物理ネットワークの状況を考慮した階層型分散ハッシュ法の提案

羽場 裕介 † 松尾 啓志 †

† 名古屋工業大学大学院情報工学専攻

分散ハッシュ法 (Distributed Hashing) は, ピアツーピア環境において高速にファイルやノードを検索 するための手法として注目されている.しかし,分散ハッシュ上のオーバレイネットワークの構築には, 物理ネットワークを一切考慮しない.そのため,検索ホップ中に無駄が含まれる可能性がある.そこで 本稿では,物理ネットワークにおける遅延時間を物理的な近さであると位置づけ,それを分散ハッシュ 上の ID に反映させる.また,それに加え階層化を行うことでより効率的に検索を行える新しい階層型 の分散ハッシュ法を提案する.

Hierarchical Distributed Hashing in consideration of physical network

Yuusuke HABA † Hiroshi MATSUO †

[†] Department of Computer Science and Engineering

Graduate School of Engineering Nagoya Institute of Technology

Distributed Hashing attracts attention of method which fast lookup a file or a node in peer-to-peer environment. However, to construction of overlay-network based on DHT is not consider of physical network, so waste may be had within a lookup hop. Therefore, in this paper, physical closeness cast delay time in a physical network. And we let ID on DHT reflect physical closeness. Additionally, to hierarchized by using the ID. We suggest a new hierarchy Distributed Hashing to lokkup more effectively.

1 はじめに

近年,インスタントメッセンジャーなと多くのピア ツーピア(P2P)アプリケーションが使用されており, インターネット上でオーバレイネットワークが構築 されている.しかし,ピュアピアツーピアネットワー クでは,中央サーバが存在しないという利点は存在 するものの,目的のファイルやノードを高速に検索 することが困難である.そのため,ピュアピアツーピ アネットワーク上で高速に検索することを可能にす る分散ハッシュ法(Distributed Hashing)や分散ハッ シュテーブル(Distributed Hash Table : DHT)と 呼ばれる手法が注目されている.分散ハッシュテー ブルを使用し,その上でオーバレイネットワークを 構築することで高速な検索が可能となる.

しかし,分散ハッシュテーブルを使用したオーバ レイネットワークは,実際の物理ネットワークの近 さが分散ハッシュ上での近さとは無関係である.その ため,オーバレイネットワーク上では無駄の無い検 索ホップであっても,実際の物理ネットワークでは, 一度遠くのネットワークに存在するノードまで検索 が進んだ後に,検索を開始したノードの近くのノー ドに検索が戻って来る無駄が含まれる可能性がある. そこで本研究では,物理ネットワークでのノード間 の遅延時間を物理ネットワークでの近さであると定 義し,物理ネットワークでの近さを分散ハッシュ上 での ID に反映させ,その ID を使用して階層化を 行うことで,先に挙げたような問題点を改良する手 法を提案する.

本稿は,2節で分散ハッシュテーブルの関連研究, 本研究の提案手法のもととなる Chord,提案手法の 比較となる HIERAS の説明を行う.3節で提案手法 の説明をし,4節で提案手法の効果を確認するため 評価実験を行い,最後に5節でまとめる.

2 関連研究

分散ハッシュテーブルを用いる手法の代表的なものとして Pastry[2], Tapestry[3], CAN[4] などがあり, 少数のノード情報をテーブルに保存することで, 目的のノードを $O(\log N)(N: J - Fo数)$ で検索を行うことが可能である. Pastry は N 分木を使用することで,検索にかかるホップ数は $O(\log_b N)$ (b は ID の基数) である.

本研究で提案手法のもととした Chord[1] では,各 ノードは SHA-1 などのハッシュ関数により自分の ID を生成する.自分の ID と他のノードの ID を もとに管理すべき領域や保存すべきノードの情報を 決定することができるため,完全に分散して動作す ることが可能である. Chord では, ID 空間上で自 分の ID から見て時計周りの方向に最初に存在する ノードを successor として記憶する.また,同様に反時計周りに最初に存在するノードを predecessor として記憶し,その predecessor と自分との間の ID 空間を責任領域として管理を行う.また,各ノードは,ルーティングに用いる m 行 (m:IDのビット数)のテーブル(フィンガーテーブル)を保持し,その i 行目には $ID + 2^{i-1}$ の successor の情報が入る.つまり,各ノードは自分の 1,2,4,..., 2^{m-1} 先のノード情報を保持していることになる. Chord はこれらの情報のみを用い,検索にかかるホップ数を $O(\log N)$ に抑えている.

