

ITRC meet23

2008年 5月 15-16日, 名古屋大学

# Overlay Weaver と その PlanetLab 上での運用

首藤一幸

ウタゴエ / NICT



# 内容

- Peer-to-Peer コンテンツ配信技術 (ウタゴエ社)
  - ライブストリーミング
  - IPTV
- オーバレイ構築ツールキット
  - Overlay Weaver
    - 産総研グリッド研究センター由来。  
2006年度からは個人で継続。

# Peer-to-Peer ライブ配信

実ネットワークに適応するオーバレイマルチキャスト 放送基盤

- 映像・音声の**ライブ**配信が可能。

ここが技術的に面白い。メッシュ型 (⇔ツリー型)  
cf. YouTube, ○○動画

- 配信元の**ネットワーク帯域幅**をほとんど**食わない**。

⇒ 極めて**低い配信コスト**

⇒ **誰でも発信**

cf. サーバクライアント, CDN

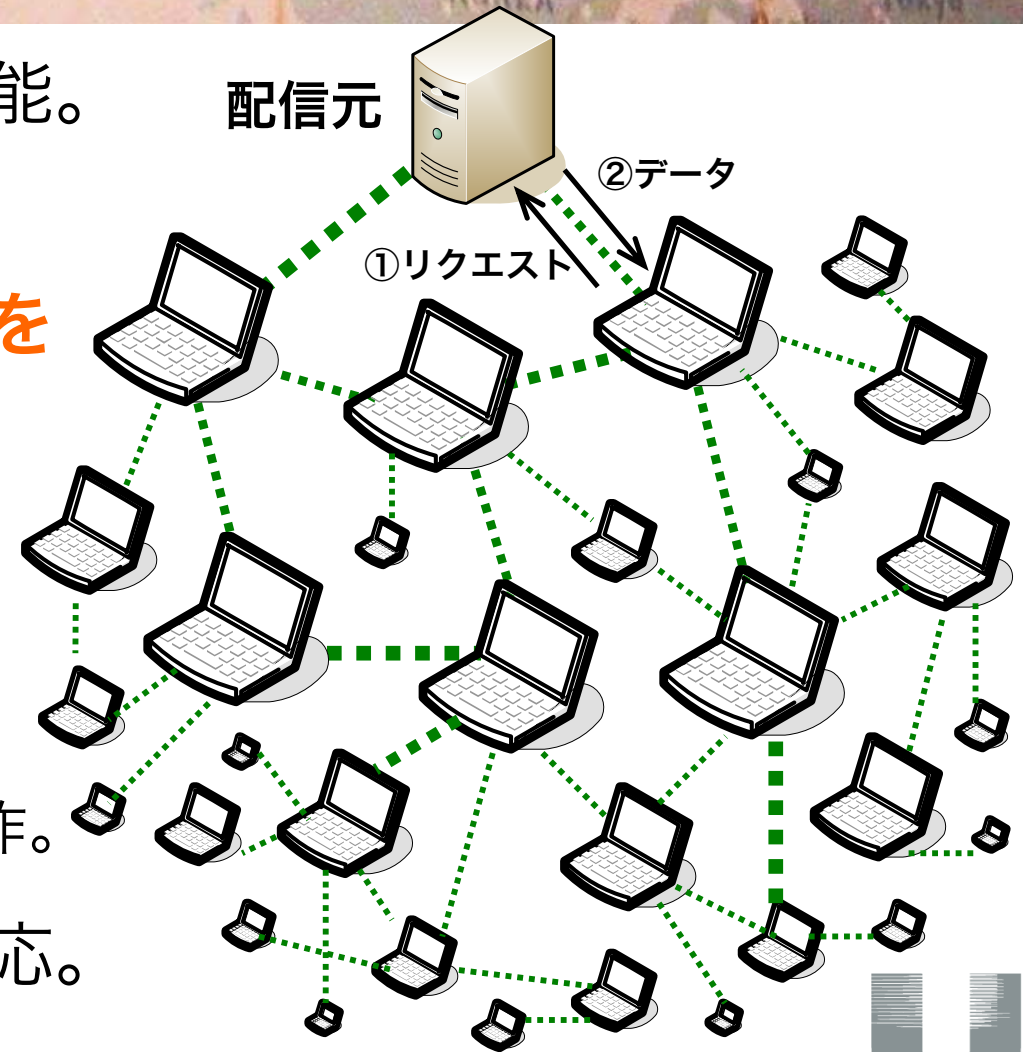
- 1万ノード**での動作実績。

(PC クラスタ上)

