

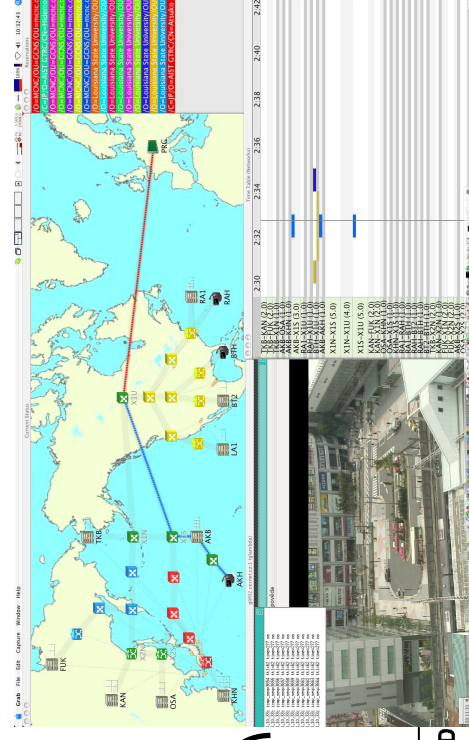
ネットワーク帯域予約とOS仮想化機構を用いた 分散アプリケーションシミュレーション実行環境に向けて

産業技術総合研究所 情報技術研究部門

中田 秀基, 高野 了成
竹房 あつ子, 工藤 知宏

背景

- 計算資源とネットワーク資源を同時に確保して
広域にまたがるアプリケーション実行を行う試
み
 - G-lambda 実験
 - 日米欧のサイトを光パスネットワークで接続してアプリケー
ションを実行
- 複数サイトにまたがるアプリケーション実行の
困難さが浮き彫りに
 - 特殊なアプリケーション/
ミドルウェアの利用が必須
 - 既存のアプリケーションを
そのまま動かすことはできない



背景 (2)

- なぜアプリケーション実行が困難なのか？
 - ネットワークが非対称
 - NAT/Private IP
 - G-lambdaではこれは問題にならない
 - ネットワークもプロビジョンするため
 - 実行に必要な情報が実行時まで決まらない
 - ネットワークを動的に取得する場合のIPアドレスなど
 - 各サイトが不均質
 - 各サイトの管理者は異なる
 - 多くのアプリケーションに共通する不備
 - マルチホームネットワークへの不対応

本研究の目的と成果

- 計算機とネットワークを同時に確保するグリッドにおいて、任意のアプリケーションが容易に実行できる環境を実現する
 - アプリケーション実行環境のアーキテクチャを提示
 - Linux-VMとaufsを用いたプロトタイプ実装によりアーキテクチャの有効性を確認

発表の概要

- アプリケーション実行フレームワークの必要性
- 要件の整理と設計
- 実装
- 関連研究
- おわりに

G-lambda 実験

The screenshot displays a network simulation interface for a G-lambda experiment. The main window shows a map of Asia and Europe with various nodes and connections. The nodes are labeled with codes such as FUK, KAN, OSA, KHN, XZIN, XZS, XIN, XIS, XIU, X1S, X1U, AKH, AKB, LAI, BT2, RAI, RAH, and PRG. The connections are represented by lines of different colors (blue, red, yellow, green) and thicknesses, indicating different network paths or states.

At the top of the interface, there is a "Reservations" section with a list of reservation codes:

- OU=MCNC/OU=GCNS/OU=mcnc.org
- C=JP/O=AIST/GTRC/CN=Hidekoto
- OU=MCNC/OU=GCNS/OU=mcnc.org
- OU=MCNC/OU=GCNS/OU=mcnc.org
- OU=MCNC/OU=GCNS/OU=mcnc.org
- OU=MCNC/OU=GCNS/OU=mcnc.org
- OU=MCNC/OU=GCNS/OU=mcnc.org
- OU=MCNC/OU=GCNS/OU=mcnc.org
- OU=Louisiana State University/OU=C
- OU=Louisiana State University/OU=C
- OU=Louisiana State University/OU=C
- OU=Louisiana State University/OU=C
- OU=Louisiana State University/OU=C
- OU=Louisiana State University/OU=C
- OU=Louisiana State University/OU=C
- C=JP/O=AIST/GTRC/CN=Atsuko Ta

