



**VINCI: Virtual Intelligent Networks for  
Computing Infrastructures**

**An Integrated Network Services System  
to Control and Optimize Workflows  
in Distributed Systems**

**CHEP February 2006**

**Harvey Newman  
and Iosif Legrand  
California Institute of Technology**



# OUTLINE



- ◆ **Introduction**
- ◆ **The MonALISA framework**
  - **Monitoring**
  - **Support for distributed services and Agents**
- ◆ **The VINCI architecture**
- ◆ **Main Services**
  - **End System Agent (LISA)**
  - **Discovery & AAA**
  - **Control of Optical planes**
  - **Interfaces with GMPLS, MLPS, SNMP ...**
- ◆ **Prediction, Learning and Self Organization**



# The Need for Network Services



- **The main objective of the VINCI project is to enable users' applications, at the LHC and in other fields of data-intensive science, to effectively use and coordinate network resources**
- **VINCI dynamically estimates and monitors the achievable performance along a set of candidate (shared or dedicated) network paths, and correlates these results with the CPU power and storage available at various sites, to generate optimized workflows for grid tasks**
- **This should significantly improve the overall performance and reduce the effective costs of global-scale grids**
- **The VINCI system is implemented as a dynamic set of collaborating Agents in the MonALISA framework, exploiting MonALISA's ability to access and analyze in-depth monitoring information from a large number of network links and grid sites in real-time**



## VINCI: A Multi-Agent System



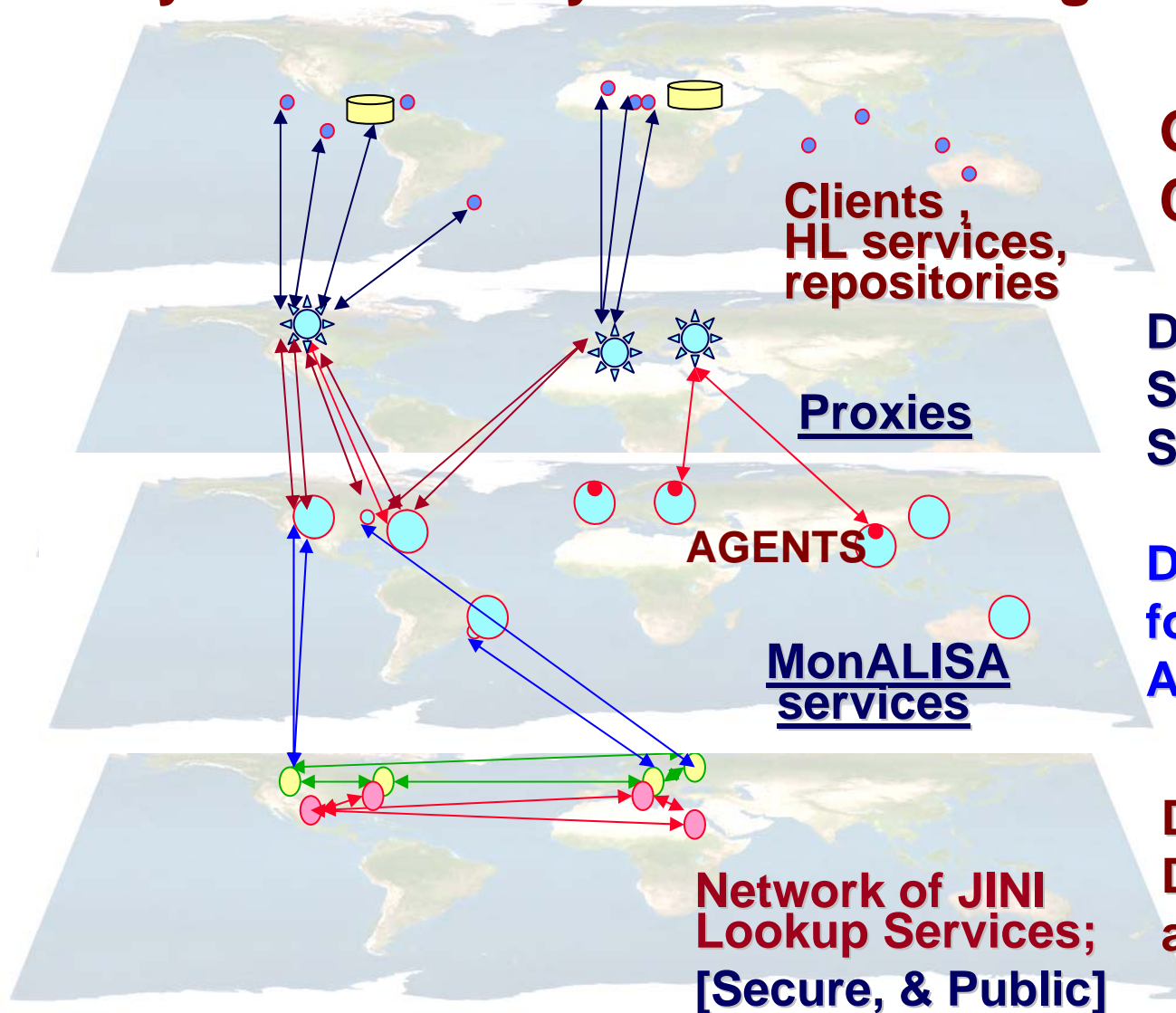
- VINCI and the underlying MonALISA framework use a system of autonomous agents to support a wide range of dynamic services
- Agents in the MonALISA servers self-organize and collaborate with each other to manage access to distributed resources, to make effective decisions in planning workflow, to respond to problems that affect multiple sites, or to carry out other globally-distributed tasks
- Agents running on end-users' desktops or clusters detect and adapt to their local environment so they can function properly. They locate and receive real-time information from a variety of MonALISA services, aggregate and present results to users, or feed information to higher level services
- Agents with built-in “intelligence” are required to engage in negotiations (for network resources, for example), and to make pro-active run-time decisions, while responding to changes in the environment



# MonALISA : An Agent-based System of Distributed Services



**Fully Distributed System with no Single Point of Failure**



**Global Services or Clients**

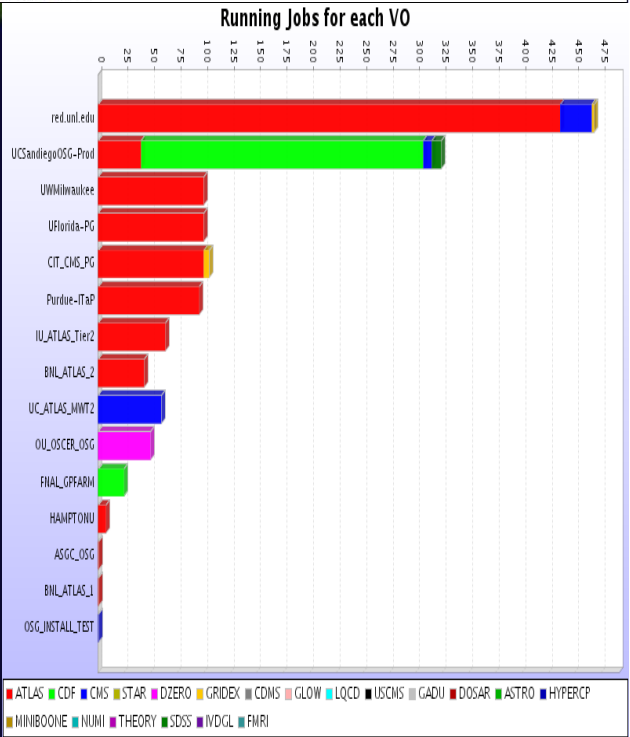
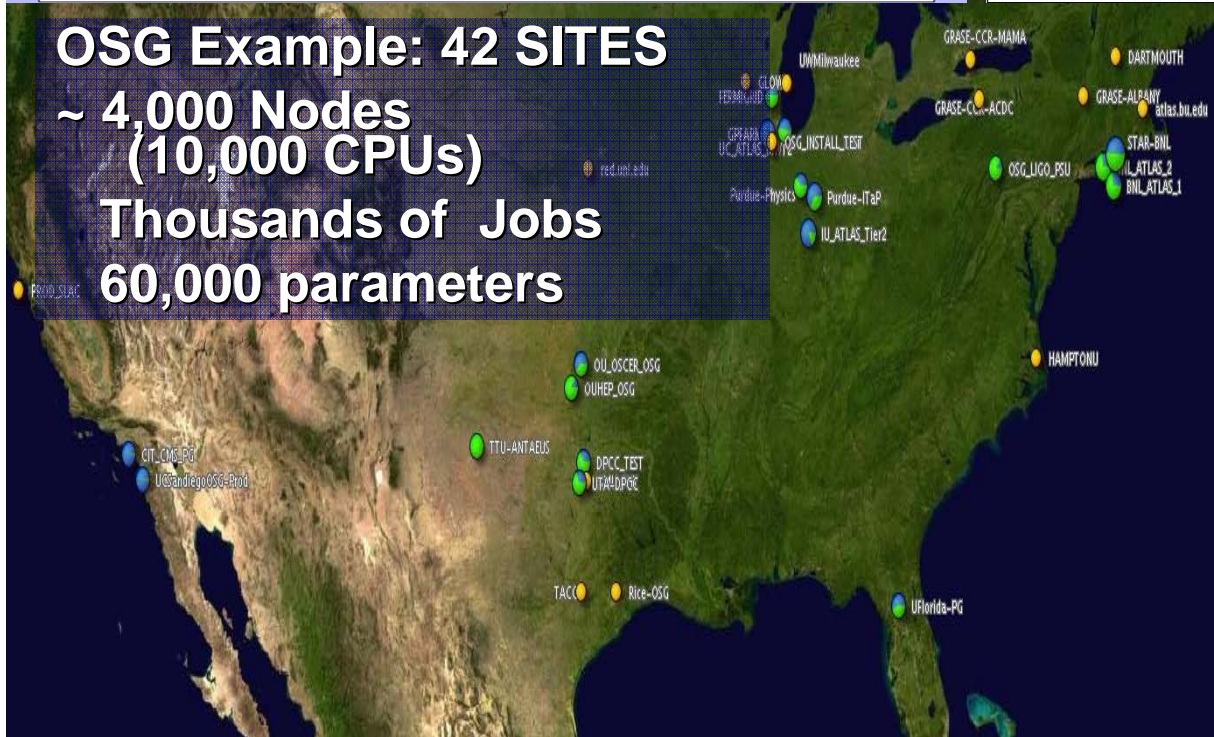
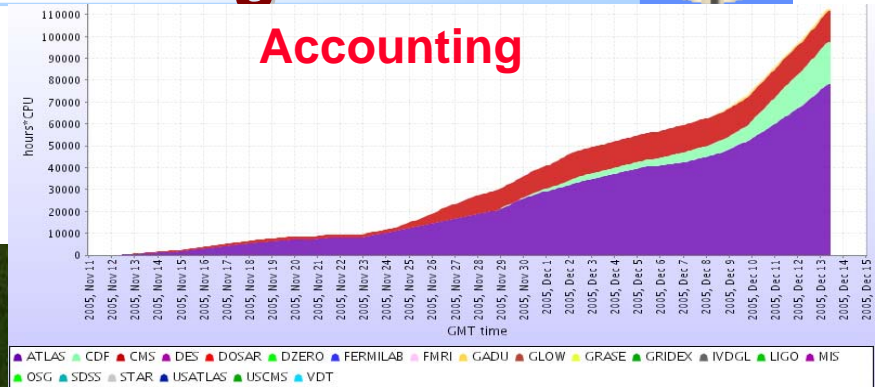
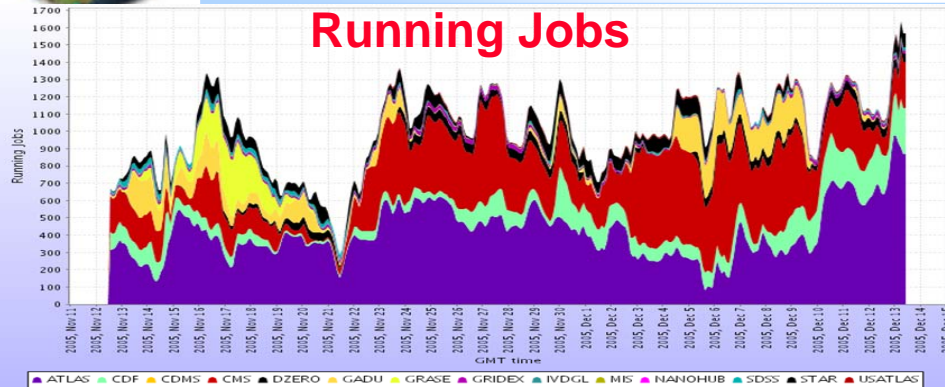
**Dynamic load balancing  
Scalability & Replication  
Security AAA for Clients**