実地 & PlanetLab上でも大規模動作。

- 様々なネットワーク環境に適応。

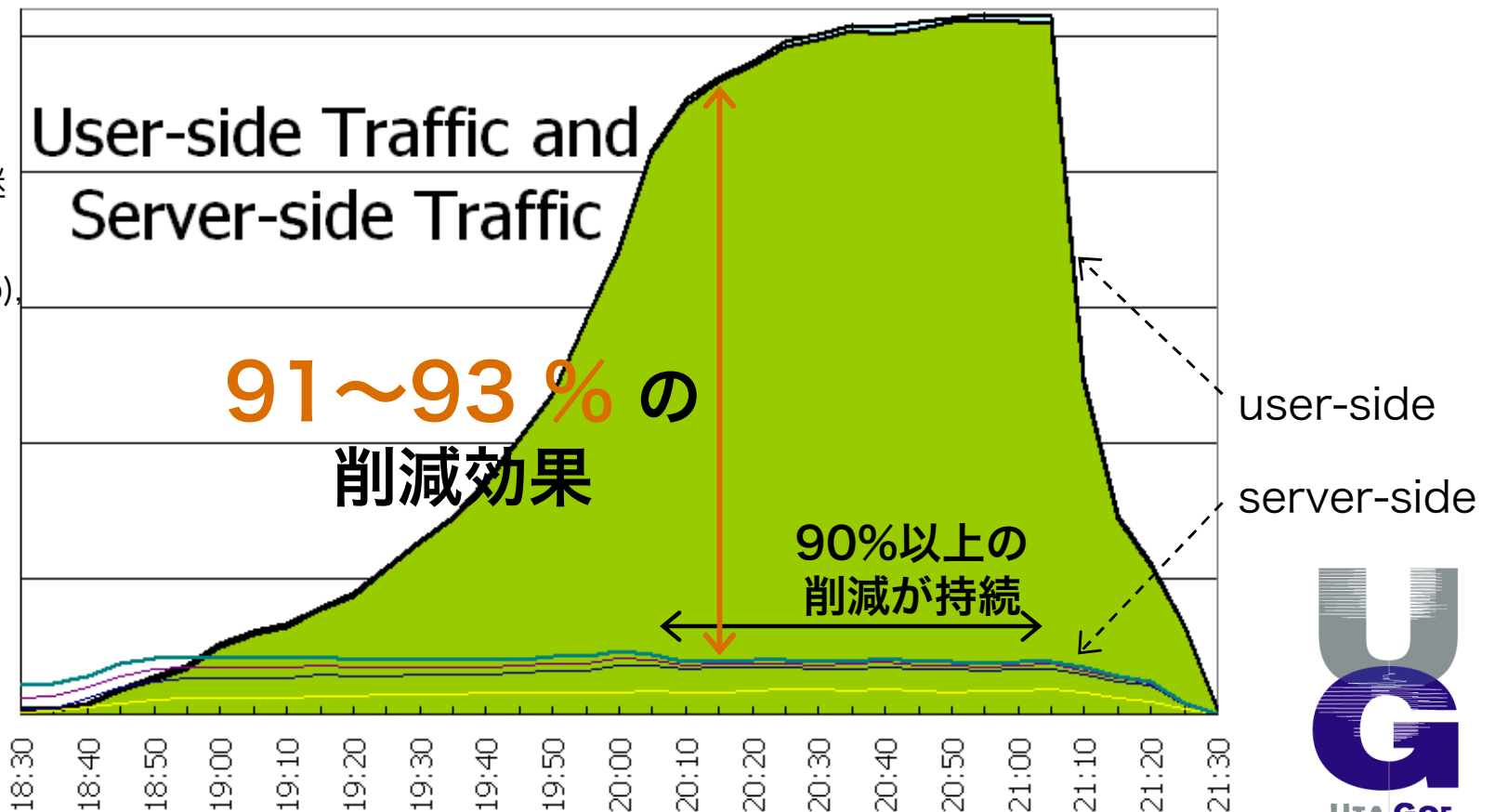
xDSL~光ファイバ, ファイアウォール/NAT



# トラフィック削減効果

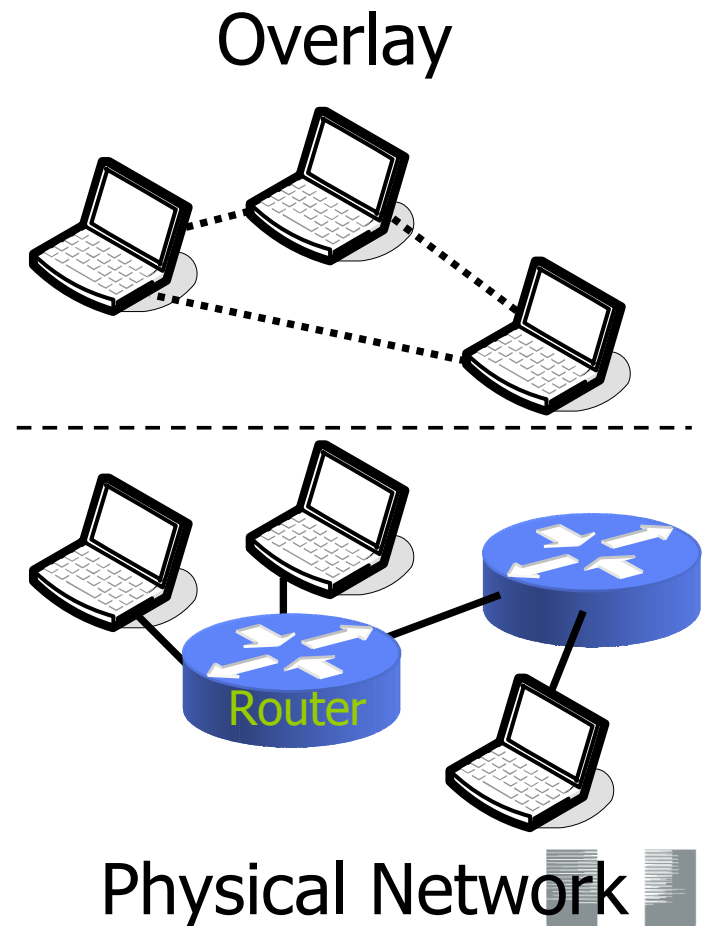
- トラフィック削減効果 90～95 %
  - ウタゴエ社 UG Live 使用

