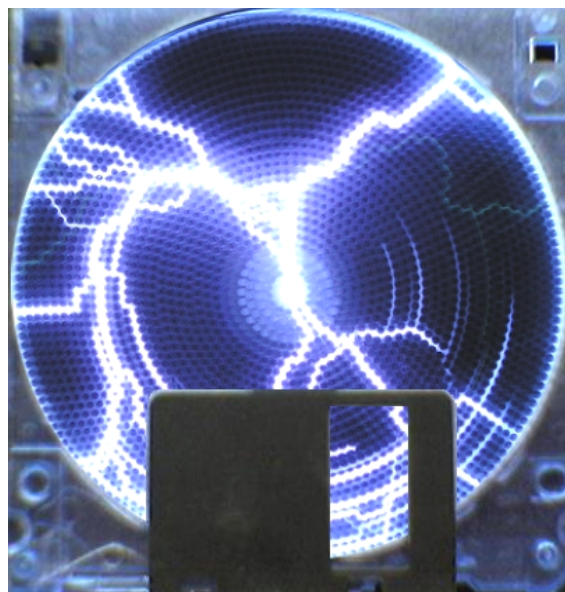


Indiana University's Lustre WAN: Empowering Production Workflows on the TeraGrid



Stephen C. Simms
Manager, Data Capacitor Project
TeraGrid Site Lead, Indiana University
ssimms@indiana.edu

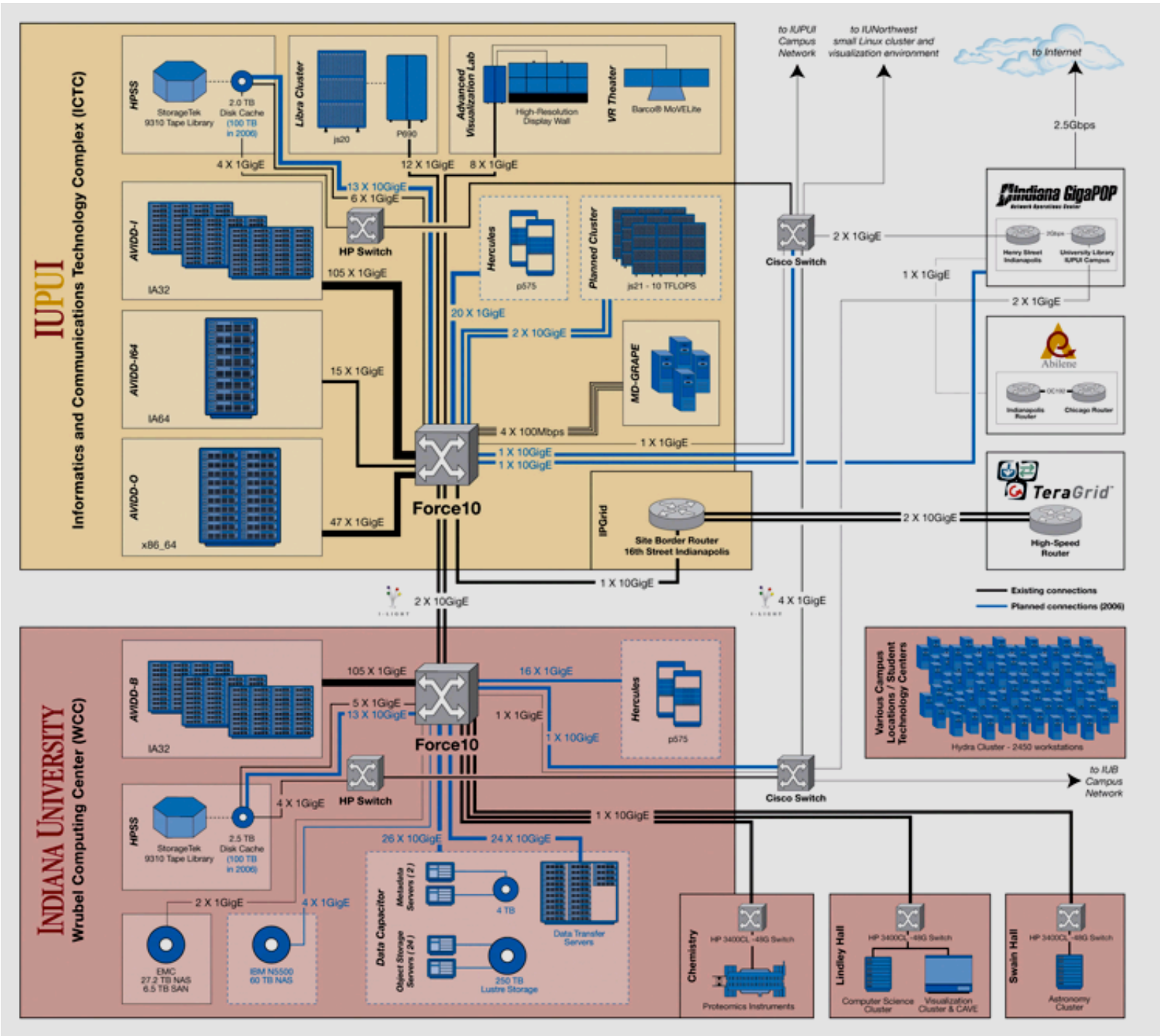
Got Lustre WAN?