Chord を使用し,先に挙げたような分散ハッシュを 使用したオーバレイネットワークの問題点を改良した 手法として HIERAS[5] が提案されている . HIERAS 上の各ノードは, ランドマークノード(ネットワーク 中のどのノードからも存在と IP アドレスなどが知 られているノード)への遅延時間をもとにした,物 理ネットワークでの位置を表すリング ID と呼ばれ る特別な識別子を持つ.同じリング ID を持つノー ド同士で生成した Chord のネットワークを通常の Chord でのネットワークとは別に保持することで階 層化を行っている. HIERAS ではこのように複数 の Chord ネットワークを使用することで先に挙げ た問題を改良している.しかし, HIERAS では各 階層で Chord と同じ情報を保持しなければならな いため,維持コストが増加する.また,階層化を行 うことにより,物理層での遅延時間は減少する反面, アプリケーション層でのホップ数が増加することが 報告されている[5].

3 提案手法

本研究では, Chord で使用される ID に物理ネット ワークの近さを反映させ, その ID を利用し階層化を 行う LCLV (Layerd Chord with Landmark Vector) を提案する.HIERAS では,リング ID と呼ばれる分 散ハッシュ上の ID とは別の識別子を新たに生成し, それをもとに階層化を行っていた.しかし,LCLV では分散ハッシュ上の ID そのものに物理ネットワー クでの近さを反映させる.この点が HIERAS との 違いである.

3.1 ランドマークベクトルによる ID の生成法

LCLV は物理ネットワーク上での距離を通信に必要な遅延時間により近似可能であるという前提のもとに,遅延時間の短いノード同士は,分散ハッシュ上のネットワークにおいても近くなるような ID を 生成する.そのような手法として,近傍の情報を使用して生成されるランドマーククラスタリングが広



く使用されている [7][8][9].また,近年,近隣情報と ランドマークの情報のみを利用して負荷分散を行う DHT ベースの P2P システムが提案されている [6].

LCLV における新しいノード N の ID の生成の流 れを以下に示す.

- 1. 各ランドマークの2次元空間での位置 (座標) を 求める.
- 各ランドマークの座標と、Nから各ランドマー クへの遅延時間から、Nの2次元空間での座標 を求める。
- 3. 決定された座標から ID を生成.

まず,ランドマークノードを図1のように2次元 空間にマッピングする(1.).その方法として,各ラ ンドマーク間の物理ネットワーク上での遅延時間を 測定し,測定された遅延と2次元空間での距離との 誤差の和が最小となる点にランドマークを配置する. 次に同様の方法を用いて,今IDを決定したいノー ド N の2次元空間での座標を求める(2.).以上の方 法で,物理ネットワークで近いノード同士は,2次 元空間においても近似的に近い座標を得ることが可 能となる.

次に,2次元空間をいくつかの空間に分け,そこ に図2のように空間充填曲線 (Space Filling Curve: SFC)を引く.ここでは,空間充填曲線として Hilbert 曲線を用いた.そして,空間充填曲線が通った順に 数字を割り当てる.このとき,ノード N が存在する 座標上の SFC によって割り当てられた数字をノー ドの ID として新たに割り当てる (3.).



この方法を用いることにより,物理ネットワーク で近いノード同士を,分散ハッシュ上での ID にお いてもおおよそ近い ID に割り当てることが可能と なる.

3.2 階層化

LCLV における階層化は, Chord で用いる ID を 数ビット毎に区切ることで行う.例として, ID の 大きさが12ビット,階層数が2,第一階層の ID が3 ビット,第二階層の ID が9ビットとし,あるノード N の ID が8進数表示で ID = 7049 と与えられた 場合を考える.この場合,第一階層の ID は7とな り,第二階層の ID は049 となる.つまり,図3の ように階層を分けることになり,ノードN は図中の 位置に配置される.このような階層化において,そ れぞれの階層は ID が小さいため,もとの Chord よ りも小さい規模のネットワークとなる.このように LCLV は,小さな規模のネットワークに分けること での階層化である.HIERAS は各階層で Chord と 同じ大きさのネットワークを生成するため, LCLV の方が維持コスト,テーブル量において優位である.

