

Logistical Networking: Developments and Deployment

Micah Beck, Assoc. Prof. & Director

*Logistical Computing &
Internetworking (LoCI) Lab*

Computer Science Department

University of Tennessee

mbeck@cs.utk.edu

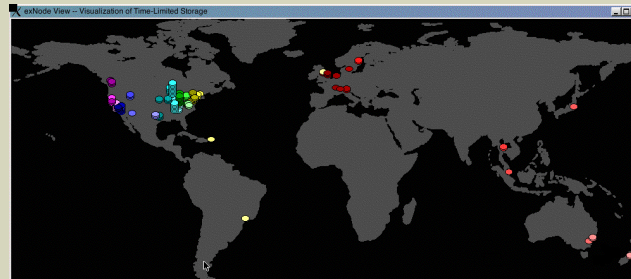
AMPATH Astronomy WG, Miami

Jan 31, 2003



LoCI

LOGISTICAL COMPUTING AND
INTERNETWORKING LAB



UNIVERSITY OF TENNESSEE

Logistical Networking Research

» University of Tennessee

- Micah Beck
- James S. Plank
- Jack Dongarra

» University of California, Santa Barbara

- Rich Wolski

» Funding

- Dept. of Energy SciDAC
- National Science Foundation ANIR
- UT Center for Info Technology Research



What is Logistical Networking

- » A scalable mechanism for deploying shared storage resources throughout the network
- » An general store-and-forward overlay networking infrastructure
- » A way to break long transfers into segments and employ heterogeneous network technologies
- » P2P storage and content delivery that doesn't using endpoint storage or bandwidth



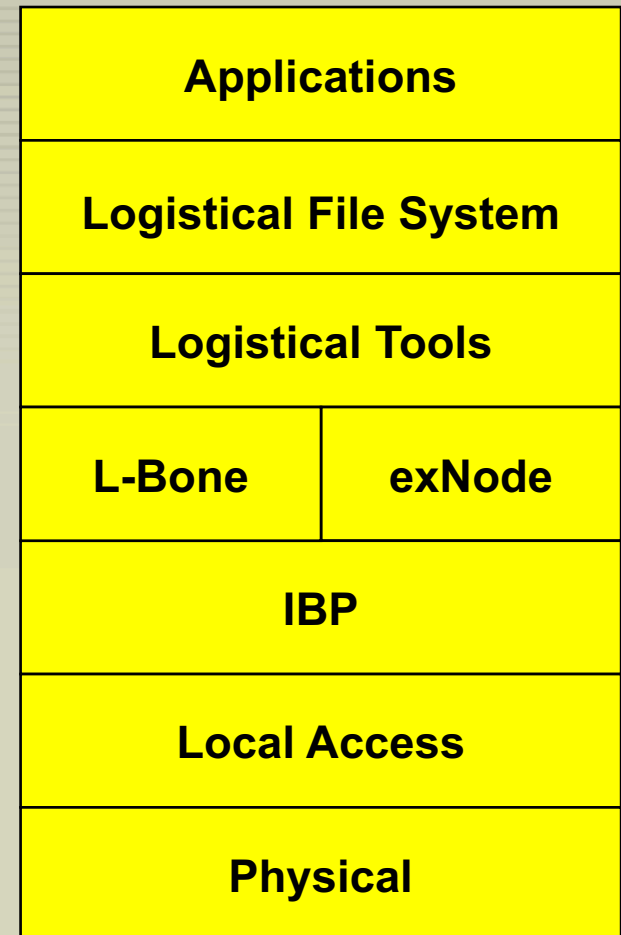
Why “Logistical Networking”

- » Analogy to logistics in distribution of industrial and military personnel & materiel
- » Fast highways alone are not enough
 - Goods are also stored in warehouses for transfer or local distribution
- » Fast networks alone are not enough
 - Data must be stored in buffers/files for transfer or local distribution
- » Conventional vs logistical networking
 - Datagram routers make *spatial* choices
 - Storage depots enable *temporal* choices



The Network Storage Stack

- Our adaption of the network stack architecture for storage
- Like the IP Stack
- Each level encapsulates details from the lower levels, while still exposing details to higher levels



IBP: The Internet Backplane Protocol

- » Storage provisioned on community “depots”
- » Very primitive service (similar to block service, but more sharable)
 - Goal is to be a common platform (exposed)
 - Also part of end-to-end design
- » Best effort service – no heroic measures
 - Availability, reliability, security, performance
- » Allocations are time-limited!
 - Leases are respected, can be renewed
 - Permanent storage is too strong to share!



