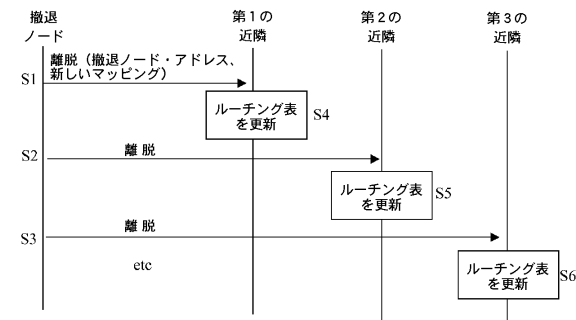


Figure 4

【図5】

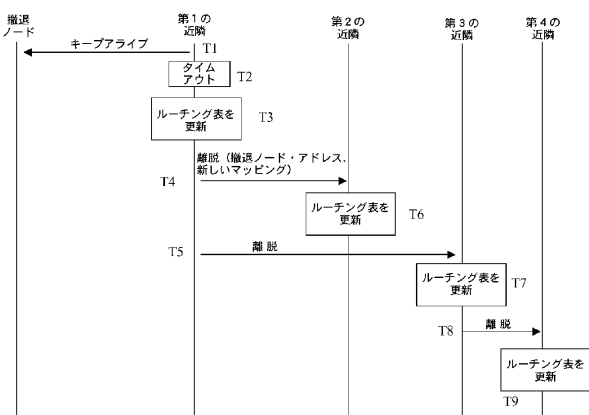


Figure 5

フロントページの続き

(72)発明者 メエンペー, ヨウニ

フィンランド国 ヌメラ エファイ - 03100, ナーランパコンティエ 32 エー5

審査官 永井 啓司

(56)参考文献 Ion Stoica, Chord: A Scalable Peer-to-Peer Lookup Protocol for Internet Applications, Networking, IEEE/ACM Transactions on (Volume:11, Issue: 1), IEEE, 2003年 2月, 17-32頁, URL, http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=1180543&tag=1
中村 元紀, ネットワーク分割に対するDHTの可用性向上, 電子情報通信学会2005年総合大会講演論文集 通信2, 日本, 社団法人電子情報通信学会, 2005年 3月 7日, S-26, S-27

白石 俊之, Chordネットワークにおけるルーティング情報の効率的な維持管理方法, 電子情報通信学会技術研究報告 Vol. 106 No. 577, 日本, 社団法人電子情報通信学会, 2007年 3月 1日, 第106巻, 137-142頁

西谷 智広, 分散ハッシュテーブル入門, UNIX magazine Vol. 21 No. 6, 日本, 株式会社アスキー, 2006年 9月21日, 第21巻, 26-33頁

Sean Rhea, Handling Churn in a DHT, Appears in Proceedings of the USENIX Annual Technical Conference, 2004年 6月

(58)調査した分野(Int.Cl., DB名)

H04L 12/00 - 12/26, 12/50 - 12/955