LCLV では上位の階層の ID のみを 3.1 節で説明 したランドマークベクトルによる ID を使用し,最 下層の ID は通常の Chord と同様にハッシュ関数を 使用して割り当てる.

このように ID を用いることで,図3においては, 第一階層でのネットワークは物理的に遠いノード同 士のネットワークとなるので,この階層における通 信遅延は大きくなる.しかし,第二階層でのネット ワークは物理的に近いノード同士で作られているの で,通信遅延は小さくなる.



LCLV では、次の理由によりホップ数減少の効果 も期待できる.例として,図4にIDの大きさを8 ビットとしたときの Chord と LCLV を比較した図 を示す.図4の(a)にChordの全体を,(b)に提案 手法の第一階層のみを示す.なお,ノードの ID は 16 進数で表す. Chord では,最も遠方のリンクと して $80(=128_{(10)})$ 先のノードの successor の情報 を有する. つまり図 4 で, ID = 00 であるノード は $ID = 82(= 130_{(10)})$ であるノードの情報を持つ. LCLV においても,同様に最も遠方のリンクとして 80 先のノードの successor を知っていることになる. しかし, LCLV における第一階層に注目すると, ID が 8 から始まるノードの successor は ID が 9 か ら始まる (もしくは, それ以降) ノードになる. つま り,図4中でID=00であるノードはID=8*の successor である ID = 9* のノード情報を持つ.こ れは, Chord では $ID = 96(=150_{(10)})$ のノード 情報を有することと同じ効果がある.このように, LCLV はより遠方のノード情報を持つ可能性が高い ため検索にかかるホップ数が減少する効果が期待で きる.

3.3 LCLV における検索

LCLV において key を検索する手順を以下に示す.

- 第一階層において key に最も近い predecessor まで, Chord と同様の手法で検索を進め,次 の階層に進む.
- 第二階層においても同様に, key に最も近い predecessor まで, Chord と同様の手法で検索 を進め, 次の階層に進む.
- 3. 以下同様に,最下層まで検索を進める.
- 4. 最下層において, *key* を管理するノードを Chord と同様の手法で検索する.

以上のような検索の手順により,検索初期のホッ プは,上位の階層で検索が進む.3.1,3.2節により, 上位の階層は物理的に遠いネットワークとなるので, 検索初期のホップでは1ホップあたりの通信にかか

	テーブル量	維持コスト
Chord	m	2+m
HIERAS	$l \times m$	$l \times (2+m)$
proposal	m	3n(l-1) + (2+m)

Table 1 コスト比較

る遅延時間が大きくなる.また,検索中期以降は下 位の階層で検索が行われる.同様に3.1,3.2節より, 下位の階層は物理的に近いネットワークとなるので, 検索の中期以降のホップでは1ホップあたりの通信 にかかる遅延時間が小さくなる.

また, Chord の検索は,検索の初期や中期以降に 関わらず遅延時間の偏りは無いため,検索の全過程 において1ホップあたりは同じ程度の遅延時間にな る.しかし,LCLVでは,検索の初期数ホップのみ が大きい遅延時間になり,それ以降は小さい遅延時 間になるためトータルでの遅延時間は小さくなる.

3.4 LCLV における保持すべき情報

LCLV において各ノードが保持する必要のある情報は、Chord と同様にフィンガーテーブル、successor, predecessor である.しかし、LCLV においては、上位の層でsuccessor, predecessor を最大でn 個保持する必要がある.これは、階層化による経路の集中を防ぐことと、ノードが故障したときの冗長化が目的となる.また、同様の理由により各階層で代表となるノードをn 個覚える必要がある.また、一番下位の階層での情報は Chord と同じ情報を保持する.LCLV、Chord、HIERASのコストの比較を表1に示す.ここで、mは ID のビット数、lは階層数,nは LCLV における保存するノードの数である.

LCLV におけるテーブル量は,階層数に依存しない ため Chord と同じ大きさである.しかし,HIERAS では各階層おいて Chord と同じ大きさのテーブルが 必要であるため,階層数倍のテーブルが必要となる.

維持コストについて比較すると,LCLV はm > 3nならば HIERAS よりも維持コストが小さい.現実的なネットワークの規模を考慮した場合,m = 160程度であり,仮にnを比較的大きな50としても維持コストは HIERAS よりも小さくなる.また,実際のピアツーピアネットワークではn = 5程度で充分であると考えられる.