Models of Sharing: Logistical Networking

- » Moderately valuable resources
 - Storage, server cycles
- » Sharing enabled by relative plenty
- » Internet-like policies
 - Loose access control
 - No per-use accounting
- » Primary design goal: scalability
 - Application autonomy
 - Resource transparency
- » Burdens of scalability
 - The End-to-End Principles
 - Weak operation semantics
 - Vulnerability to Denial of Service



The Network Storage Stack

LoRS: The Logistical Runtime System:
Aggregation tools and methodologies

The L-bone:
Resource Discovery
& Proximity queries

The exNode:
A data structure
for aggregation

IBP: Allocating and managing network
storage (like a network malloc)

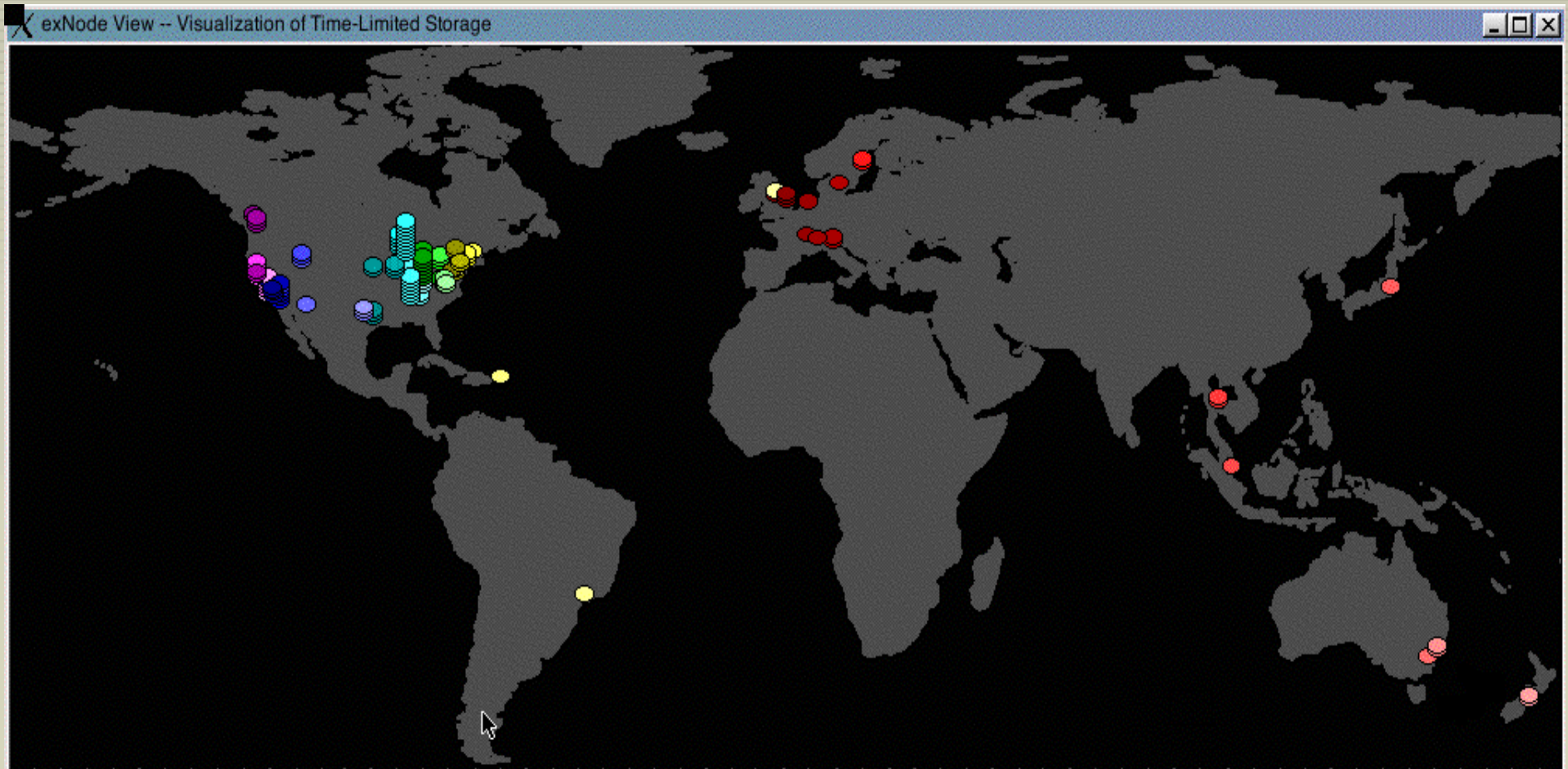


The Logistical Backbone (L-Bone)