2007年 11月 26日  
accessのtalk about生中継  
配信: Jストリーム,  
Castella (www.castella.jp)  
ウタゴエ



# オーバレイネットワーク (overlay network)

- 下位ネットワークとは独立したトポロジを持つ上位ネットワーク
  - UUCP, NetNewsサーバのトポロジ
  - 電話網上のインターネット
  - インターネット上のVPN, CDN
  - MBone, 6bone, ...
  - GENI
    - 米国 NSF のテストベッド構築プロジェクト
    - 研究者は、自身の“net”を構築
    - PlanetLab からの強い影響
  - ここでは特に、**大規模な自律分散系**
    - それって (pure) P2P ?
      - 多ノード/低管理コスト
      - 自律的にオーバレイを構成

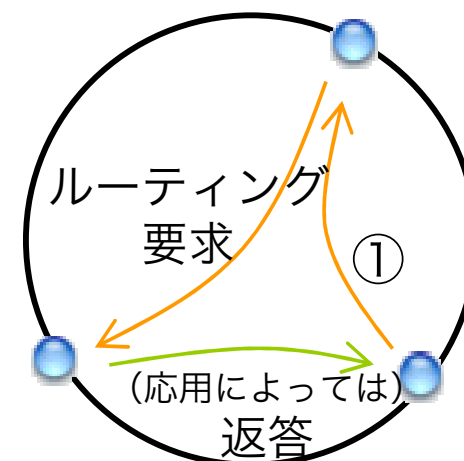
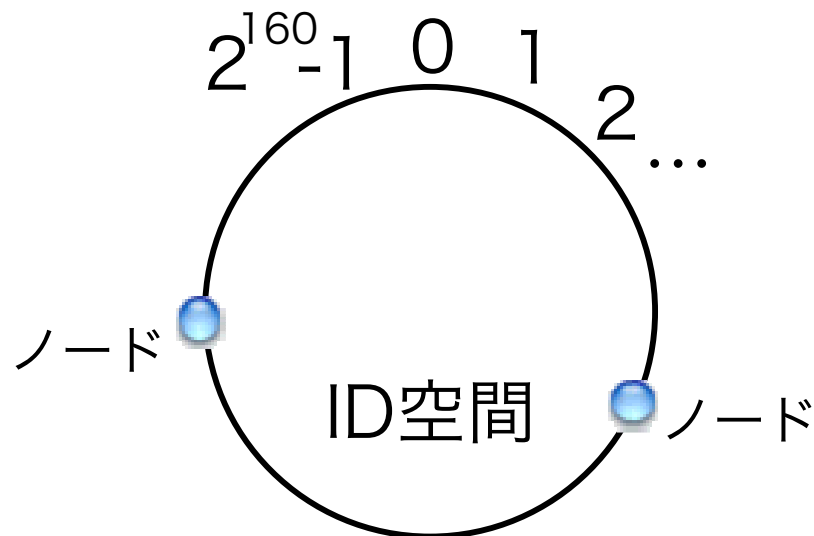


# Unstructured と Structured

- 非構造化 (unstructured) オーバレイ
  - 例： **Gnutella** ネットワーク, Winny ネットワーク
  - 誰を隣接ノードとするか、 **トポロジに制約がない**。
  - 存在するオブジェクトは、 **発見できる可能性がある**。
  - 一般に、効率は良くないが、柔軟な検索が可能。
- 構造化 (structured) オーバレイ
  - 例： **DHT** (分散ハッシュ表) のネットワーク
  - 誰を隣接ノードとするか、 **トポロジに制約がある**。
  - 存在するオブジェクトは、 **(たいてい) 発見できる**。
  - 一般に、効率は良いが、柔軟な検索が苦手。