4 評価実験

LCLV の効果を評価するために, Chord, HI-ERAS, LCLV のそれぞれをネットワークシミュレー タである NS-2[12] 上に実装し評価実験を行った.実



Fig. 5 Average Application Hop(TS-model 1)

験では, 各ノードがランダムな key を検索し, 検索 にかかった遅延時間の平均, アプリケーション層で のホップ数の平均, 物理層でのホップ数の平均を測 定する.

また, LCLV においては, ランドマークベクトル による ID の割り当ての効果を確認するため, 階層 化のみを行った場合, 階層化とランドマークベクト ルによる ID の両方を行った場合の二種類の実験を 行った.

ID の大きさを 24 ビットとし, ノード数を 512 台 から 2048 台まで変化させて実験を行った. LCLV における種々のパラメータは, ランドマークの数を 15 台, 階層数を 3 とし, 第一階層の ID は 4 ビッ ト, 第二階層の ID を 4 ビット, 第三階層の ID を 16 ビットに設定した.また,実験に使用したネッ トワークは GT-ITM(TS モデル)[11] を使用した. TS モデルの設定において, HIERAS の論文 [5] では Transit-Transit(TT) 間の遅延時間を 100ms, Transit-Stub(TS) 間を 20ms, Stub-Stub(SS) 間を 5ms としていることを考慮し,本研究では,以下の ような 3 つの設定で実験を行った.

- 1. TS モデル1 (設定の変更無し)
- 2. TS モデル 2 (HIERAS と同じ設定 TT 間 100ms, TS 間 20ms, SS 間 5ms)
- 3. TS モデル 3 (遅延時間の比率を変更 TT 間 100ms, TS 間 80ms, SS 間 60ms)

3. については, HIERAS での設定の比率を変更し, より均質なネットワークに近づけた場合の HIERAS と LCLV の性能を比較する目的で用いた.

4.1 TS モデル1 での結果

図 5 にアプリケーション層でのホップ数の比較の グラフ,図 6 に物理層でのホップ数の比較のグラフ,



Fig. 6 Average Physical Hop(TS-model 1)



図 7 に遅延時間を比較したときのグラフをそれぞれ 示す.図5と図6より, HIERAS は Chord よりも アプリケーション層でのホップ,物理層でのホップ 共にホップ数が増加していることが確認できる.ま た,3.2節で説明したように, LCLV にはホップ数 減少の効果があるため, Chord, HIERAS よりも検 索にかかるホップ数が減少していることも同様に確 認できる.

図7より, HIERAS は Chord よりも検索にか かる遅延時間がわずかながら減少している.これは, HIERAS の階層化の方法では, TS モデル1のよう な遅延時間に規則性の無いネットワークを適切に階 層化出来なかったからと考えられる.

また,LCLV とその他の手法を比較すると,LCLV の方が他の手法よりも遅延時間が大きく減少してい る.これは,階層化のみの場合においても遅延時間 が減少していることから,図5や図6においてホップ 数が減少した分だけ遅延時間が減少したからである.



LCLV 同士を比較してみた場合,階層化のみの場合,階層化とランドマークベクトルの両方を用いた 場合の両方において同程度の遅延時間になっている. このときのランドマークベクトルを用いた場合の検 索の過程を追跡すると,検索の初期,中期以降のい ずれにも関わらず同程度の遅延時間になっていた. つまり,物理的に近いノード同士であるにも関わら ず,分散ハッシュ上の ID では遠くなってしまう場 合が多数存在していた.これは,図2において数字 が2と7の位置は座標的には近いが分散ハッシュ上 の ID 的には遠くなってしまっていることが原因で あると考える.

4.2 TS モデル2 での結果

図 8 に, HIERAS の論文と同じ設定である TS モデル 2 で実験を行ったときの遅延時間のグラフを 示す.また,ホップ数に関してのグラフは 4.1 節で のグラフと同様の傾向であった.

図8に示すとおり, HIERAS は Chord よりも大 きく遅延時間が減少している.このことは,先の4.1 節の結果と違い,TS モデル2のネットワークは遅 延時間に規則性があるため,ネットワークの階層化 が適切に出来たためである.

また, LCLV と他の手法を比較すると, LCLV の方が他の手法よりも遅延時間が大きく減少している.これは,4.1 節と同様の理由である.