**Distributed System  
for gathering and  
Analyzing Information.**

**Distributed Dynamic  
Discovery- based on  
a lease Mechanism**

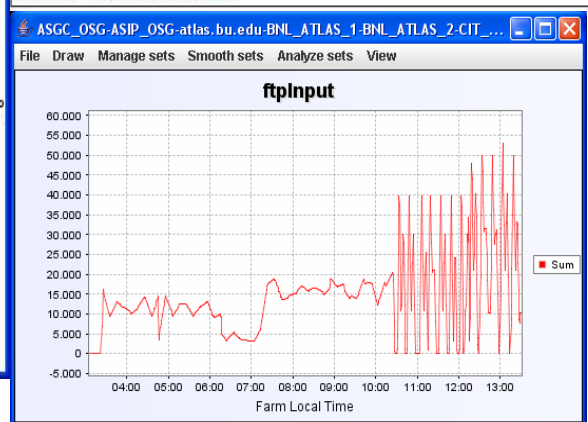
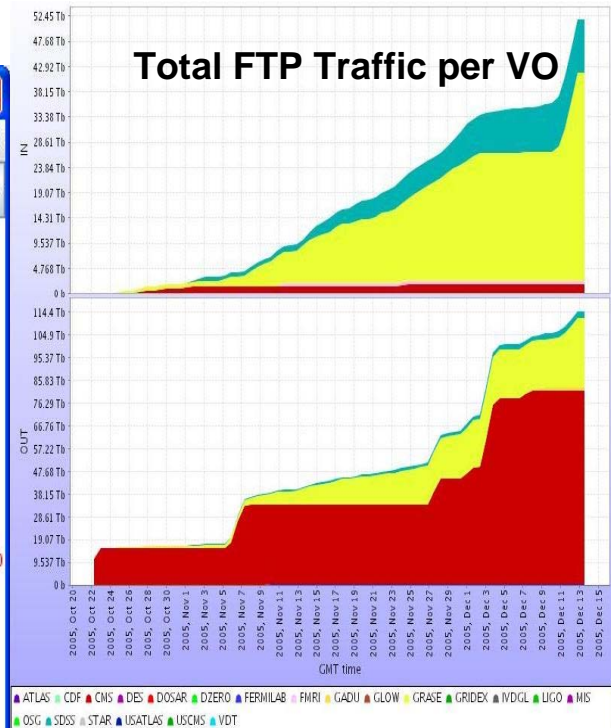
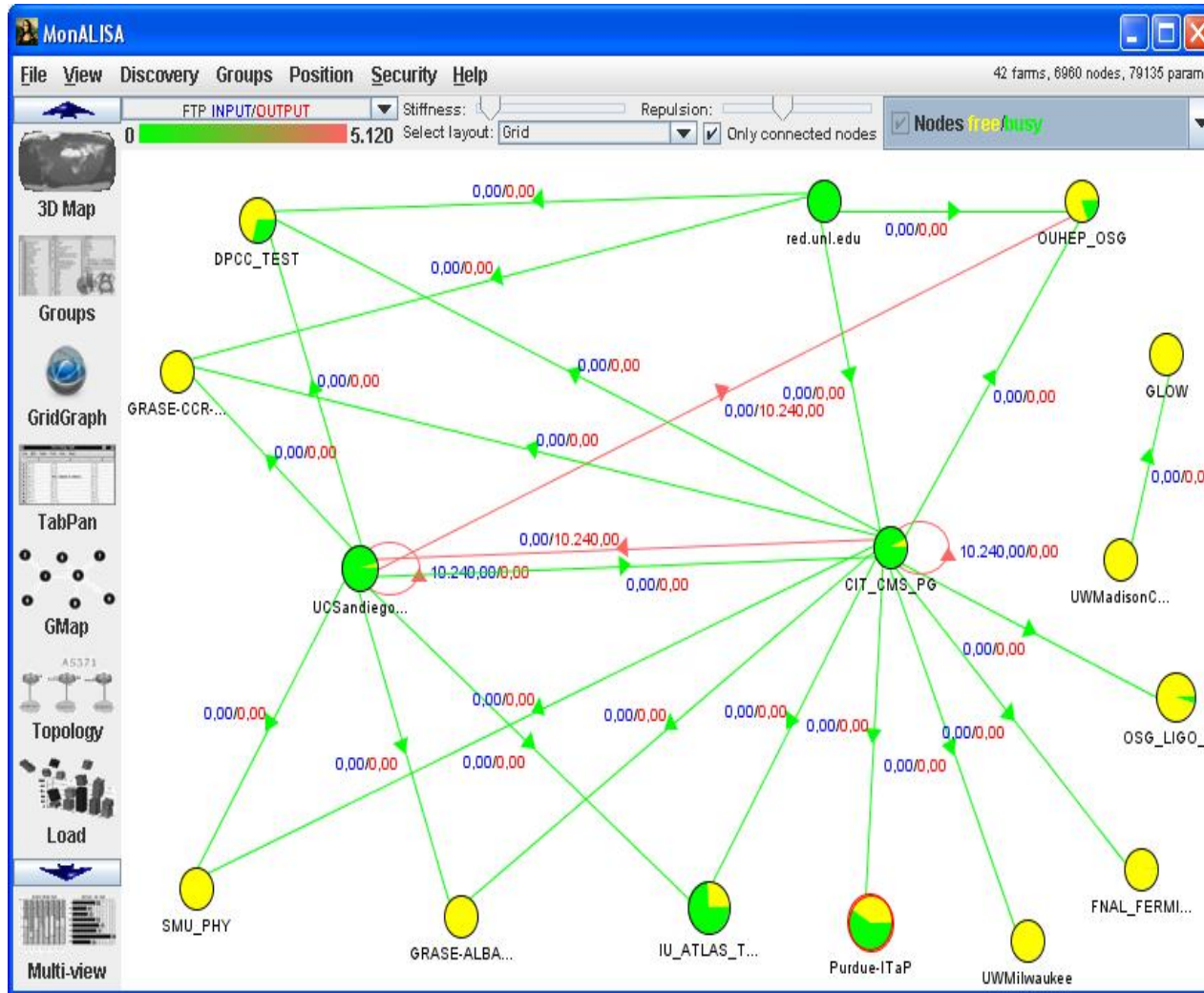