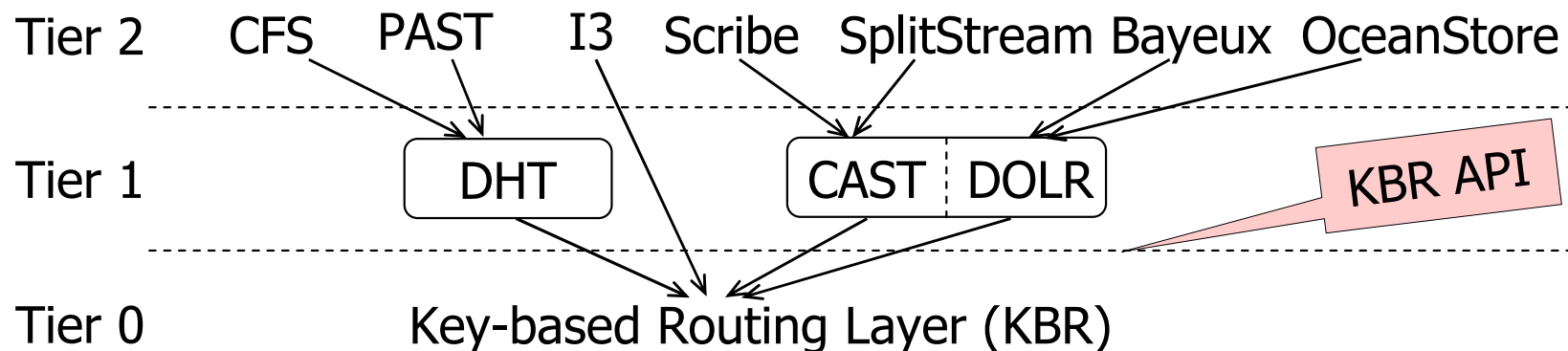
# 構造化オーバレイの基本

- ノード（計算機）とオブジェクトの両方にIDが振られる。
  - IDはたいてい整数値。160ビットだったり 128ビットだったり。
  - オブジェクト：任意の文字列だったり、ファイルだったり、プロセスだったり。
- ノードは、ID空間中のある範囲を受け持つ。
  - だいたい、ノードのIDと数値的に近い範囲を受け持つ。
- IDを宛先としてルーティングが行われ、その行き着く先は、そのIDを受け持つ担当ノードとなる。



# 構造化オーバーレイの基本

- 肝は **routing / forwarding**
- 資源にかかる負担が、ノード数  $n$  として  $O(\log n)$ 。
  - forwarding 時のメッセージ数など。
  - cf. unstructured オーバレイ上の flooding
- 様々なアルゴリズムが提案されてきた。
  - CAN, Chord, Tapestry, Pastry, Kademlia, Koorde, Broose, Accordion, ORDI, DKS, D2B, Symphony, Viceroy, ...
- 構造化オーバーレイ上に、いろいろなサービスが載る。
  - 分散ハッシュ表 (DHT), マルチキャスト, メッセージ配送, ...



Cited from [Dabek03]



# 構造化オーバーレイの応用例

- BitTorrent クライアント
  - eDonkey2k (hybrid P2Pプロトコル) から派生した Kad Network。eMule 等が実装。
  - BitTorrent、Azureus のトラッカーなし動作
    - Azureus: 100万以上のノード
- DNS
  - “A Comparative Study of the DNS Design with DHT-Based Alternatives”, Proc. INFOCOM 2006.
- その他、もろもろの名前解決や位置解決
  - ホスト名 → IPアドレス, 名前 → 電話番号, 曲名 → 楽曲ファイルやそのURL などなど。

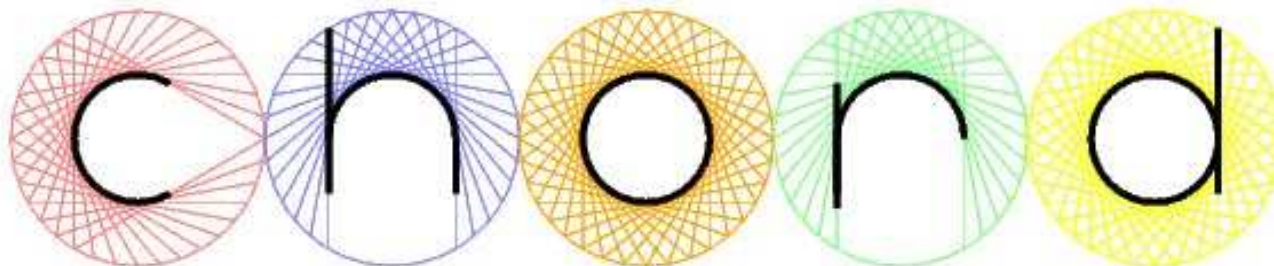
# オーバレイ織り器

Overlay  
Weaver

- (構造化) オーバレイ関係の名前には、糸を織るというアナロジが見られる。
  - Chord -- 弦
  - Tapestry -- 壁掛けの織物
  - Pastry -- 練り粉, パイの皮 ??
- Overlay Weaver
  - (structured) オーバレイの織り器
  - Weave -- 織る, 編む
  - Weaver -- 織り手, 織工, 編む人