The Data Capacitor Project

NSF Funded in 2005
535 Terabytes Lustre storage
14.5 GB/s aggregate write
Short term storage



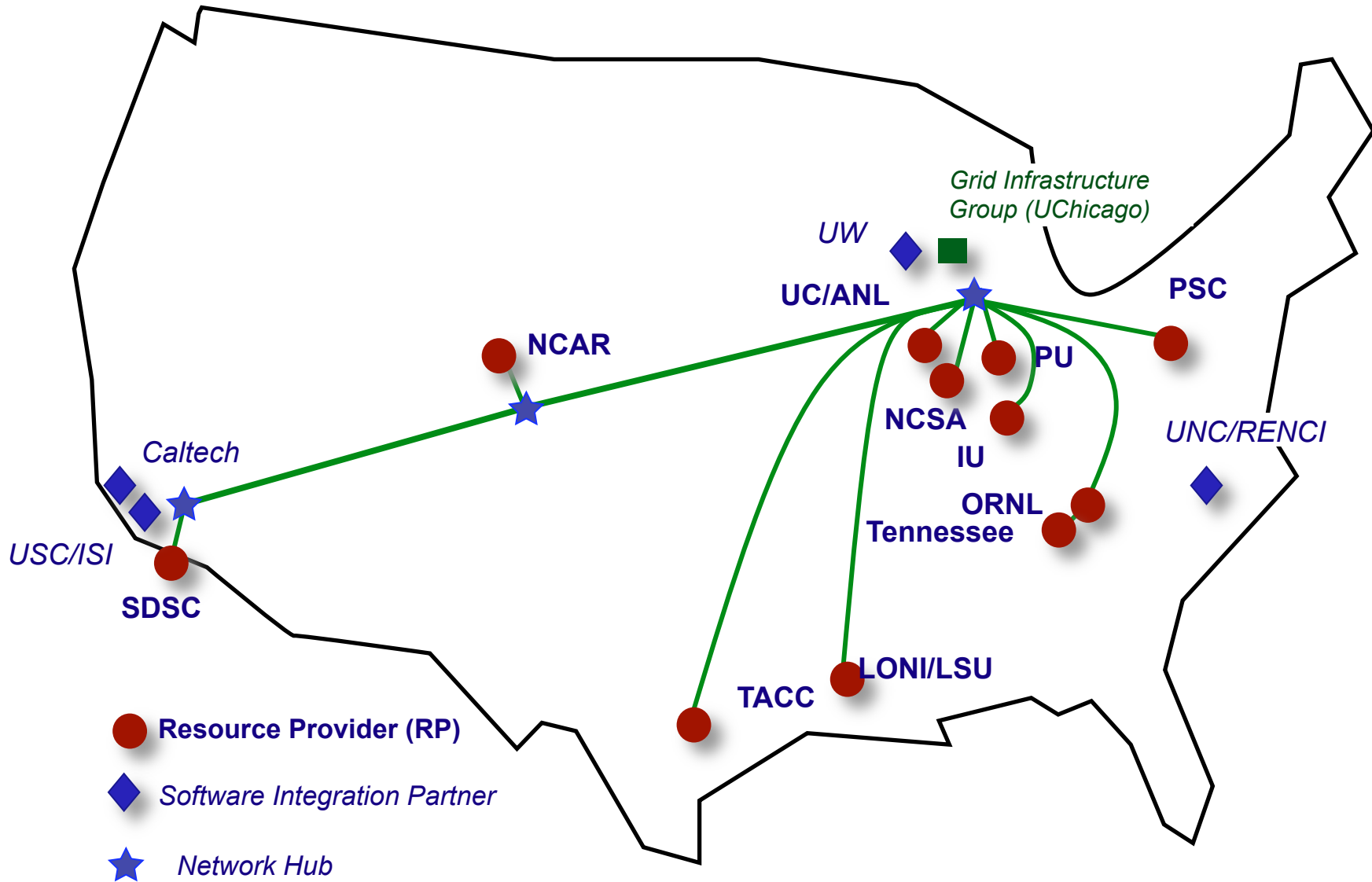
<http://www.flickr.com/photos/shadowstorm/404158384/>
<http://www.flickr.com/photos/dvd5/163647219/>
<http://www.flickr.com/photos/vidiot/431357888/>