# Monitoring OSG: Resources, Jobs & Accounting





# FTP Data Transfers Among Grid Sites

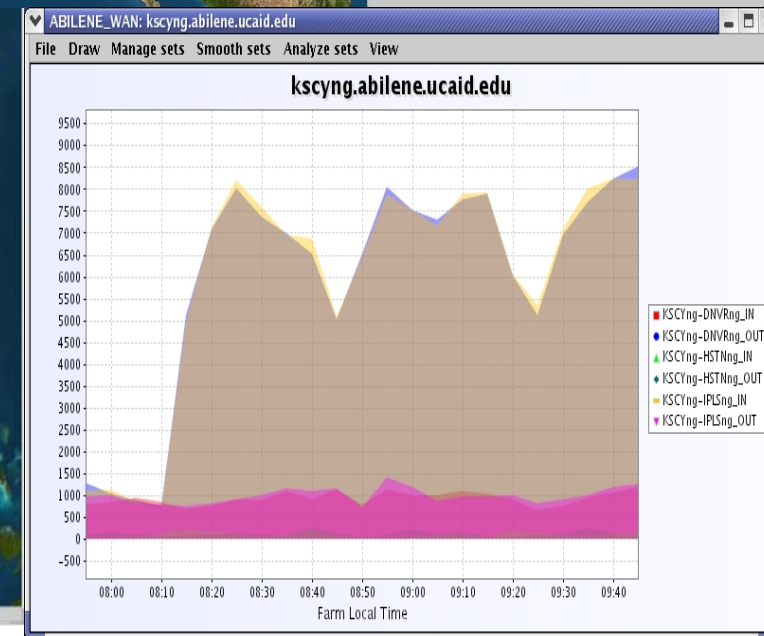
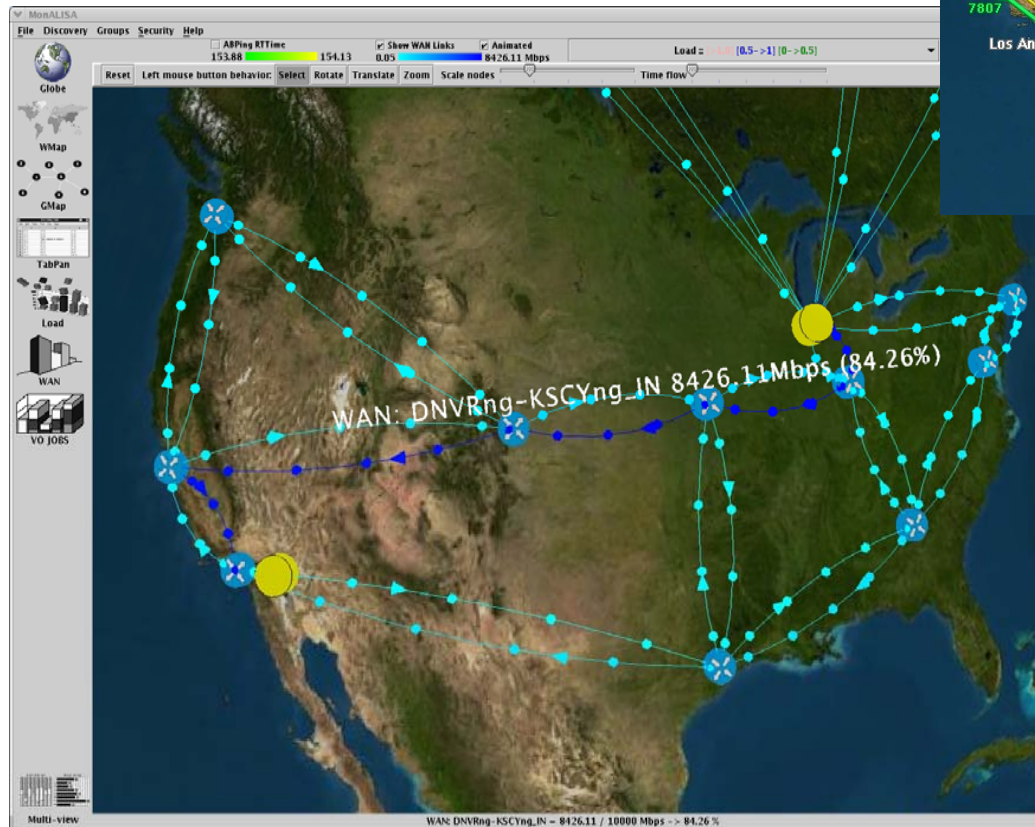




# Monitoring the Abilene Backbone Network



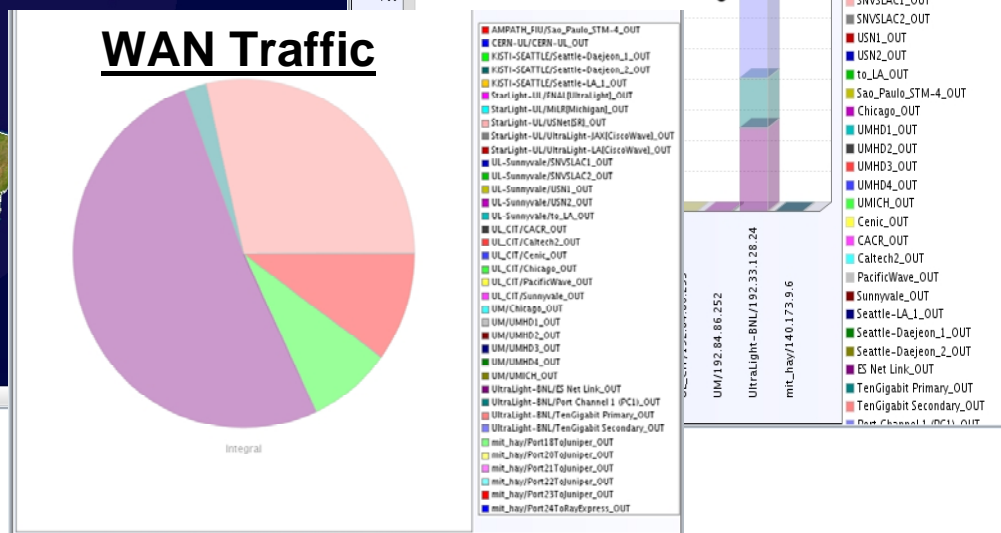
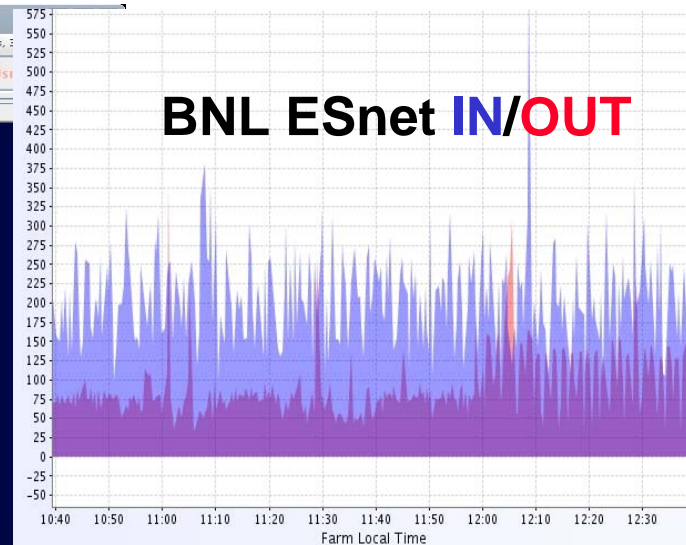
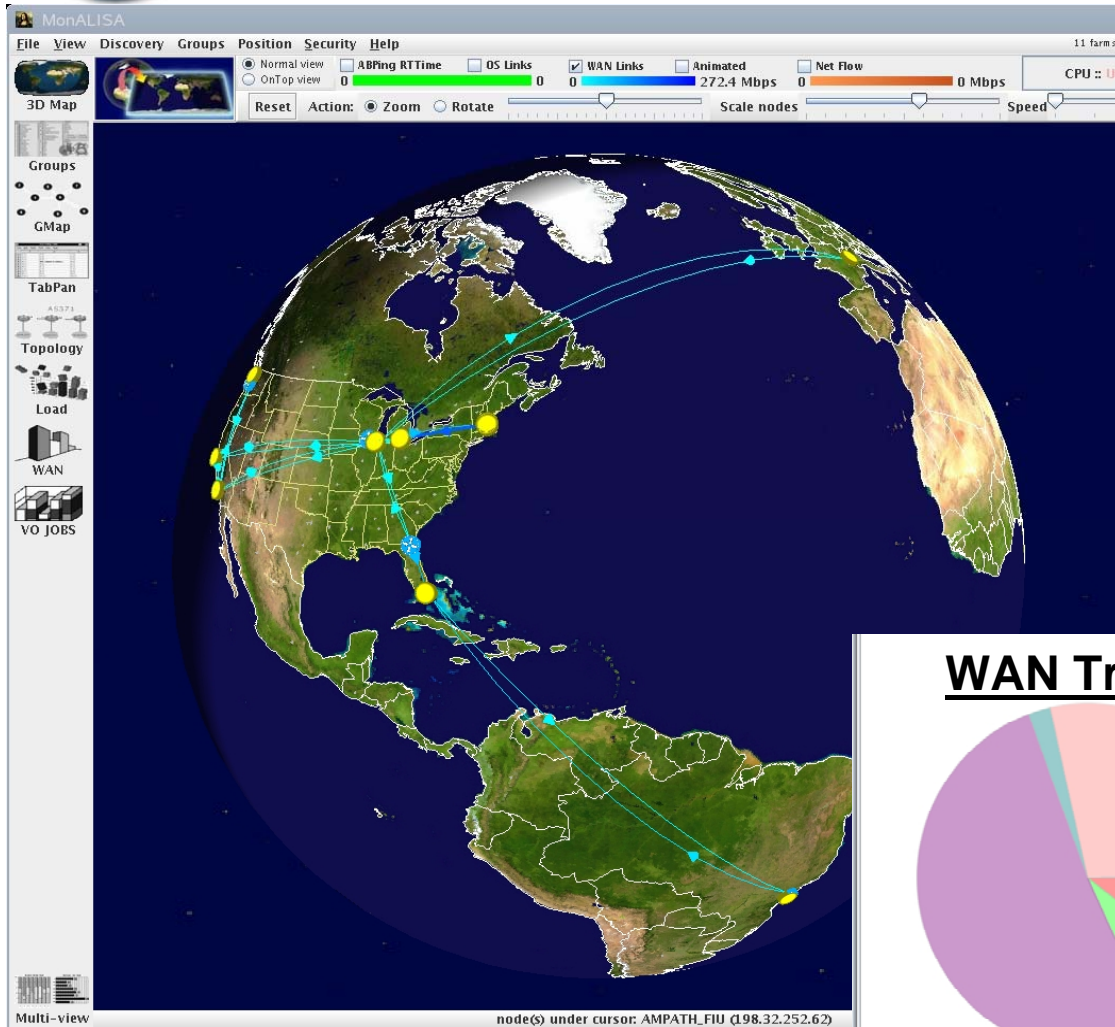
- ◆ Test for a Land Speed Record
- ◆ ~ 7 Gb/s in a single TCP stream from Geneva to Caltech