Weaver

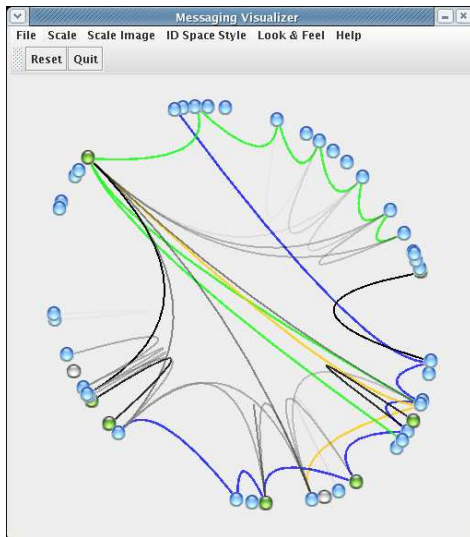


Tapestry of the Grid

# Overlay Weaver

## ● 構造化オーバーレイライブラリ

- 記述言語は Java。
- 約 3 万ステップ。
  - 4 万 8 千行。
- Apache License 2.0。
  - いろいろな目的に使いやすい。



可視化ツール

## ● 特徴

- 分散ハッシュ表 (DHT) だけじゃない。アプリケーション層マルチキャストも。
- ルーティングアルゴリズムを差し替え可能：  
Chord, Kademlia, Koorde, Pastry, Tapestry
- コーディングなしで運用や実験が可能。
  - サンプルツール (DHTシェル等) を使った運用。
  - エミュレータを使った実験: メッセージ数やホップ数の計測。PC 1 台で 70,000 ノード動作。
- 動作状態を可視化して楽しめる。デモできる。
- XML-RPC 経由で DHT を使える。
  - Bamboo, OpenDHT と同じプロトコル → 同じクライアントが使える。
- 地球規模テストベッド PlanetLab 上で運用している。
  - 35 ヶ国 418 台

# オープンソースソフトウェア、 研究プラットフォームとして

Overlay  
Weaver

- <http://overlayweaver.sf.net/> (SourceForge)

- 2006年1月17日、リリース。
- Apache License 2.0

- 状況 (2008/5/17 23時)

- ダウンロード 9,500件。
- メーリングリスト登録数
  - 英語 71名, 日本語 90名
- Mixi コミュニティのメンバー数 152名

- 研究プラットフォームとして

- ブラジル, ロシア, 英国, ハンガリー, ...
- 第三者が、ルーティング方式を追加実装:  
Symphony, EpiChord, ...
- 第三者による利用 / 応用多数:  
RDFのDB(産総研), XMLデータ検索(筑波大), 複製数の影響  
(U Federal do Rio Grande do Sul), 複数サーバの横断検索  
(お茶大), ウェブアクセス動向集計(東工大), 複製配置手法  
(早大), コミュニケーション向けユーザ管理サーバ(日立),  
多次元範囲検索システム(NEC), ...

オーバーレイ構築ツールキット

## Overlay Weaver

[English | Japanese]

概要

Overlay Weaver はオーバーレイ構築ツールキットです。アプリケーション開発に加えて、オーバーレイのアルゴリズム設計もサポートします。

アプリケーション開発者に対しては、分散ハッシュ表 (DHT) やマルチキャストといった高レベルサービスに対する共通 API を提供します。

ウェブサイト

### スクリーンショット

#### Messaging Visualizer

円 直線 曲線 渦巻 格子

50 のノードと、それらの間の通信が可視化されています。マルチキャストのための配送木も色付きの線として描かれています。ここで全ノードと Messaging Visualizer は、エミュレータの上で計算機 1 台上で動作しています。Messaging Visualizer は実ネットワーク上でも動作します。

#### Messaging Visualizer と 300 の (仮想) ノード

円 格子

# 研究プラットフォームとして

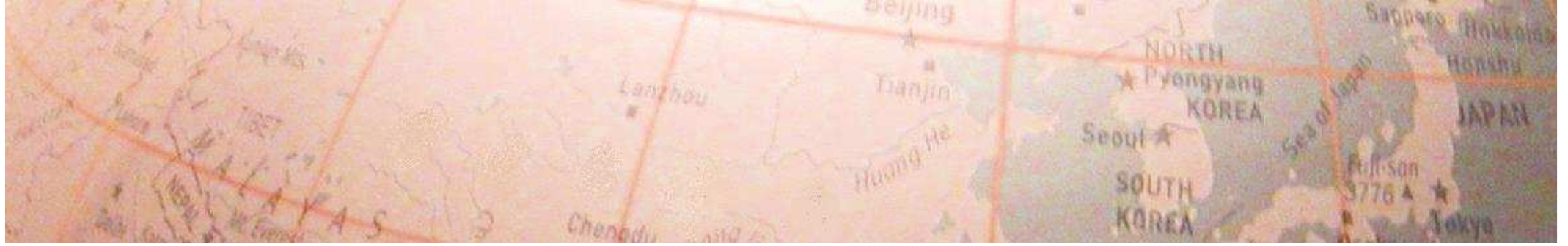
- churn 耐性の評価と向上手法
  - 首藤, “下位アルゴリズム中立なDHT実装への耐churn手法の実装”, 情処ACS論文誌, 2008年3月.
- collective routing
  - 首藤, 中尾, “構造化オーバーレイでの一括ルーティング”, OS研究会 (SWoPP), 2008年 8月.