On the left side, there is a "Current Status" section with a list of network statistics:

```

1,10,33: lomp_seq=3554 t1152 time=277 ms
1,10,33: lomp_seq=3555 t1152 time=277 ms
1,10,33: lomp_seq=3556 t1152 time=277 ms
1,10,33: lomp_seq=3557 t1152 time=277 ms
1,10,33: lomp_seq=3558 t1152 time=277 ms
1,10,33: lomp_seq=3559 t1152 time=277 ms
1,10,33: lomp_seq=3560 t1152 time=277 ms
1,10,33: lomp_seq=3561 t1152 time=277 ms
1,10,33: lomp_seq=3562 t1152 time=277 ms
1,10,33: lomp_seq=3563 t1152 time=277 ms
1,10,33: lomp_seq=3564 t1152 time=277 ms
    
```

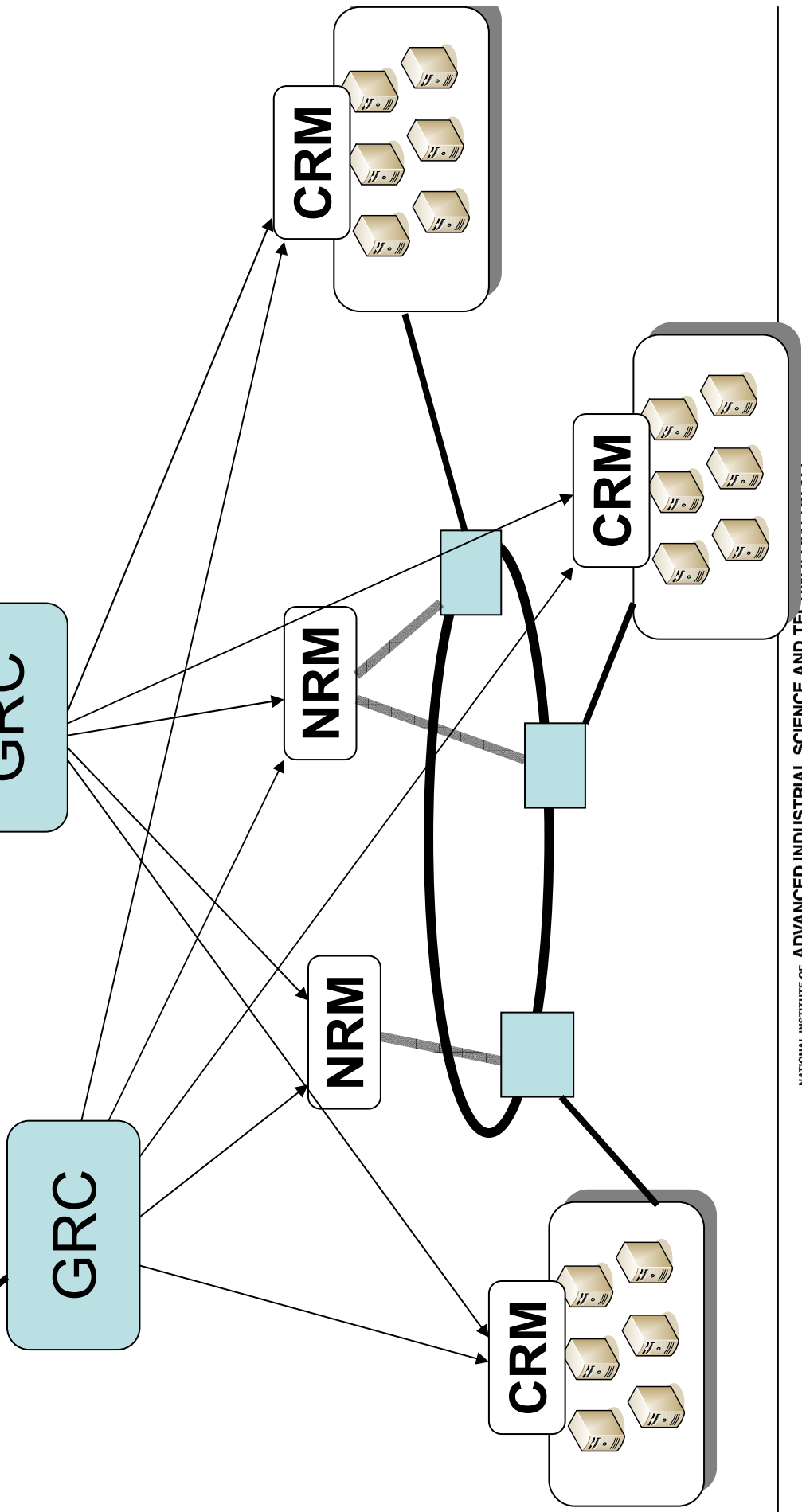
At the bottom, there is a "Time Table (Networks)" section showing a timeline of network activity from 2:30 to 2:42. The timeline is divided into segments for different nodes and connections, with colors corresponding to the network paths shown in the map.