# The UltraLight Network





# Available Bandwidth Measurements



## Embedded Pathload module.

**MonALISA Repository UltraLight**

**MonALISA Client**  
Click on the button below to start the Monalisa Client.

**Client**

MonALISA Repository

- Global Views
- Node Info
- Available Bandwidth
  - Spider View
  - Matrix View
  - History
  - Sites Status
  - Site Info

close all

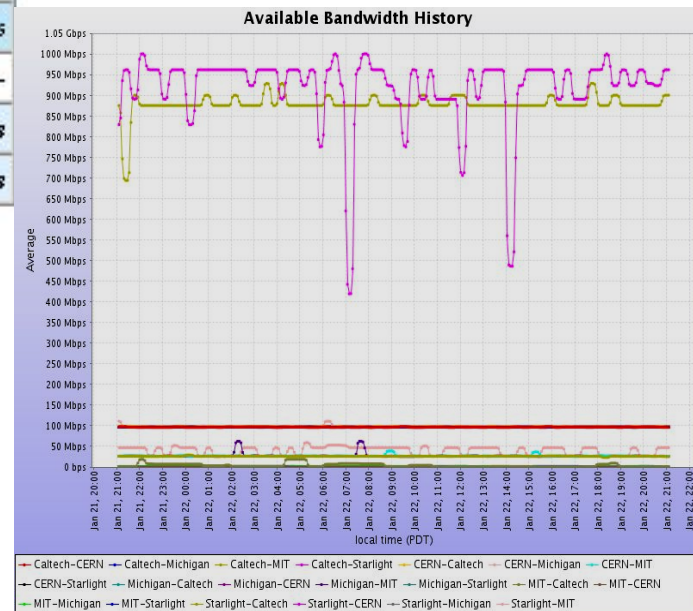
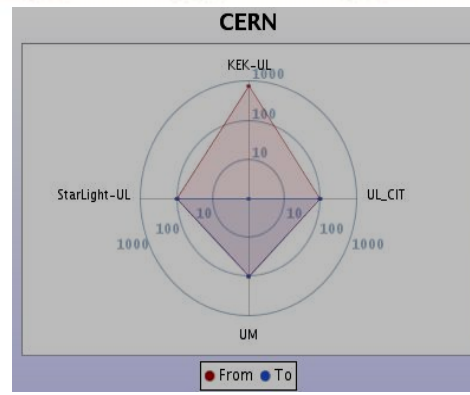
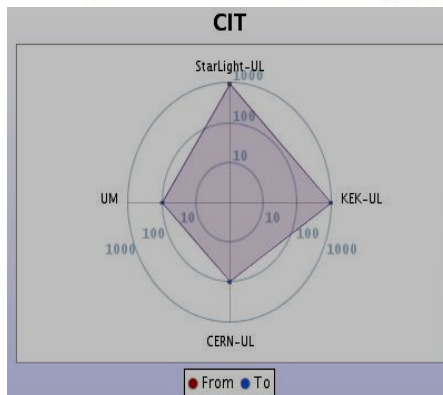
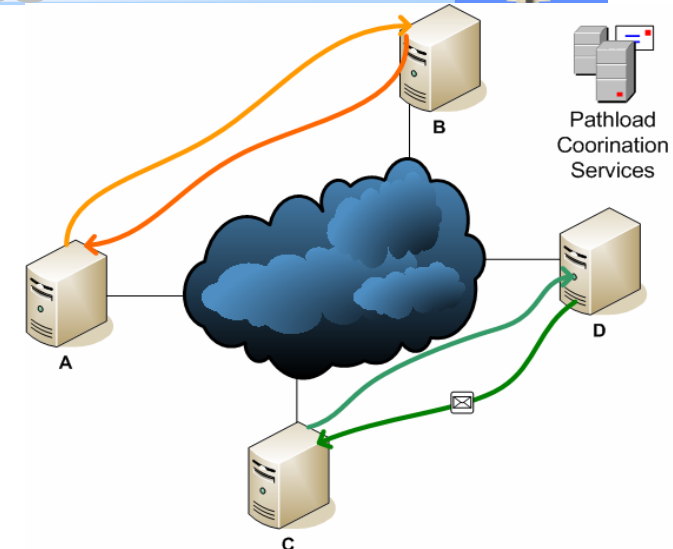
[ABPing Configuration](#)

[Site Administration](#)



## Available bandwidth measurements using pathload

Available bandwidth between UL sites (average)					
Site (from->to)	CERN-UL	KEK-UL	SPRACE-UL	UL_CIT	UM
CERN-UL	-	794.4 Mbps	97.37 Mbps	97.73 Mbps	97.85 Mbps
KEK-UL	750 Mbps	-	96.32 Mbps	993.2 Mbps	96.34 Mbps
StarLight-UL	97.4 Mbps	97.53 Mbps	-	875 Mbps	97.5 Mbps
UL_CIT	97.5 Mbps	993.2 Mbps	876 Mbps	-	96.63 Mbps
UM	97.48 Mbps	96.84 Mbps	96.55 Mbps	96.85 Mbps	-
<b>Min</b>	<b>97.4 Mbps</b>	<b>97.53 Mbps</b>	<b>96.32 Mbps</b>	<b>96.85 Mbps</b>	<b>96.34 Mbps</b>
<b>Max</b>	<b>750 Mbps</b>	<b>993.2 Mbps</b>	<b>876 Mbps</b>	<b>993.2 Mbps</b>	<b>97.85 Mbps</b>





# Coordination Service for Available Bandwidth Measurements



- ◆ Enforces measurement fairness
- ◆ Avoids multiple probes on shared network segments
- ◆ Dynamic configuration of measurement timing
- ◆ Logs events
- ◆ Provides service redundancy by using a master-slave model

The screenshot shows the MonALISA PathloadConfig web interface. The browser address bar displays `http://rb.rogrid.pub.ro:8081/PathloadConfig/`. The page features the MonALISA logo and navigation tabs: HOME, CLIENTS, REPOSITORIES, DOWNLOADS, and LOOKING GLASS. The main content area is titled "Pathload PeerCache Status" and displays the following information:

Current peers count: 3  
Nr. of queued peers: 0  
Token Status: Owner: 141.85.99.143 Src: 141.85.99.143 Dest: 141.85.99.144  
Mon Oct 24 2005 12:14:40 GMT+0300

	FarmName	FarmGroups
1.85.99.144	141.85.99.144	test
1.85.99.143	141.85.99.143	test
99.143	141.85.99.144	
99.133	141.85.99.143	141.85.99.144
	x	x
	-	x
	o	-

Below the table, there is a log of token acquisition and release events:

```
CW9XFqkRCjq12Fwg== from [141.85.99.143/141.85.99.143] to .85.99.144] acquired by [141.85.99.143/141.85.99.143]
iZyLykanj7J1vx36A== from [141.85.99.133/141.85.99.133] to .85.99.144] released.
iZyLykanj7J1vx36A== from [141.85.99.133/141.85.99.133] to .85.99.144] acquired by [141.85.99.133/141.85.99.133]
e2I9j/d/59k01bHw== from [141.85.99.144/141.85.99.144] to
```

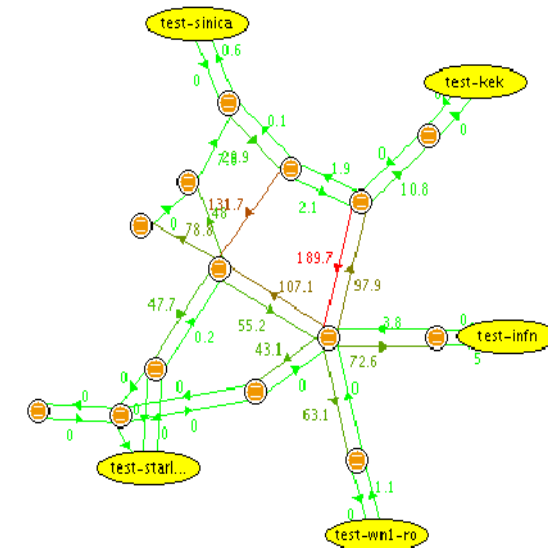
The diagram below the screenshot illustrates the coordination service architecture. It shows four nodes: Sender (A), Receiver (B), and two other nodes (C and D). A central "Coordination Service" cloud is connected to all nodes. A "Token request" is sent from node C to the Coordination Service, which then grants the token to node B. Node B then sends the token to node A. A box indicates that the token is for measuring A->B and B->A.