# PlanetLab 上での運用

- 実機上、できればインターネット上で、実際に大規模な試験・運用を。
- PlanetLab: 地球規模テストベッド
  - PC 2台を供出すれば、使える。
    - プロジェクトごとにVM (Linux-VServer) を用意できる。
  - ただし、企業の fee は \$25,000 / 年～！
  - 客員研究員として、大学を PlanetLab に参加させた。
  - 最大で、35ヶ国 418台で動作。

Overlay Weaver Node Status	
<b>Node Information</b>	
URL:	<a href="http://planetlab. ....it:3998/">http://planetlab. ....it:3998/</a>
Node ID:	9e0074f964d11454587484e4bb5ca356658da1c8
Lookup algorithm:	Chord
Lookup style:	Iterative
# of stored items:	0
<b>Routing Table</b>	
<b>Predecessor</b>	
<a href="http://p. ....fr:3998/">http://p. ....fr:3998/</a> 902cd17e540547a50b5c12c680a4b2b985947972	
<b>Successor List</b>	
<a href="http://. ....cn:3998/">http://. ....cn:3998/</a>	9e2a85116bbe229e69399dd8f4470ad8a24259a
<a href="http://. ....cn:3998/">http://. ....cn:3998/</a>	a4b5eda65b7726e858d8fbd0fb364fbd84231a61
<a href="http://. ....cn:3998/">http://. ....cn:3998/</a>	a6c1c54334ef2a7d2c5e303094393e978692b30a
<a href="http://. ....jp:3998/">http://. ....jp:3998/</a>	ab028d1caad8482a8777bf93216515f4e02e19a4
<a href="http://. ....edu:3998/">http://. ....edu:3998/</a>	b010cfa1ec63dff04123eac81fee105f033c4ec
<a href="http://p. ....com:3998/">http://p. ....com:3998/</a>	b41d267be23ef5ecb5b1ef659f54a925d0bef37
<a href="http://p. ....jp:3998/">http://p. ....jp:3998/</a>	b48851667b142d4253aeba48381e5402e0c8ad5c
<a href="http://p. ....cn:3998/">http://p. ....cn:3998/</a>	b591843d3cc42317a8e065e80a8089c90320cd27
<b>Finger Table</b>	
1	<a href="http://. ....cn:3998/">http://. ....cn:3998/</a>
151	<a href="http://. ....cn:3998/">http://. ....cn:3998/</a>
156	<a href="http://. ....edu:3998/">http://. ....edu:3998/</a>
158	<a href="http://. ....net:3998/">http://. ....net:3998/</a>
159	<a href="http://p. ....fr:3998/">http://p. ....fr:3998/</a>
160	<a href="http://. ....cn:3998/">http://. ....cn:3998/</a>