At the bottom right, there is a small inset image showing a street view of a building, likely the National Institute of Advanced Industrial Science and Technology (AIST).



GridARS アーキテクチャ

GRC: Global Resource Coordinator
NRM: Network Resource Manager
CRM: Computing Resource Manager

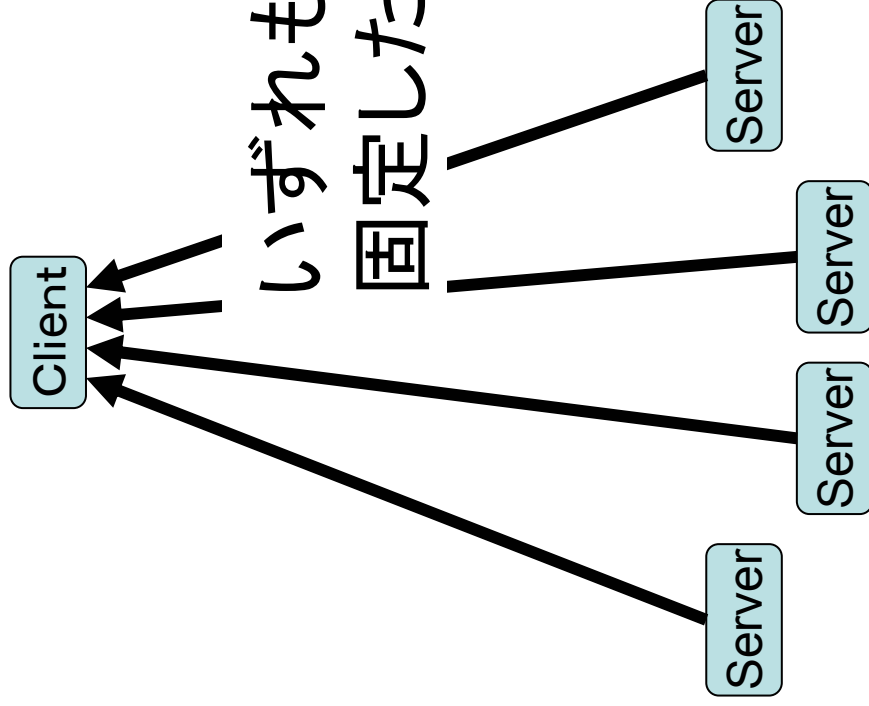


Lessons learned

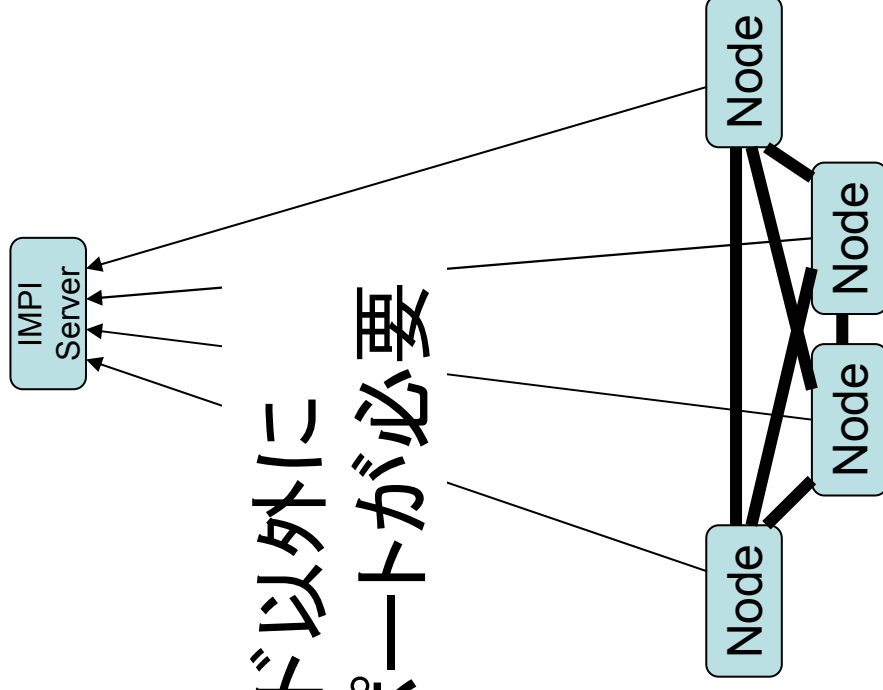
- IPアドレスの共有がむずかしい
 - 動的にCRMのレイヤでローカルに決定されるため
 - 衝突しないような制御も必要
- 各サイトの設定は煩雑
 - ssh keyの分配, ユーザ名の違いなど
 - authorized_keys, known_hosts の登録
- アプリケーションに不備も
 - 複数のインタフェースを想定していない
 - 適当なネットワークアドレスを勝手に使ってしまう
 - 性能低下

GridRPC/GridMPIの実行

GridRPC (Ninf-G)



GridMPI



いずれも予約ノード以外に