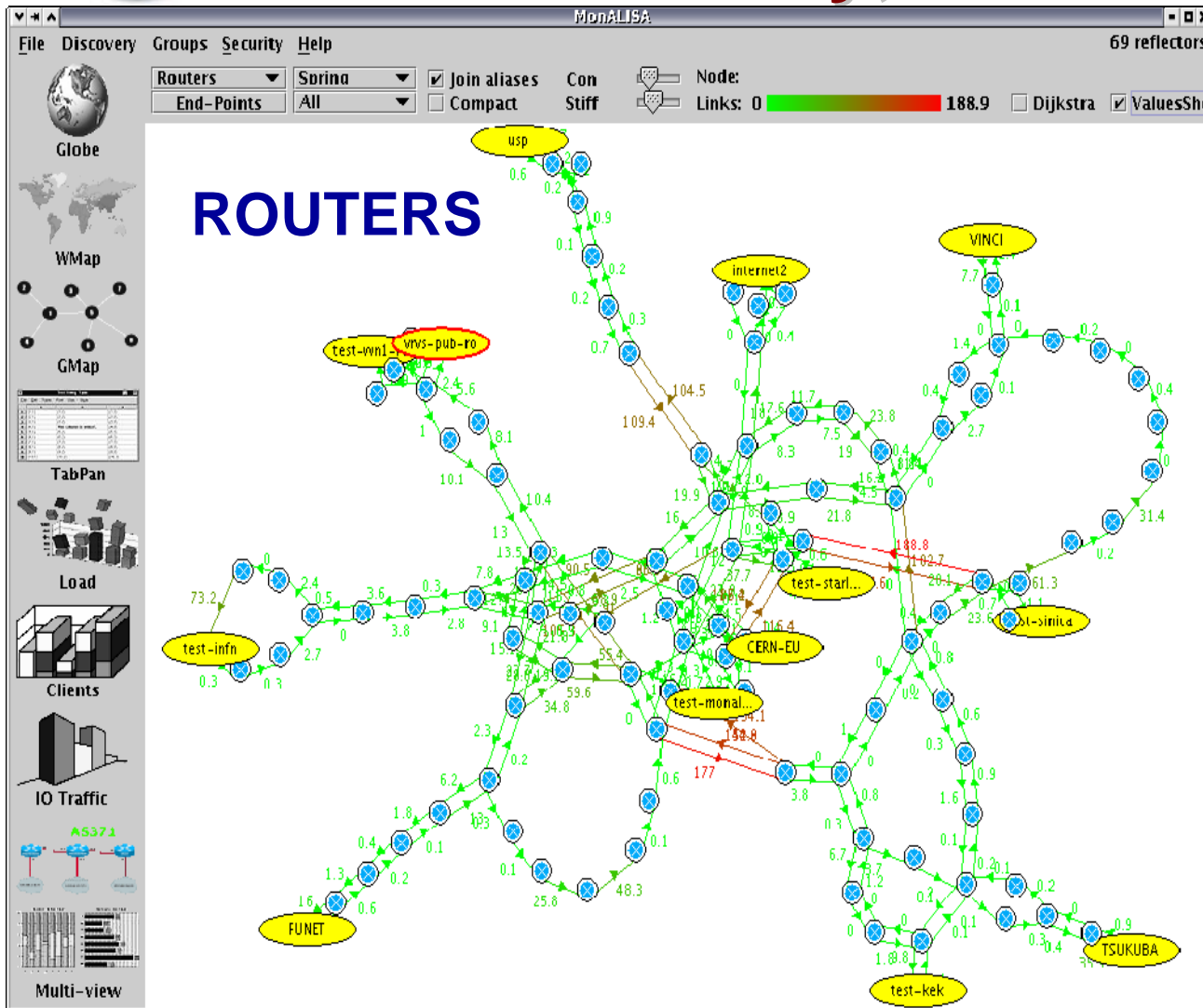
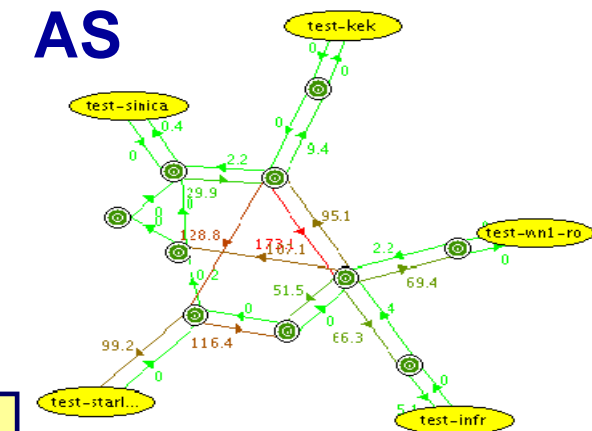
# Monitoring Network Topology, Latency, Routers



## NETWORKS



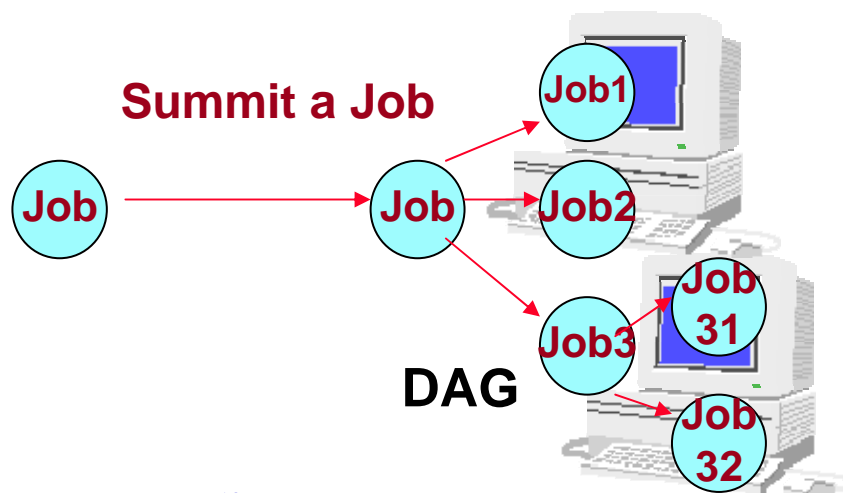
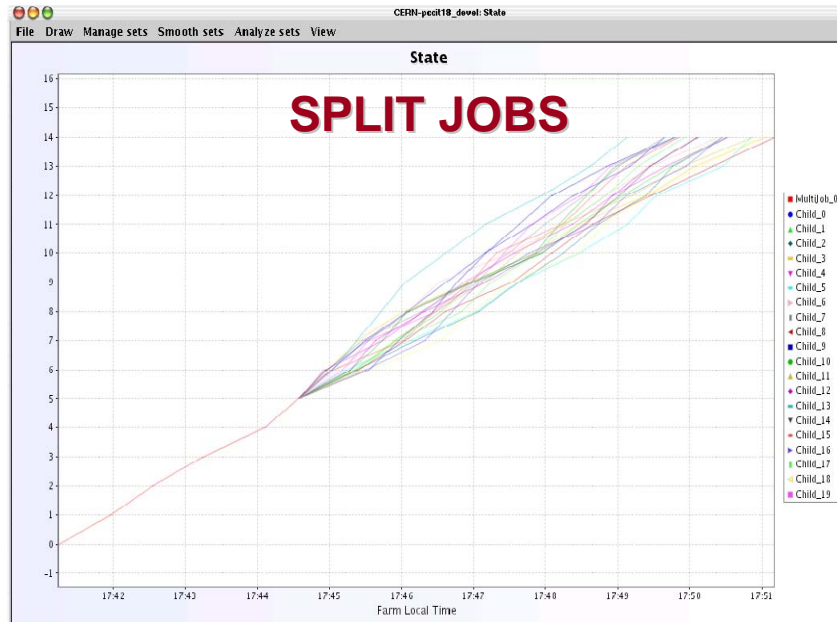
## AS