- » LDAP-based storage resource discovery.
- » Query by capacity, network proximity, geographical proximity, stability, etc.
- » Periodic monitoring of depots.
- » 10 Terabytes of shared storage. (with plans to scale to a petabyte...)



L-Bone: January 2003



LDCI



IBP Deployment

- » Logistical Backbone
 - 147 depots in 15 countries
 - 10TB of shared storage
- » Leverages Planet Lab nodes (Intel Research Labs)
- » Depots/collaborations in AMPATH region
 - Puerto Rico (Guy Cormier, Univ. of Puerto Rico)
 - Brazil (Univ. of Sao Paolo)
- » AMPATH Chicago/FIU 1 GB link test
 - 75 Mb/s to a depot attached at 100Mb/s



The Network Storage Stack

LoRS: The Logistical Runtime System:
Aggregation tools and methodologies

The L-bone:
Resource Discovery
& Proximity queries

The exNode:
A data structure
for aggregation

IBP: Allocating and managing network
storage (like a network malloc)



Logistical Runtime System

- » Basic Primitives:
 - Upload, Download, Augment, Refresh
- » End-to-end Services
 - Checksums, Encryption, Compression
- » Other Things We Can Do
 - Routing through an intermediate depot to reduce IP RTT, speeding up TCP transfers
 - Overlay multicast using either multiple TCP streams or IP multicast at tree nodes



Upload

Terminal Window:

```
login
map.cs.utk.edu
Iboneserver == dsj.sinrg.cs.utk.edu
segs == zzz
copies == 1
frags == 8
windemo == 0
offline == 0
threads == 8
xnd_launch xnd_mupload "/Users/atchley/Movies/o-brother.mpg"
-lbone-host dsj.sinrg.cs.utk.edu
entering new_exnode standard-output
ibanez.cs.utk.edu:6714
ibanez.cs.utk.edu:6714
galapagos.cs.utk.edu:6714
taylor.cs.utk.edu:6714
galapagos.cs.utk.edu:6714
ovation.cs.utk.edu:6714
taylor.cs.utk.edu:6714
ovation.cs.utk.edu:6714
627 hang zhao CFP
628 Linzhen Xuan Ban
629 Terry Moore Re:
630 <ekl@mail.inter.Re:
631 YING DING nev
632 Micah Beck Nev
```

exNode Command Window:

Mode: **upload**

Select a file to Upload:

Save the exNode as:

Choose location:

Duration (days):

Copies:

Fragments:

Max Depots:

exNode View – Visualization of Time-Limited Storage

United_States
Europe

Other
UPR

UCSB
UCSD
TAMU
WISC
IUPUI
UNC
Harvard
CUA
Surfmnet
Stuttgart
ENS
UNIPMN

standard-output
126638084 bytes

Uploading

Augment

XDarwin File Edit Window Help Mon 09:33 AM

```
cache ==
cmd == augment
location == hostname= ucsb.edu
inputfile == /Users/atchley/Movies/o-brother.mpg.xnd
blocksize == 256 K
maxdepot == 12
duration == 1
outputfile == /Users/atchley/Movies/o-brother.mpg.xnd
prebuf == 1
lboneserver == dsj.sinrg.cs.utk.edu
segs == zzz
copies == 1
frags == 10
windemo == 0
offline == 0
threads == 10
xnd_launch xnd_maugment "/Users/atchley/Movies/o-brother.mpg.xnd" -dir /Users/atchley/Movies -F 10 -l thread 10 -p
```

exNode View – Visualization of Time-Limited Storage

United States Europe

exNode Command

Mode: **augment**

Select an exNode file to Augment:

Save new exNode to:

Choose location:

Duration (days):

Copies:

Fragments:

Max Depots:

Augmenting

(0) Aug 20	██████████
(1) Aug 20	██████████
(2) Aug 20	██████████
(3) Aug 20	██████████
(4) Aug 20	██████████
(5) Aug 20	██████████
(6) Aug 20	██████████
(7) Aug 20	██████████