固定したノードとポートが必要

アプリケーションシヨン実行環境の設計

- 目的：
 - 既存の並列アプリケーションを変更なしに動作させることができる環境の構築
- 前提：
 - 並列アプリケーションの各サブジョブはsshを用いて起動される
 - ホスト名 (IPアドレス) はファイルから取得
 - アプリケーションはマルチホームネットワークに未対応
 - 各サイトはそれぞれ使用可能なIPアドレスの領域を持つ。この領域は相互に重なっている可能性がある。

システムへの要請

- すべてのノードのIPアドレスのリストがすべてのノードに配備されている
- すべてのノードのsshユーザ公開鍵，ホスト公開鍵が交換されている
- すべてのノードのIPアドレスは重複してはいならない
 - Private IPアドレスを用いる場合これは自明ではない
- アプリケーションプロセスからは単一のネットワークインターフェイスのみが見える
 - マルチホーム未対応アプリケーションへの対応

設計の概要

- 実行環境を整えるパイロットシヨブ (NodeManager) を利用
 - ローカルバッチキューイングシステムから直接ユーザシヨブを実行せず，パイロットシヨブを起動
 - パイロットシヨブが実行環境を整備，ユーザシヨブを実行
 - IPアドレス，鍵の交換など
- 実行環境には仮想化環境を利用
 - ユーザシヨブにみせるネットワークインターフェイスを制御
 - ユーザのホスト環境を変更せず，公開鍵の変更が可能