**Real Time Topology Discovery & Display**

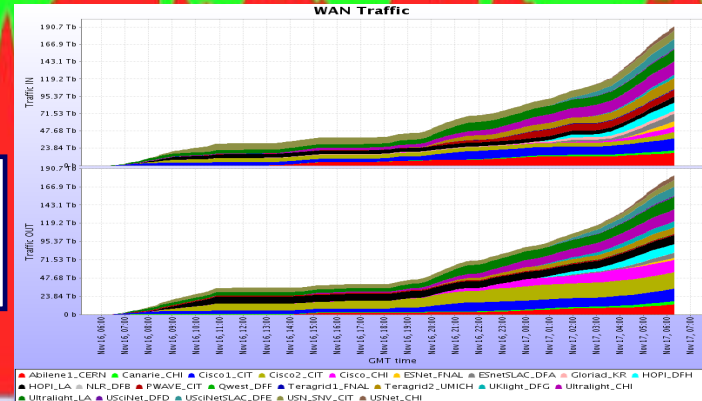
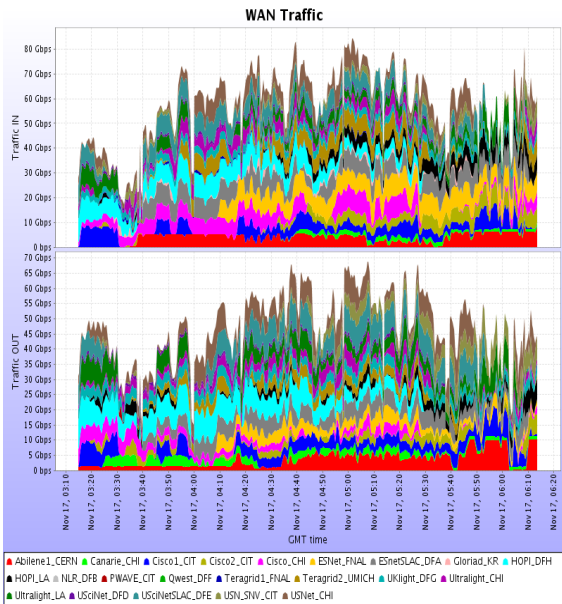
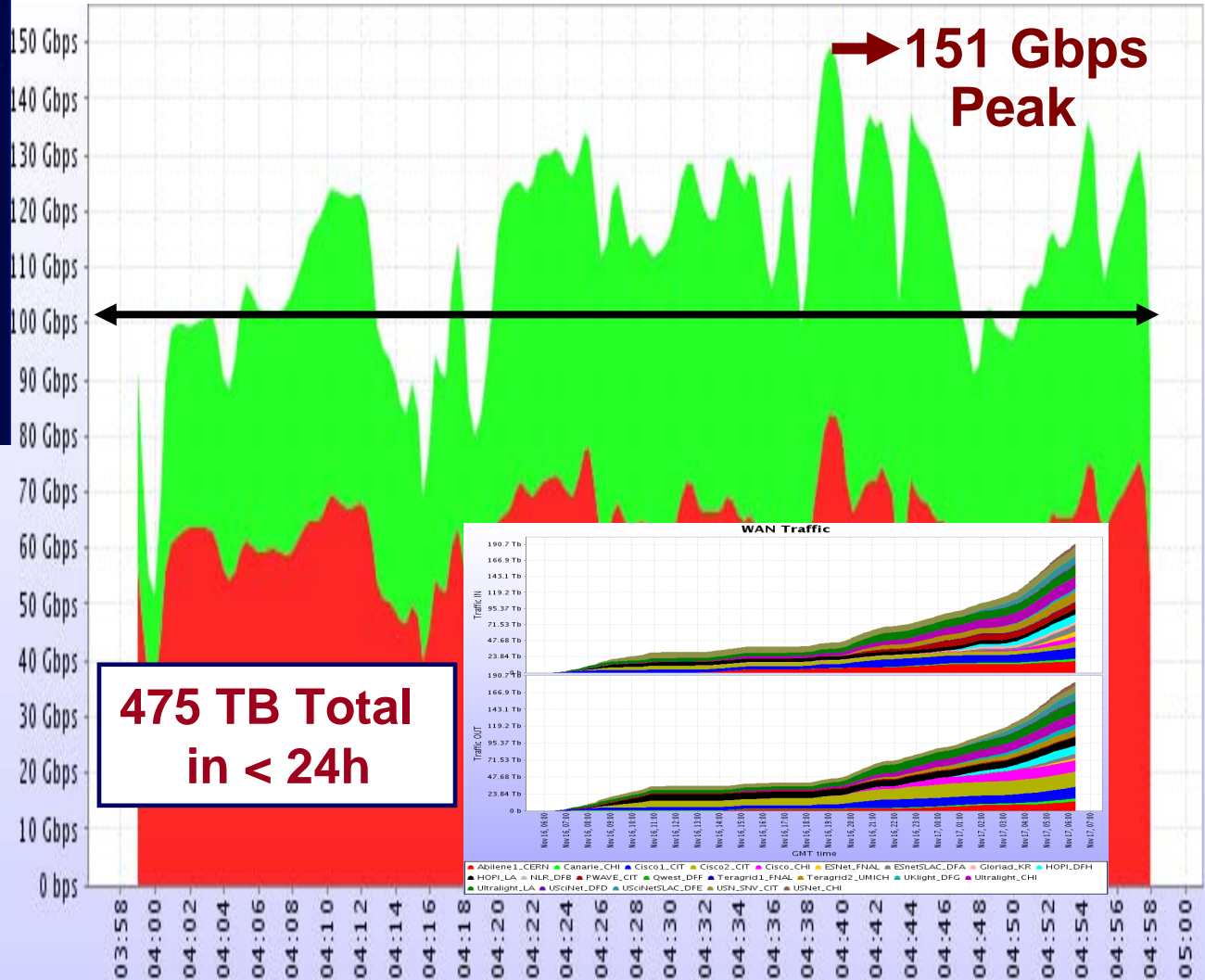


# Monitoring the Execution of Jobs and their Time Evolution



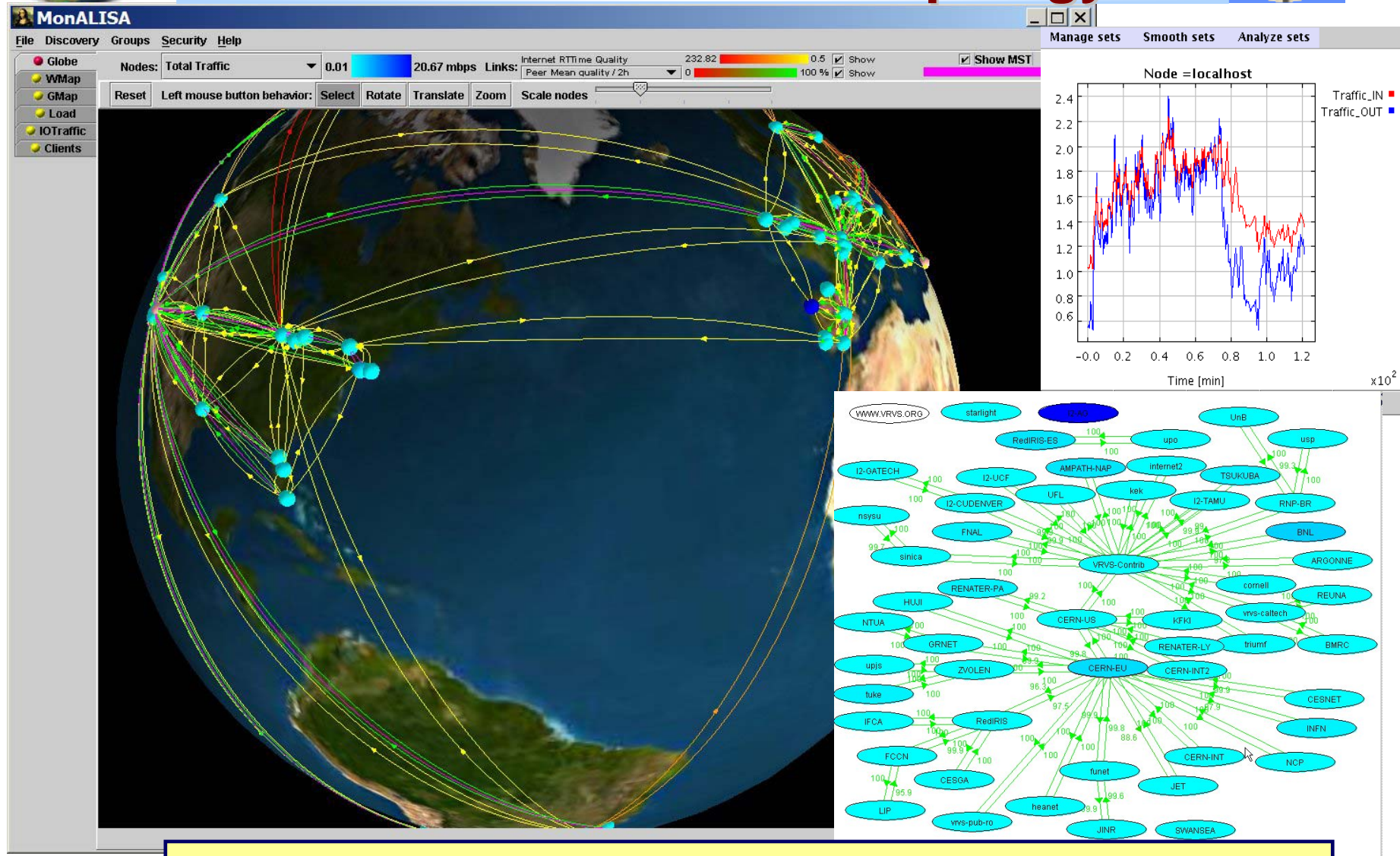


# Bandwidth Challenge at SC2005





# Monitoring VRVS Reflectors and Communication Topology



**Real Time Topology Discovery and Optimization**



# Communities using MonALISA



## Major Communities

- OSG
- CMS
- ALICE
- D0
- STAR
- VRVS
- LGC RUSSIA
- SE Europe GRID
- APAC Grid
- UNAM Grid (Mx)
- ABILENE
  
- ULTRALIGHT
- GLORIAD
- LHC Net
- RoEduNET

## MonALISA Today

**Running 24 X 7  
at 250 Sites**

- Collecting 250,000 parameters in near real-time
- Update rate of 25,000 parameter updates per second
- Monitoring
  - 12,000 computers
  - > 100 WAN Links
- Thousands of Grid jobs running concurrently