LCLV 同士を比較してみると, 階層化のみに比べ, ランドマークベクトルも使用している方は, 遅延時 間が大きく減少している.これは,3.3 節で説明し たように,ランドマークベクトルによる ID の生成 により,物理的に近いノード同士が分散ハッシュ上 でも近い ID を得ることが出来たことによる,遅延 時間減少の効果である.



4.3 TS モデル3 での結果

図 9 に,遅延時間の比率を HIERAS の設定から 変更した TS モデル 3 で実験を行ったときの遅延時 間のグラフを示す.

図9より, HIERAS は Chord よりも遅延時間 が大きくなった.このことは,TS モデル3のよう に,遅延時間に規則性はあるものの差が少ないネッ トワークでは, HIERAS によるネットワークの階 層化を適切に行うことができず,逆に遅延時間を増 加させる結果となったためである.

LCLV と他の手法を比較すると,図7,図8と同 様に他の手法より遅延時間が減少している.また, LCLV 同士を比較すると,HIERAS のように遅延時 間が増加してしまうことは無く,ランドマークベク トルを使用した場合の方が遅延時間が減少している. これは,ランドマークベクトルによる ID の生成が, TS モデル3のように,遅延時間の差が少ないネット ワークにおいても適切に働いたためである.

5 まとめ

分散ハッシュを使用したオーバレイネットワーク は、物理ネットワークでの遅延時間などを考慮した ものになっていないため、検索ホップ中に無駄が含 まれてしまう可能性があった.本研究では、ランド マークベクトルを使用し物理ネットワークにおける 近さを分散ハッシュでの ID に反映させた.また、そ の ID を使用して階層化を行うことで、検索にかか るホップ数、遅延時間共に減少させることができた. また別の手法による階層化よりも、より効率良く検 索が可能であることを実験により確認した.

参考文献

 I. Stoica, R. Morris, D. Karger, M.Kaashoek, and H. Balakrishnan, "Chord: A Scalable Peer-to-Peer Lookup Service for Internet Applications". Proc. ACM SIGCOMM, pp.149-160, Aug. 2001.

- [2] A. Rowstron and P. Drushel, "Pastry: Scalable, Decentralized Object Location, and Routing for Large-Scale Peer-to-Peer Systems". Proc. 18th IFIP/ACM Int'l Conf. Distributed System Platforms (Middleware), pp.329-350, Nov. 2001.
- [3] B.Y. Zhao, J.D. Kubiatowicz, and A.D. Joseph, "Tapestry: An Infrastructure for Fault-Tolerance Wide-Area Location and Routing". Technical Report UCB/CSD-01-1141, Computer Science Division, Univ. of California, Berkeley, Apr. 2001.
- [4] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker, "A scalable content addressable network". Technical Report, TR-00-010, U.C.Berkeley, CA, 2000.
- [5] Z. Xu, R. Min, and Y. Hu,1 "HIERAS: A DHT Based Hierarchical P2P Routing Algorithm". Proc. of hte 2003 International Conference on Parallel Processing, pp.187-194, Aug. 2003.
- [6] Y. Zhu, and Y. Hu, "Efficient, Proximity-Aware Load Balancing for DHT-Based P2P Systems". IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, VOL 16, No. 4, APRIL 2005.
- [7] S. Ratnasamy, M. Handley, R.M. Karp, and S. Shenker, "Topologically-Aware Overlay Construction and Server Selection". Proc. IEEE INFOCOM, vol. 3, pp. 1190-1199, June 2002.
- [8] Z. Xu, C. Tang, and Z. Zhang, "Building Topology-Aware Overlays Using Global Soft-State". Proc. 23rd Int 'l Conf. Distributed Computing Systems (ICDCS), pp. 500-508, May 2003.
- [9] Z. Xu, M. Mahalingam, and M. Karlsson, "Turning Heterogeneity into an Advantage in Overlay Routing". Proc. IEEE INFOCOM, vol. 2, pp. 1499-1509, Apr. 2003.
- [10] T. S. Eugene, Ng, H. Zhang, "Towards global network positioning". ACM SIGCOMM Internet Measurement Workshop 2001.
- [11] E.W. Zegura, K. L. Calvert, and S. Bhattacharjee, "How to model an internetwork". Proceedings of the IEEE Conference on Computer Communication, San Francisco, CA, pp. 594.602, Mar. 1996.
- [12] http://www.isi.edu/nsnam/ns/