IPアドレスの決定と収集分配

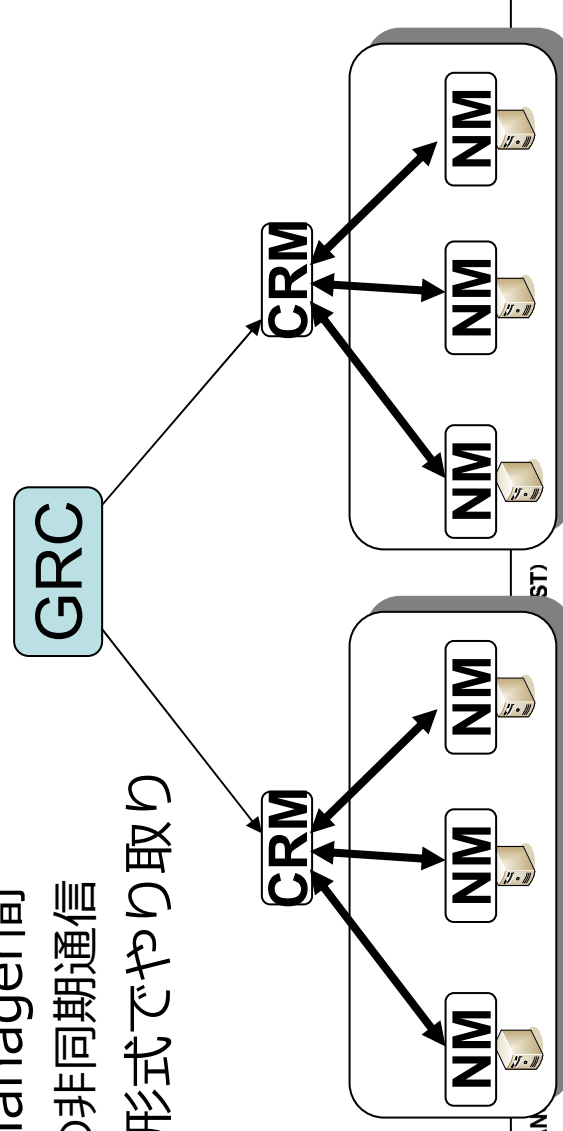
- 事前に、各ノードで利用可能なIPアドレスレンジの集合を公開
- 予約時に、GRCは各サイトで利用するアドレスレンジを決定，CRMに指定
- 実行開始時
 - CRMは利用するIPアドレスを決定，各NodeManagerに対して指定
 - CRMは利用するIPアドレスをGRCに報告
 - GRCはすべてのサイトのIPアドレスをまとめたりすトをCRMに返送
 - CRMはそれをNodeManagerに返送

sshユーザ，ホスト公開鍵の交換

- NodeManagerが実行環境構築の一環としてsshホスト鍵ペア，ユーザ鍵ペアを新たに生成
- 実行時
 - CRMが公開鍵をすべてのNodeManagerから収集，GRCに送信
 - GRCは，すべてのCRMからの情報を集計してCRM経由で分配
 - authorized_keys
 - known_hosts

プロトタイプ実装

- GRC/CRM/NodeManager から構成されるシステム全体のモックアップ
 - フレームワーク構成の妥当性を確認
 - GRCモック/CRMモックはJavaでプロトタイプ
 - 通信はRMIを用いた同期通信
 - NodeManagerはPythonで記述
 - CRMモック/NodeManager間
 - テキストベースの非同期通信
 - 構造データはJSON形式でやり取り



仮想化層の選択

- 計算機仮想化
 - いわゆる‘VM’
 - ホストとゲストは完全に分離
 - 負荷は大きい - メモリ, I/O速度など
 - Ex. VMware, Xen, KVM
- OS仮想化
 - OSレベルのサンドボックス
 - ホストとゲストはカーネルを共有
 - 軽量.
 - Ex. Jail, Solaris zone, Linux-VServer, OpenVZ

Linux向けOS仮想化機構

- Linux-VServer
 - PlanetLabで使用されている
 - ホスト・ゲスト間のネットワークの分離が不完全
 - ホストがanyでlistenするとゲストアドレス宛のconnectがホストに捕捉されてしまう
- OpenVZ
 - Parallels社の商用製品Virtuozzoのオープンソフト版
 - ネットワークの分離はより完全
 - パケットをホストがルーティングしている