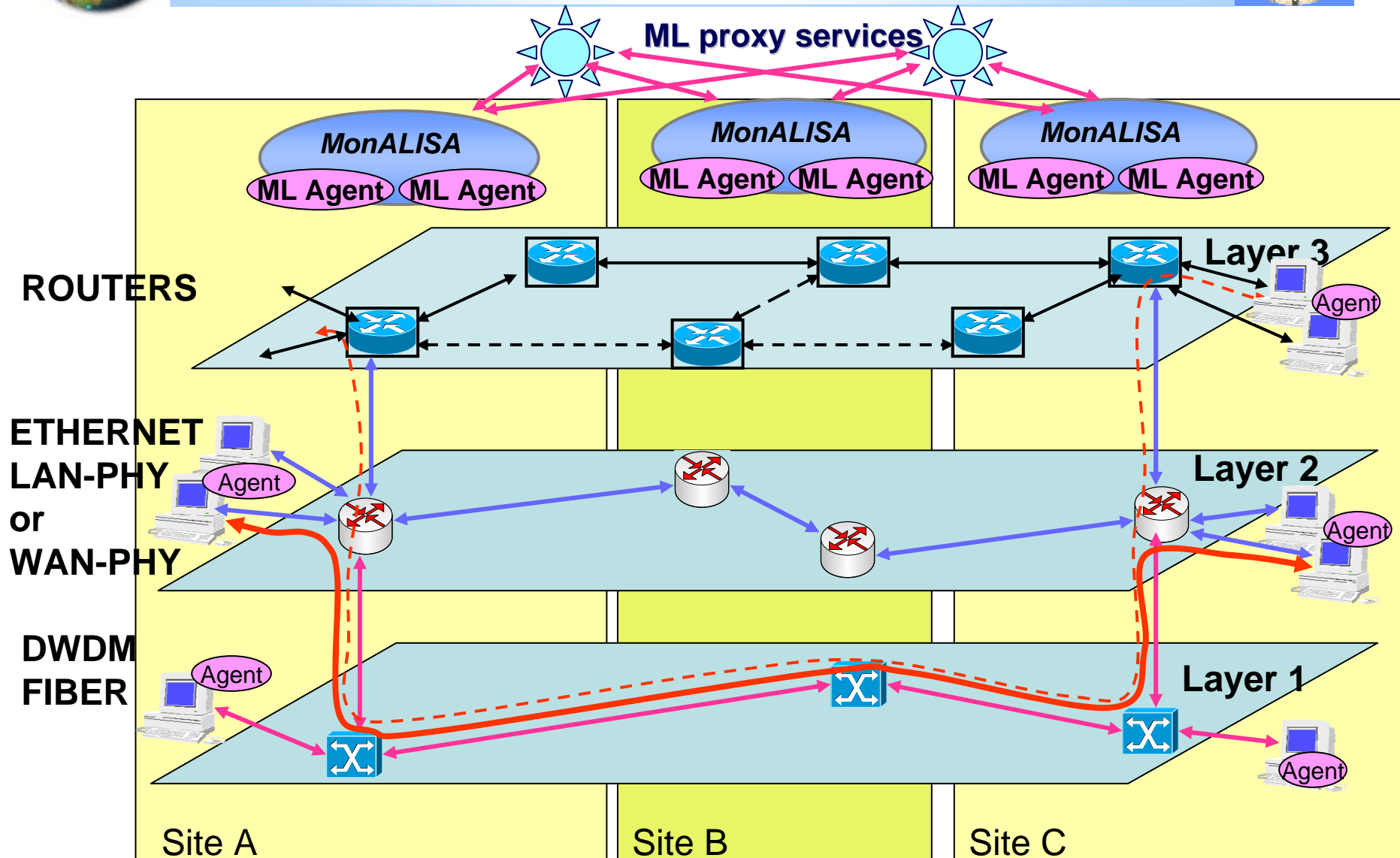
## Demonstrated at:

- ❖ SC2003
- ❖ Telecom World 2003
- ❖ WSIS 2003
- ❖ SC 2004
- ❖ Internet2 2005
- ❖ TERENA 2005
- ❖ IGrid 2005
- ❖ SC 2005



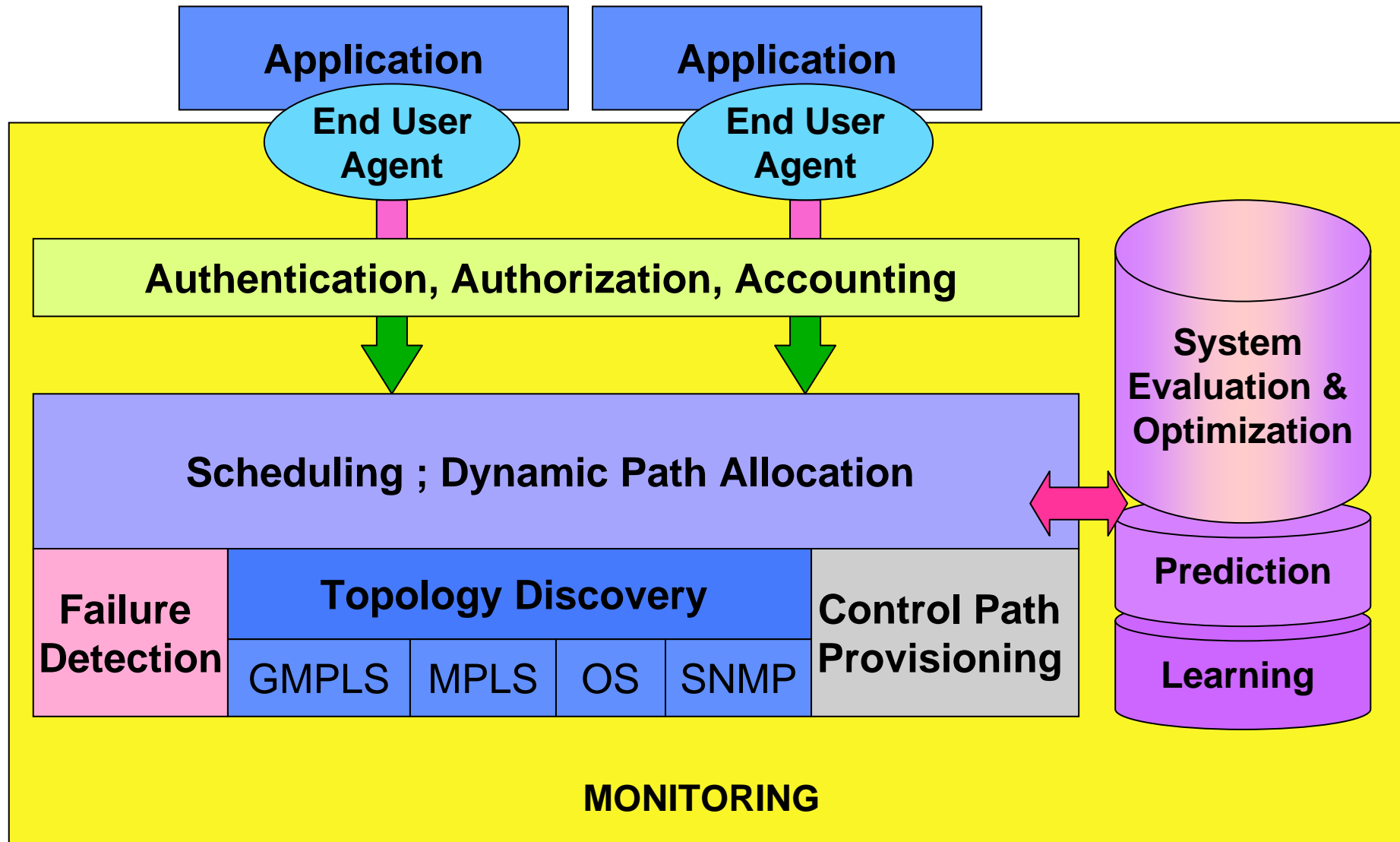


# The Functionality of the VINCI System





# The Main VINCI Services



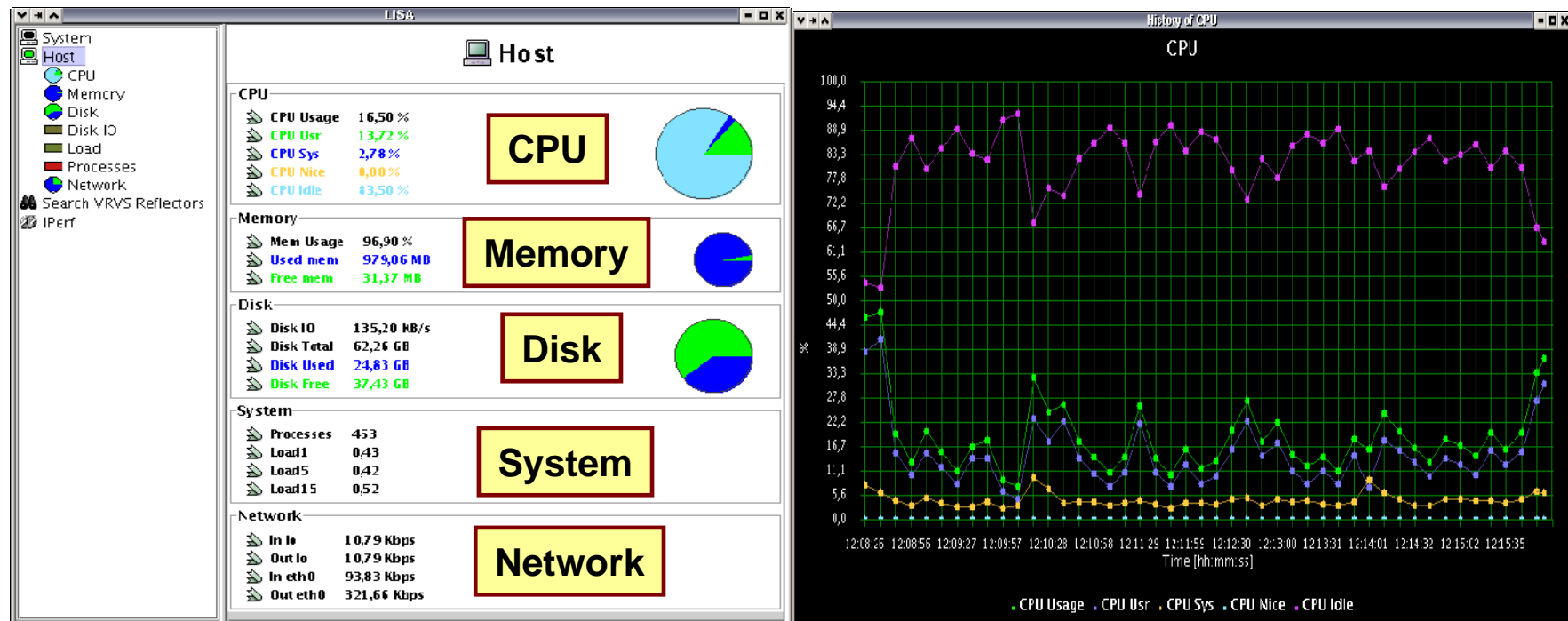


# End User / Client Agent

## LISA- Localhost Information Service Agent



- ❖ Authorization
- ❖ Service discovery
- ❖ Local detection of the hardware and software configuration
- ❖ Complete end-system monitoring: Per-process load, Disk Storage and I/O, per-port network throughputs, etc.
- ❖ End-to-end performance measurements
- ➔ Acts as an active listener for all events related to the requests generated by its local applications





# Secure Service Discovery and AAA



## Service Registration and Discovery

We use JINI Lookup Services to provide a reliable mechanism to dynamically register services, and their dynamic sets of attributes

## Authentication, Authorization and Accounting for Users

We use external AAA services supported by different Virtual Organizations. Loadable plug-in interface modules to support different protocols and services will provide the necessary flexibility to work with different grids and networks



# Topology Discovery Using Specialized Agents



- ◆ **Specialized agents are used to**
  - (1) discover the connection topology for each service**
  - (2) keep a dynamic map of how they are allocated & used, and**
  - (3) get information on the traffic on each segment.**
- ◆ **Agents running on multiple MonALISA services in parallel provide the basic information to the scheduling system**
- ◆ **These agents draw on information from MPLS/GMPLS /DRAGON/Optical Path agents, where the infrastructure provides this functionality**



# Targeted Capabilities for Topology Discovery & Path Selection



## Examples of Capabilities:

- ◆ Determine which path-options exist between two locations in the network
- ◆ List components in the path that are “manageable”
- ◆ Locate network resources and services which have agreements with a given VO
- ◆ Given two replicas of a data source, “discover” (in conjunction with monitoring) the estimated bandwidth and reliability, and hence the “estimated time of successful delivery” of each to a given destination.