The TeraGrid

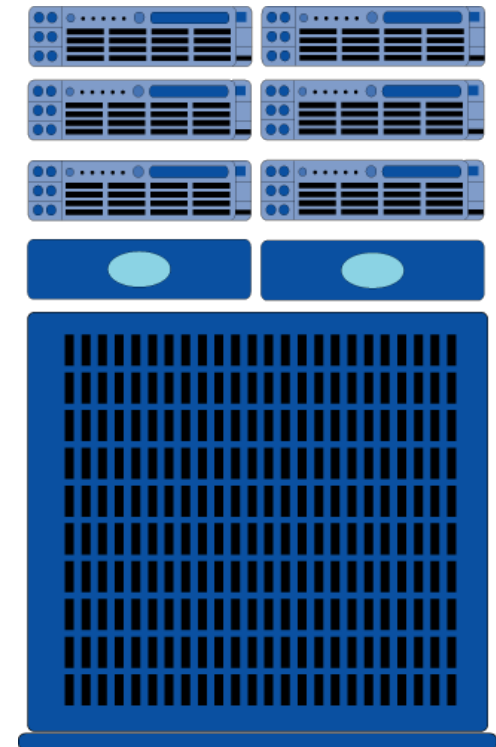
- TeraGrid is an open scientific discovery infrastructure combining leadership class resources at eleven partner sites to create an integrated, persistent computational resource.
- ANL, IU, Louisiana Optical Network Initiative (LONI), NCSA, NICS, ORNL, PSC, Purdue, SDSC, TACC, NCAR

The TeraGrid Map



IU's Data Capacitor WAN

- 1 pair Dell PowerEdge 2950 for MDS
- 2 pair Dell PowerEdge 2950 for OSS
 - 2 x 3.0 GHz Dual Core Xeon
 - Myrinet 10G Ethernet
 - Dual port Qlogic 2432 HBA (4 x FC)
 - 2.6 Kernel (RHEL 5)
- DDN S2A9550 Controller
 - Over 2.4 GB/sec measured throughput
 - 360 Terabytes of spinning SATA disk
- Currently running Lustre 1.6.7.2
 - Upgrading to 1.8.1.1 in May
- Announced production at LUG 2008
 - Allocated on Project by Project basis