ネットワーク速度比較

- 1GbEと10GbE で測定
 - 1GbE -Broadcom BCM5722, 10GbE-Myricom Myri-10G
- 使用バージョンは以下のとおり
 - Vanilla -2.6.27-9 server
 - VServer- 2.6.27-10-vz2.3.0.36.2
 - OpenVZ- 2.6.27.10-openvz
 - Xen -3.3, Dom0 2.6.26.1

1GbE [UDP/TCP]

	Mbps	送信側CPU	受信側CPU
vanilla	941.38 / 957.2	0.76 / 1.99	0.33 / -0.25
VServer	941.37 / 957.2	1.00 / 1.77	0.16 / -0.19
OpenVZ	941.07 / 957.2	1.90 / 3.02	1.57 / 2.42
Xen	949.19 / 957.2	5.99 / 0.25	4.54 / 3.50

ネットワーク速度比較

- 1GbE ではどの仮想化機構も十分な性能
- 10GbE では OpenVZ, Xen は性能が十分
 - Xen, OpenVZ は CPU 負荷が高い

測定方法が異なるのであまり参考にならない。
上3つは8コアで100.0
Xenのみ1コアで100.0

送信側のCPUが飽和
(12.5 = 1/8)

[UDP/TCP]

	Mbps	送信側CPU	受信側CPU
vanilla	9525.94 / 8051.2	8.13 / 8.13	7.84 / 7.84
VServer	9521.79 / 8149.9	9.11 / 14.11	7.51 / 7.11
OpenVZ	2049.89 / 3005.0	10.12 / 12.60	2.74 / 3.55
Xen	1011.47 / 6252.5	0.85 / 99.98	1.38 / 8.04

NodeManagerの実装

- Pythonで記述
- 実効uid rootで動作
 - VServerの操作に必須
 - 通常時は実uid nobodyで動作，必要時にrootにseteuidする
- ユーザ権限で動作するLauncherプロセスからsudoで起動
- Launcherが停止すると，自動的に停止
 - バッチキューイングシステムがNodeManagerを間接的に制御することが可能.