PlanetLab上で運用しているノードの  
ウェブインタフェース

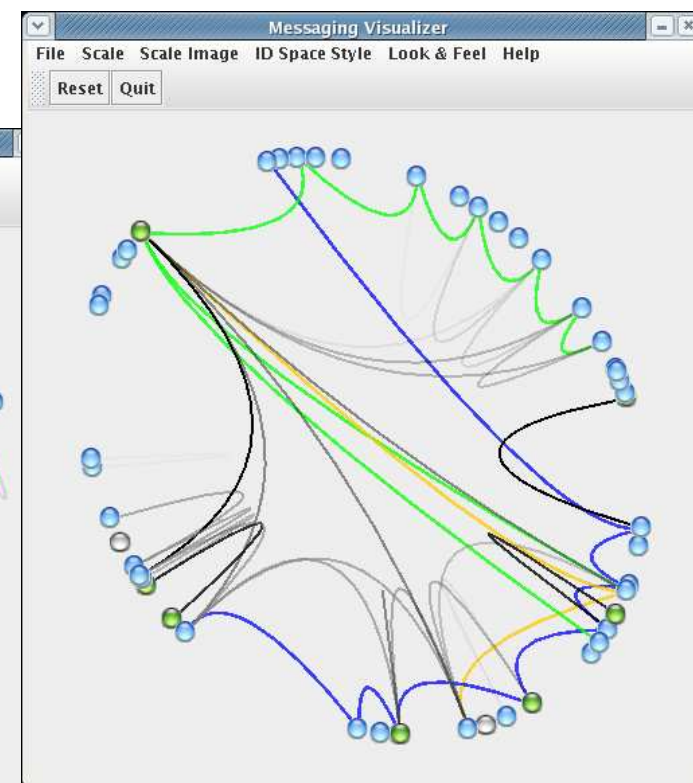
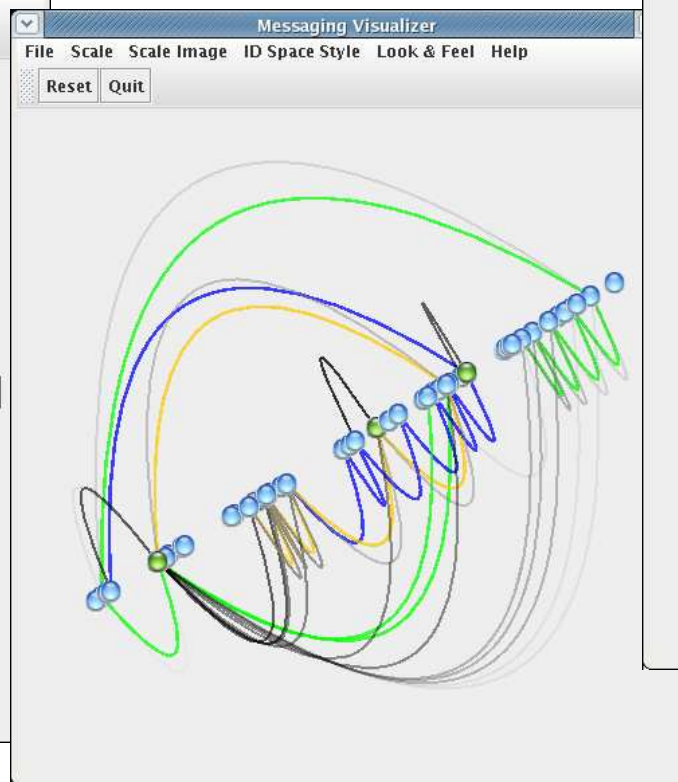
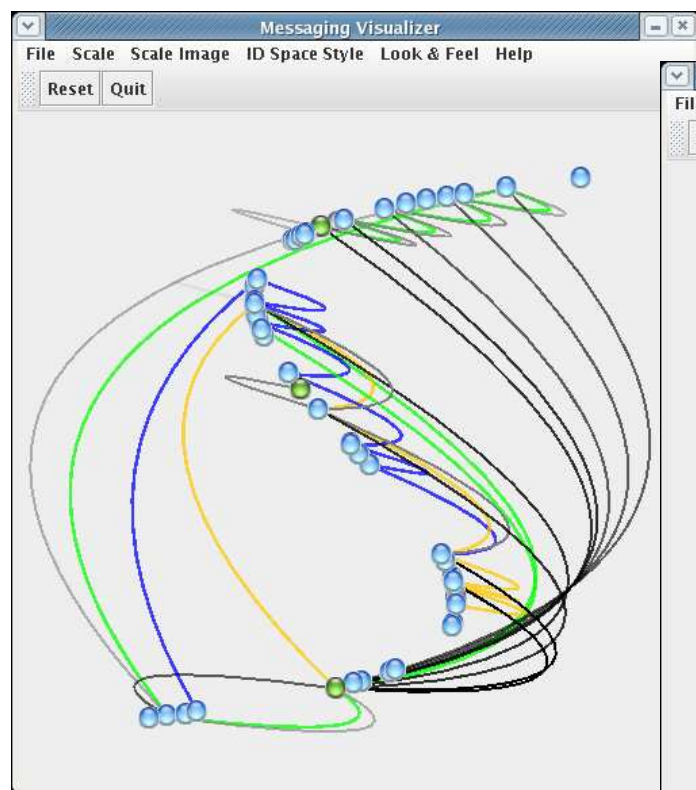


# デモ

- 可視化ツール: Overlay Visualizer
- PlanetLab 上の DHT に対する操作

# Overlay Visualizer

- ノード, メッセージ, マルチキャストの配送木を、実行中に可視化。
- 実ネットワーク / エミュレータどちらでも動作。





# ノードのウェブインタフェース

- ノードについての諸元。

## Overlay Weaver Node Status

### Node Information

URL: <http://planetlab0.otemachi.wide.ad.jp:3998/>  
Node ID: 7a17c1434ef5858f0fbfe052a08b318bbb385433  
Lookup algorithm: Chord  
Lookup style: Iterative  
# of stored keys: 0

## Routing Table

### Predecessor

# • Chord の経路表。

<http://pub2-s.ane.cmc.osaka-u.ac.jp:3998/> 793bfb7e552e859a63a2eff10d956ebbe678e226

