

## IDEON-retreat レジューム: 担当 doi 2006-04-03

Peter Pietzuch, Jeffrey Shneidman, Jonathan Ledlie, Matt Welsh, Margo Seltzer, and Mema Roussopoulos (Harvard Univ.), *Evaluating DHT-Based Service Placement for Stream-Based Overlays*, Proceedings of the 4th International Workshop on Peer-to-Peer Systems (IPTPS'05), Ithaca, New York, February 2005.

## 1 背景

この論文の背景には、著者らが SBON(Stream-based overlay networks) と呼ぶアプリケーションのクラスがある。SBON の例として挙げられているのは、Aurora[文献 2], Borealis[文献 1], PIER[文献 9], IrisNet[原論文文献 8] などである。著者らは HOURGLASS<sup>1</sup> と呼ぶセンサーネットワーク処理オーバーレイグループの一員であり、その関連研究である。なお、今回は時間の関係で HOURGLASS 自体を記述した論文までサーベイが及ばなかった。

SBON では、処理を Circuit という単位で記述する。Circuit は、継続的にデータを生成する Producer(s)、Producer(s) から生成されたデータを処理する Operator(s) そして一つの Consumer から成る。それぞれの要素間をつなぐリンクを Circuit Link と呼ぶ。Operator(s) は複数の Producer から生成されたデータに、平均やフィルタリングなどの前処理を加えてデータ量を絞り込んだ上で Consumer(おそらく実際にユーザが利用するアプリケーション/クライアント) に渡す構造になっている(図 1)。

Consumer/Producer はノード位置の制約がある。前者はユーザが選択したノードであり、後者はデータが生成されるノードである。一方、Operator は入出力を持つ一つのプロセスであり、基本的にどのノードに設置しても問題はない。

本論文では、以上の定義において、Consumer/Producer(s) が与えられたときに、Circuit により消費するネットワーク資源に対して、Operator(s) の配置を最適化する、という問題に対して、DHT<sup>2</sup>の経路情報のみを利用する方式の有用性を検

<sup>1</sup><http://www.eecs.harvard.edu/>

[%7esyrah/hourglass/index.shtml](http://www.eecs.harvard.edu/~7esyrah/hourglass/index.shtml)

<sup>2</sup>この論文では DHT で KBR を指している?

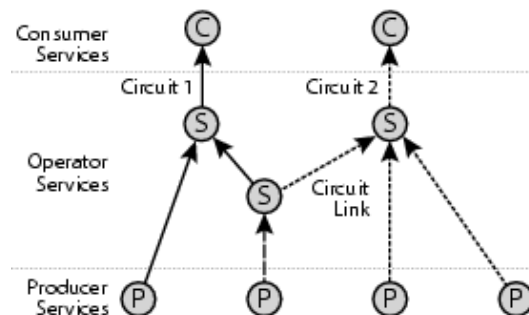


図 1: SBON の概略 (Harvard のページより引用)

証している。検証の結果、Proximity Optimization を利用した DHT アルゴリズム (Bamboo/Pastry) であっても、DHT 経路情報の利用による Operator 配置問題への寄与はわずかであり、他の方法を選択したほうが良い、と結論づけている。<sup>3</sup>著者らは spring relaxation に基づく方式 [文献 11] を提案している。

## 2 問題定義の詳細化と評価方法

著者らは、最適化の尺度として、BW-Lat product、すなわち消費する帯域と遅延を乗じたものを選択している。最適な配置とは、Producer(s)-Operator(s)-Consumer の全ての経路 (circuit link) 上の BW-Lat の和を最小化する配置としている (section 2.1)。<sup>4</sup>