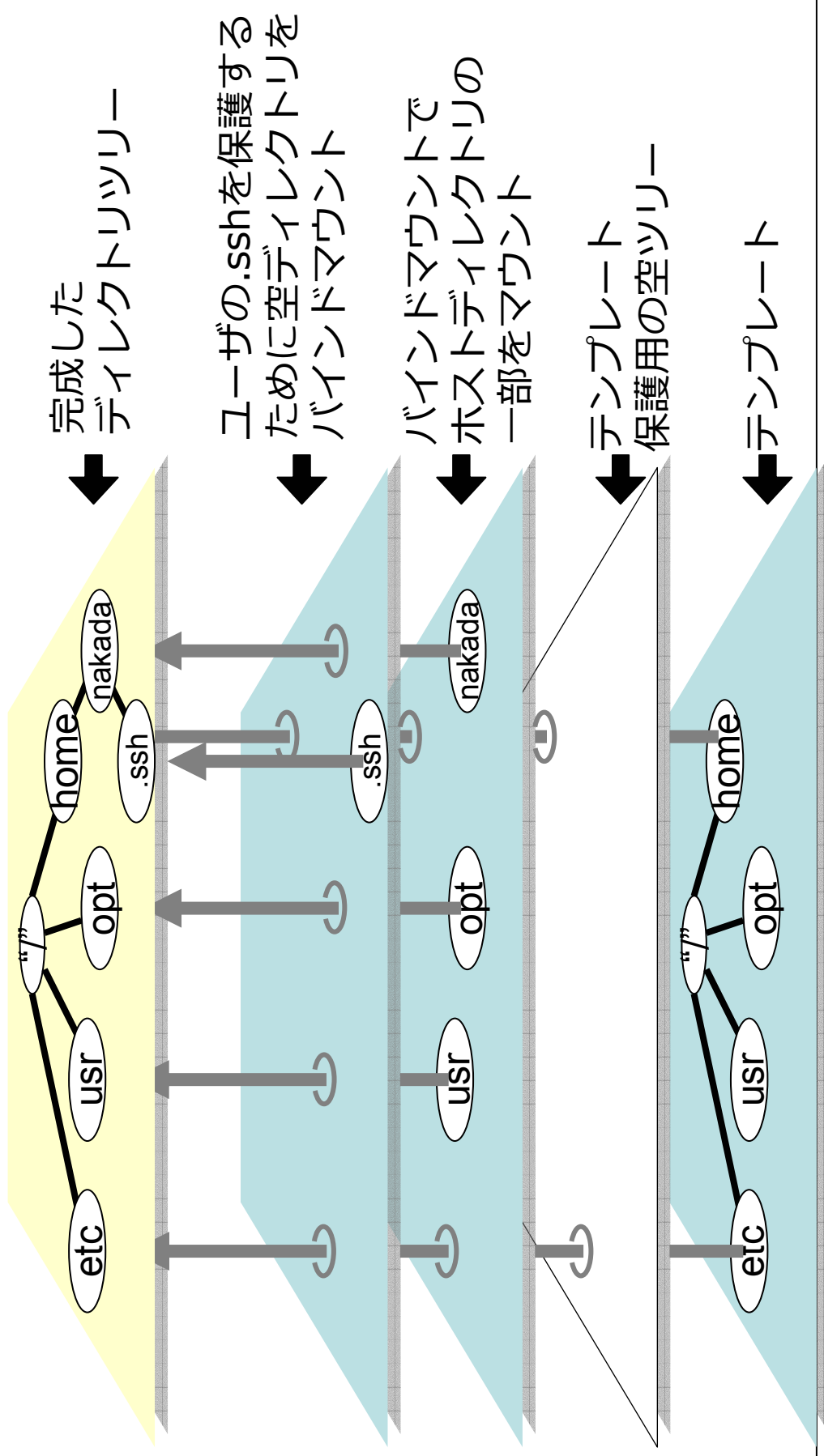
NodeManagerに対するコマンド

- INIT
 - IpAddressなど必要な情報を与えて仮想環境をセットアップ
- GET_KEYS
 - ユーザ, ホストの公開鍵を取得
- SET_KEYS
 - Hostfile, kownhosts, authorized_keysを設定
- RUN
 - 仮想環境の起動
- EXEC
 - 仮想環境内でユーザのジョブを起動
- SHUT_DOWN
 - 仮想環境を停止
- GET_STATUS
 - 現在のNodeManagerのステータスを取得.
 - ステータスは適宜Notify されるので, 基本的には不要

仮想化環境内ファイルシステムの構成

- ホスト環境と共有しながら隔離された環境が必要
 - ホスト環境でアプリケーションをテスト, そのまま実行
 - .sshはホストと分離
- aufs(Another UniFS)を使用
 - スタックカブルファイルシステムの1実装
- バインドマウントを援用
 - Mount -t bind

ディレクトリツリーの構成



システムの動作



関連研究

- PlanetLab
 - 広域ネットワークのテストベッド
 - Linux-Vserverを用いてユーザに環境を提供
 - 比較的長期間の貸与
 - 各ユーザ環境は実環境から隔離
- 尾崎ら[comsys'08]
 - NICT AKARIプロジェクト
 - KVMを利用
 - 仮想化環境でカーネルを入れ替えることを重視
 - 1Gbps程度までは対応可能

おわりに

- **まとめ**
 - 計算機とネットワークを同時に確保するグリッドにおいて、任意のアプリケーションが容易に実行できる環境を実現するためのアーキテクチャを提示
 - プロトタイプ実装を行いアーキテクチャの有効性を確認
- **今後の課題**
 - g-lambdaの各コンポーネントに実装
 - 実験
 - 上記を行ったうえで、広域の実環境で実験を行い有効性を確認

謝辞

本研究の一部は、情報通信研究機構(NICT)の委託研究「新世代ネットワークサービス基盤構築技術に関する研究開発」により実施した。