### Successor List

<a href="http://planetlab1.cs.duke.edu:3998/">http://planetlab1.cs.duke.edu:3998/</a>	7abe663cbc720375b37dd7dcb71e7153aa85da0e
<a href="http://planetlab04.cs.washington.edu:3998/">http://planetlab04.cs.washington.edu:3998/</a>	7b638aa42cf93c1f81c8caecdc5e4cf2e4cdabc0
<a href="http://planck227.test.ibbt.be:3998/">http://planck227.test.ibbt.be:3998/</a>	7bc89d2027953a7bb7c7872290afe025abfce687
<a href="http://planet1.l3s.uni-hannover.de:3998/">http://planet1.l3s.uni-hannover.de:3998/</a>	7f9b10393062377effc2785cf705ab4aad953bcb
<a href="http://planet6.berkeley.intel-research.net:3998/">http://planet6.berkeley.intel-research.net:3998/</a>	8376275b1547a93b518075bb3af9f53501408455
<a href="http://planetlab1.cesnet.cz:3998/">http://planetlab1.cesnet.cz:3998/</a>	83aa8c5cdd6dd5f9ff441ed6e0e72be2b2c8910e
<a href="http://planet1.prakinf.tu-ilmenau.de:3998/">http://planet1.prakinf.tu-ilmenau.de:3998/</a>	85ad64a1d1ce6cb55b3b1919b01f4f508dca9536
<a href="http://planetlab-1.unk.edu:3998/">http://planetlab-1.unk.edu:3998/</a>	868c3bbad9adca55d22bdc878c926087f4c7a035

### Finger Table

1	<a href="http://planetlab1.cs.duke.edu:3998/">http://planetlab1.cs.duke.edu:3998/</a>	7abe663cbc720375b37dd7dcb71e7153aa85da0e
153	<a href="http://planetlab-01.kyushu.jgn2.jp:3998/">http://planetlab-01.kyushu.jgn2.jp:3998/</a>	903f7c72ac22c1bb71ab57a0230c167517d39e96
158	<a href="http://planetlab4.cs.uiuc.edu:3998/">http://planetlab4.cs.uiuc.edu:3998/</a>	a1f4bfcdcb8d56a21ca007de183e3a24d503e3a1
159	<a href="http://planetlab02.sys.Virginia.EDU:3998/">http://planetlab02.sys.Virginia.EDU:3998/</a>	bda4cc7207314bf6ea874a04a37d8653954da9ac
160	<a href="http://pli2-pa-1.hpl.hp.com:3998/">http://pli2-pa-1.hpl.hp.com:3998/</a>	fd04d9785c525f614ae3384d5f181815d93b827d

# ウェブインタフェース

- 分散ハッシュ表 (DHT) への put, get, remove

## Put, Get and Remove Operations

operation	key	value	TTL (sec)	secret	
get	<input type="text" value="foo"/>				<input type="button" value="submit"/>
put	<input type="text" value="foo"/>	<input type="text"/>	<input type="text" value="600"/>	<input type="text"/> (option)	<input type="button" value="submit"/>
remove	<input type="text" value="foo"/>	<input type="text"/>	<input type="text" value="(option)"/>	<input type="text"/>	<input type="button" value="submit"/>

# ウェブインタフェース

- DHT からの get の結果。

## Results

Get results:

key: foo

value: bar

## Route

Hop	Node	ID	time
0	<a href="http://planetlab0.otemachi.wide.ad.jp:3998/">http://planetlab0.otemachi.wide.ad.jp:3998/</a>	7a17c1434ef5858f0fbfe052a08b318bbb385433	0
1	<a href="http://planetlab02.sys.Virginia.EDU:3998/">http://planetlab02.sys.Virginia.EDU:3998/</a>	bda4cc7207314bf6ea874a04a37d8653954da9ac	416
2	<a href="http://planetlab2.engr.uconn.edu:3998/">http://planetlab2.engr.uconn.edu:3998/</a>	de9aab0095f26bdf3125a0843dd8e0c83676ae38	1245
3	<a href="http://planet2.cs.ucsb.edu:3998/">http://planet2.cs.ucsb.edu:3998/</a>	f131d7ec5a7b8b2ce6631506151dcf9af4d8ef16	2074
4	<a href="http://planetlab-1.cs.uh.edu:3998/">http://planetlab-1.cs.uh.edu:3998/</a>	00e64d0fcab4640596d7d63e73aa0842ce6a49e1	2231
5	<a href="http://PLANETLAB-2.CMCL.CS.CMU.EDU:3998/">http://PLANETLAB-2.CMCL.CS.CMU.EDU:3998/</a>	080d283be8537508023ad9e070b4778bc1cc51b3	2558
6	<a href="http://planetslug3.cse.ucsc.edu:3998/">http://planetslug3.cse.ucsc.edu:3998/</a>	0b25d99513184b13cacc4ca18f43d3079a8f002c	3273
7	<a href="http://planetlab0.dojima.wide.ad.jp:3998/">http://planetlab0.dojima.wide.ad.jp:3998/</a>	0c66909552301a7d33275cb1716cf418e87b18d3	3715

# PlanetLab の資源利用状況

- CoVisualize