また、Operator 配置問題を二つに分割し、そのうち discovery の部分の評価を中心に議論している。

分割された問題の一つは node discovery であり、candidate set を選択するメカニズムである。もう一つは node selection であり、candidate setの中から一つのノードを選択するメカニズムである。この二つの定義上の明確な区別は論文上では示されていないが、selection ではノード間のネットワークの品質測定を行う場合があるのに対して、discovery では事前知識などから副次的に得られるノードの組を選出している。

discovery のメカニズムとして評価されているのは次の 6 つ (variant 除く) である。

<sup>3</sup>が、そうと言いきれるほど強いデータは示されていないと思う。

<sup>4</sup>doi: なお、ここでの著者らの BW-Lat product が最適化の尺度として適切であるかという議論は弱いように思える。

All 全てのノードを評価の対象とする

Consumer Consumer ノード (1つ) を選択する

Producer Producer ノード (複数) を選択する

RandomSet(k)  $k$  個のランダムなノードを選択する

DHT Union 全てのProducer から Consumer までの Overlay 経路を構成するノード集合の和

DHT Intersection " の積

また、node selection のメカニズムとして評価されているのは次の3つである。

Random ランダムに一つ選択。

Optimal candidate 全てから最適なもの (組み合わせ) を選択

Relaxation [文献 11,5] により提案される方式

最後に、問題の評価の条件を次のように設定している。評価のフィールドとなるのは、PlanetLab を構成する 186 ノード、あるいは GATech topology generator から生成された 600 ノード。評価には、4 つの Producer、1 つの Operator (unpinned service)、1 つの Consumer からなる 1000 個の Circuit を生成し、BW-Lat product を集計する。なお、Producer はそれぞれ 2kb/s のデータストリームを生成し、Operator により 1kb/s のデータストリームに統合されると仮定している。

実装およびアルゴリズムは、Bamboo を利用した OpenHash および Pastry をベースにした独自開発の Pan を用いている。いわゆる  $b$  (DHT key base) の影響の評価を  $\text{Pan}(b = 64)$  を用いて行っている。

### 3 評価結果

原論文 Table 1 にそれぞれの方式における 80th percentile の、最適配置時における結果に対する比が示されている。CDF を見ると分布の傾向はそれぞれの方式に差は少ないため、Table 1 のみでも評価結果を読み取ることができる。<sup>5</sup>

<sup>5</sup> 著者による議論は別論文における提案方式 [文献 11] を推すバイアスがかかっているように思えるので、ここから下の議論は無視して Table 1 のみを見る。

ここで、上の 6 段に示されている非 DHT 方式はそれぞれ有効な順に並んでおり、まったくランダムに選択しても最適な利用の 2 倍未満である。ここで、DHTUnion 方式による寄与を、DHTUnion discovery の結果の candidate 数の平均に近いランダムなノード集合である RandomSet(6) と比較する。Selection 方式に Optimal を利用したとすると、DHTUnion:1.13 ~ 1.18 に対して RandomSet(6):1.26 ~ 1.36 であり、DHTUnion はオーバーヘッドを 1 割以上削減するのに貢献している。ただし、DHTUnion には Producer Set が含まれており、これを削除する (DHTUnion-NoProd) と、1.17 ~ 1.31 となり、DHT そのものによる貢献は相対的には小さくなる。また、DHTUnion の結果を Random に選択すると、1.52 ~ 1.63 となり、GATech トポロジでは単純に Producer からランダムに選択したほうが有利となる (1.43) ケースすらあった。

一方、DHTIntersection は、DHT 上の経路を Scribe 的な ALM tree とみなした場合、ALM tree の結節点に Operator を置く、という方式である。しかし、Producer が 4 つ存在する今回のケースでは、candidate node がほとんど Consumer になってしまい、Consumer の 1.63 ~ 1.70 に比べて 1.63 ~ 1.68 とほとんど差がなくなっている。もしも Operator が任意の 2 つの Producer との間に配置できるとする (DHTIntersec.-Split) と、1.54 ~ 1.61 と多少改善するが、相対的には小さい。

なお、著者らの提案方式は GATech トポロジにおいて 1.09 の効率の配置実現するとしている。が、提案方式を実現するためのコストが不明 (仮想空間にノードを配置した上でその中でのノード間の距離をもとに計算を行うらしい)。