# Monitoring and Controlling Optical Planes



File View Discovery Groups Position Security Help 170 farms, 16781 nodes, 218730 params

Local Time: 13:28 (CET) M

**glimmer2**

Input

22	25	28	31	32

Output

22	25	28	31	32

Device parameters  
Device name: GLIMMERGLASS\_glimmer2  
IN Port Labels:  
OUT Port Labels:

Port parameters

	Port ID	Power
Input	25	-8.874 dBm
Output	28	-10.129 dBm

FDX Legend Connect Disconnect

Site info ML Admin App Control OS Admin

File Draw Manage sets Smooth sets Analyze sets View

**glimmer3**

File Draw Manage sets Analyze sets View

**OS\_Ports** Current time: 12/13/05 1:33 PM

Values

Port-Power

**glimmer2**

In

22	25	28	31	32

Out

22	25	28	31	32

**glimmer1**

**Port-Power**

28\_Out  
28\_In  
25\_Out  
25\_In

Farm Local Time

**gva-x3**

**gva-x5**

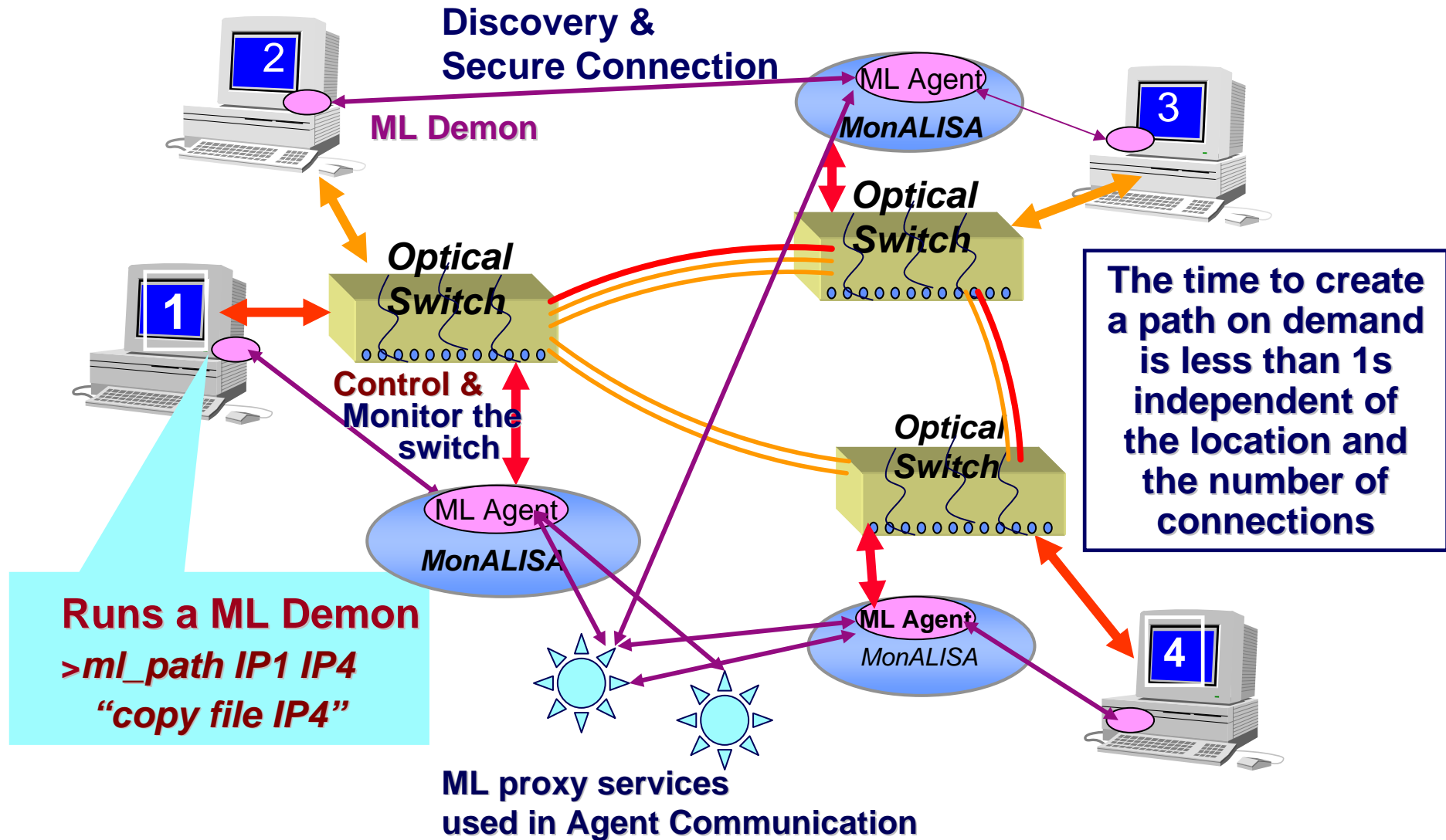
**Controlling**

**Port power monitoring**

**Glimmerglass Switch Example**



# Agents to Create on Demand an Optical Path or Tree







# A Real-World Working Example: Agents Create an Optical Path on Demand



**Dynamic restoration of lightpath if a segment has problems**

The screenshot shows a software interface for network simulation. On the left is a 3D globe map with several colored lines (representing optical paths) connecting various locations across North America. The top menu includes 'File', 'Discovery', 'Groups', 'Position', 'Security', and 'Help'. Below the menu are various controls like 'Normal view', 'OnTop view', 'ABPing RTTme', 'OS Links', and 'Show v'. A toolbar contains 'Reset', 'Left mouse button: Select', 'Rotate', 'Translate', 'Zoom', and 'Sca'. On the right side of the interface, there is a 'Multi-view' panel with icons for 'Groups', 'GMap', 'TabPan', 'Topology', 'Load', 'WAN', 'VO JOBS', and 'OS GMap'. In the center-right, a network topology diagram shows nodes like 'calient1', 'calient2', 'calient3', 'glimmer1', 'glimmer2', and 'glimmer3' connected by lines. A detailed panel for 'glimmer2' is visible, showing input and output port configurations and a 'Sys300Power' graph. The graph shows power consumption for various ports (IN1-P, IN2-P, etc.).



# The Workflow Scheduler



## Scheduler is implemented as a set of collaborating agents

- It provides complete autonomy to each provider of resources, who can implement his own policy
- There is no single point of failure

## “Market Model” Scheduling Scheme

- ❖ Each agent uses policy-based priority queues; it negotiates for an end-to-end connection using a set of cost functions
- ❖ A lease mechanism is implemented for each offer an agent accepts from its peers
- ❖ Two phase commit and periodic lease renewal are used for all agents; this allows a flexible response of the agents to task completion, as well as to application failure or network errors
- ❖ If network errors are detected, supervising agents cause all segments to be released along a path
- ❖ An alternative path then may be set up: rapidly enough to avoid a TCP timeout, so that the transfer can continue uninterrupted



# Interfaces with Network Services

## MPLS, GMPLS



**We are developing agents capable of interacting with MPLS and GMPLS controllers, to provide in near real-time topology maps for other services and to generate connection requests. These agents will continuously monitor and supervise the connections they created.**

- **The MPLS agents can be used together with the optical path agents to create an end-to-end network configuration.**
- **For networks where GMPLS is supported, the agents only need to interface with the head-end devices. The GMPLS standard protocols provides topology discovery (LMP), routing (OSPF) and provisioning (RSVP) and allow interoperability across domains. For example, special routing can be done with VINCI agents if not included in OSPF, but the topology discovery and provisioning can be done with GMPLS.**
- **We also provide agents capable of configuring routers or switches using SNMP or TL1 directly**
- **We will provide agents capable to interacting with other network services systems (for example those in the DRAGON project)**



# Learning and Prediction



- ❖ Learning algorithms (e.g. Self Organizing Neural Networks) will be used to evaluate the traffic created by other applications, to identify major patterns, and dynamically setup effective connectivity maps
- ❖ It is very difficult if not impossible to assume that we could predict all possible events in a complex environment like a grid in advance
- ❖ **Heuristic learning is thus the only practical approach, where agents can acquire the necessary information to describe their environments**
- ❖ The multi-agent learning task includes two levels:
  - ➔ the local level of individual learning agents
  - ➔ the global level, exploiting inter-agent communication
- ❖ We need to ensure that each agent can learn to optimize its actions locally, while the global monitoring mechanism acts as a 'driving force' that causes the agents' behavior to evolve collectively, based on the accumulated experience



# Mumbai-Japan-US Links