Download

The screenshot shows a Mac OS X desktop with several windows. The top window is a terminal displaying a file download progress log. The middle window is a map titled "exNode View - Visualization of Time-Limited Storage" showing various nodes across the United States and Europe. The bottom window is the "exNode Command" interface, which includes a "Mode" dropdown set to "download", a file path, a destination, and a progress bar. The progress bar shows the file is being downloaded in 26 chunks, with the first 10 chunks completed and the remaining 16 in progress.

```
DIR
OFFSET 0
SIZE 0
BLOCKSIZE 5120 K
CACHE_SIZE 1
PREBUFFER_SIZE 1
free ing..
Downloading 1,0 5242880 - 5242880 from rave
display_rtdownload downloadarrow_1,0
Downloading 2,0 10485760 - 3585142 from i2-
display_rtdownload downloadarrow_2,0
Downloading 0,0 0 - 5242880 from dsj2.uits.
display_rtdownload downloadarrow_0,0
Downloading 3,0 15728640 - 5242880 from cha
display_rtdownload downloadarrow_3,0
Downloading 4,0 20971520 - 5242880 from tau
display_rtdownload downloadarrow_4,0
Downloading 5,0 26214400 - 1927404 from ova
display_rtdownload downloadarrow_5,0
```

628 Linzhen Xuan Bar
629 Terry Moore Re:
630 <ekl@mail.inter Re:
631 YING DING nev
632 Micah Beck Nev
633 Debashis Talukr SCC

Mode: **download**

Select an exNode to download:
/Users/atchley/Movies/o- **Browse..**

Download the file to:
/dev/null **Browse..**

Threads: **6**
Prebuffer: **1**
Blocksize: **5120 K**

Kill
List
Run

Downloading
/Users/atchley/Movies/o-brother.mpg.xnd
/Users/atchley/Movies/o-brother.mpg
126638084 bytes

(0) Aug 21	(1) Aug 21	
(2) Aug 21	(3) Aug 21	(5) Aug 21
(4) Aug 20	(6) Aug 21	
(7) Aug 20	(8) Aug 21	(9) Aug 21
(10) Aug 21	(11) Aug 21	
(12) Aug 21	(13) Aug 21	(14) Aug 21
(15) Aug 21	(16) Aug 21	(18) Aug 20
(17) Aug 21	(19) Aug 21	(20) Aug 21
(22) Aug 21	(21) Aug 21	(23) Aug 21
(26) Aug 21	(25) Aug 21	(24) Aug 20

Clear Slow Redraw Quick Redraw Draw Arrows Active Connection Quit

IBP Enables Data Intensive Collaboration

- » Large files can be uploaded to nearby depots, then managed by movement between depots
 - End systems are not involved in long distance transfers
- » Data can be moved near to distant collaborator without being downloaded into their end system
 - Direct access to collaborators private storage is not required
- » Depot-to-depot transfers can take advantage of multithreading, UDP transfer, Net/Web 100, other high-performance optimizations



Example Application: IBPvo

- » Web interface allows television shows to be recorded in U.S., uploaded to IBP depots
- » Resulting AVI files are O(1GB) in size
- » ExNode is delivered to user by mail
- » Multithreaded transfer to APAN region depots
- » Users watch programs by downloading to their own workstations, viewing locally
- » A reciprocal service would allow users in U.S. direct access to AMPATH region television
- » <http://promise.sinrg.cs.utk.edu/votest>



Other Areas of Application

- » Management of massive data sets
 - Produced by simulation
 - Captured from experimentation
 - Generated by sensors and instruments
- » Caching and staging of data in high-performance wide area (e.g. Grid) computation
- » Content Distribution of highly popular content
- » Overlay routing, multicast
- » Digital Libraries
- » Checkpoints and backups
- » Wide area file systems



LoCI Lab Online

<http://loci.cs.utk.edu>

- » IBP server and clients for Unix/Linux/OS X
 - Additional clients for Java, Win32
- » Logistical Runtime System libraries and tools
 - Run under Unix/Linux/OS X natively
 - Ported to Windows under Cygwin
 - Includes visualization (Tcl/tk)
 - Web interface
- » Logistical Backbone resource discovery server
 - Unix/Linux/OS X only
- » Publications, documentation, L-Bone status



LoCI