Networking



IU UID Mapping

Lightweight

- Not everyone needs / wants kerberos

- Not everyone needs / wants encryption

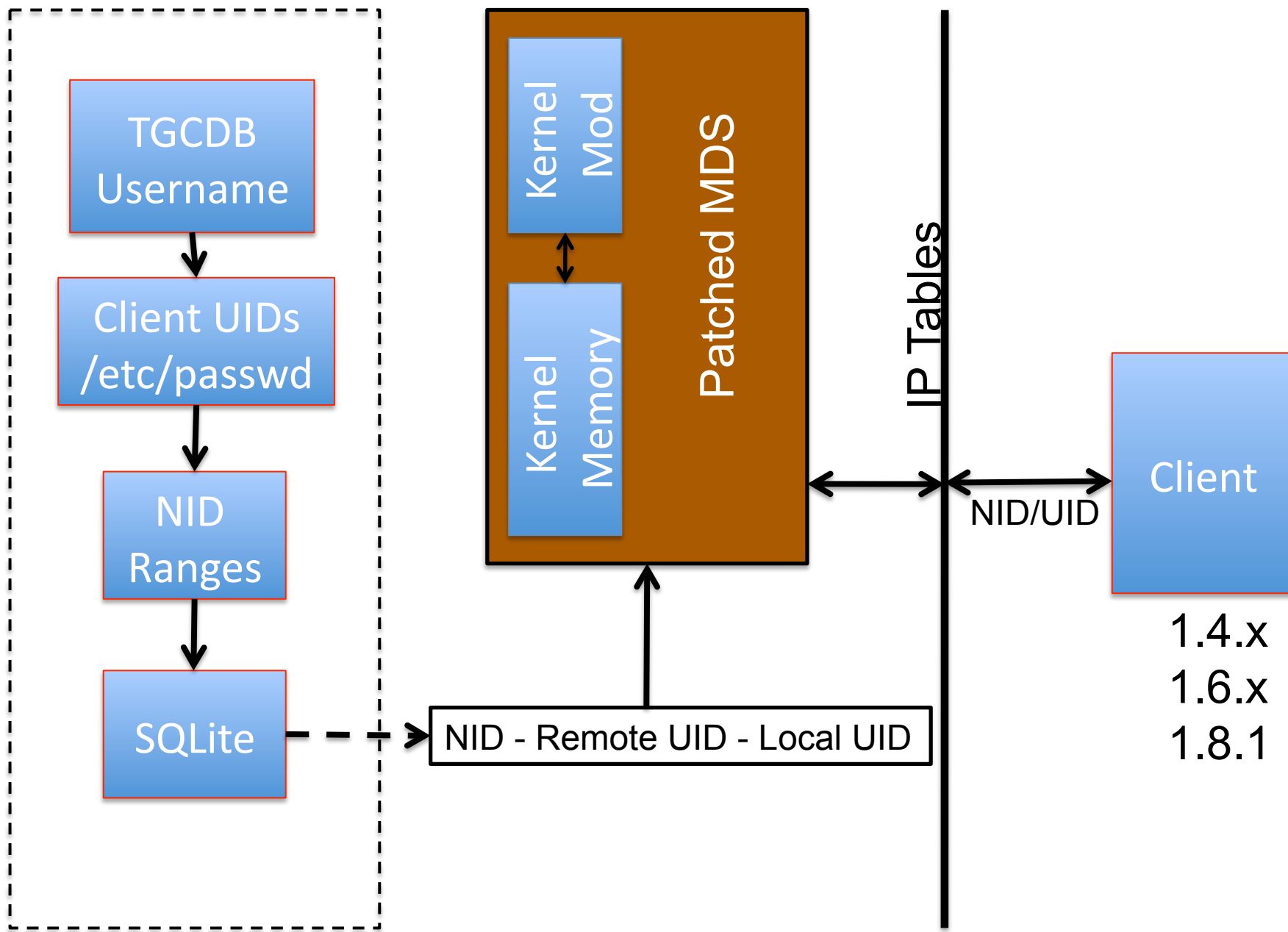
Only change MDS code

- Want to maximize clients we can serve

Simple enough to port the code forward

IU UID Mapping cont'd

- UID lookups on the MDS call a pluggable kernel module
 - Binary tree stored in memory
 - Based on NID or NID range
 - Remote UID mapped to Effective UID



UID Mapping

- Userspace – Kernel Space Barrier
 - Only crossed when we update the table
- Create a Forest of Binary Trees
 - Forward and Inverse Lookups for each UID
 - Time consumed for lookup is predictable
- Speed over Space
 - Consume memory rather than on the fly lookups
 - Every UID node consumes 6 Ints
 - 300 Users approximately 300KB

IU's Lustre WAN on the TeraGrid

- 8 Sites currently mounting IU DC-WAN
 - IU, LONI, NCSA, NICS, PSC, Purdue, SDSC, TACC
- 5 Sites mounting on compute resources
 - IU, LONI, NCSA, PSC, TACC
- Average of 93% capacity for the last quarter
- 2009 uptime of 96%
 - Filesystem availability to users

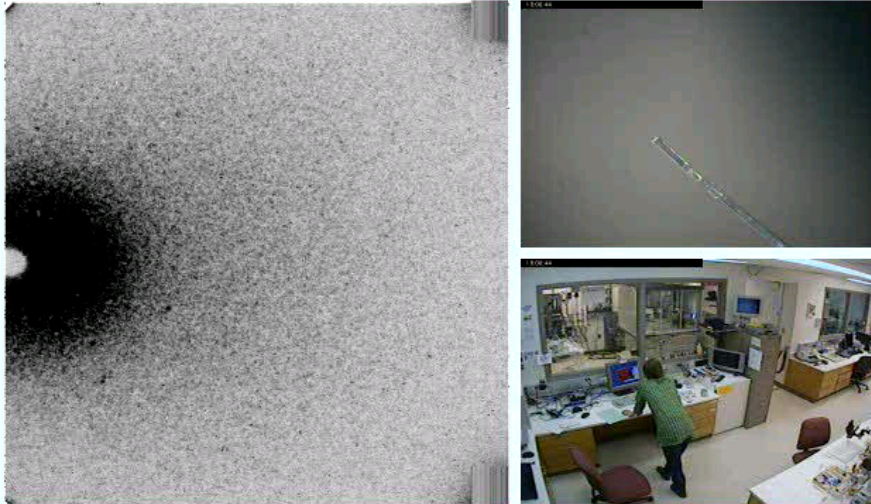
Data In and Data Out



<http://www.flickr.com/photos/davesag/4307240/in/set-799526/>

Common Instrument Middleware Architecture

CIMA

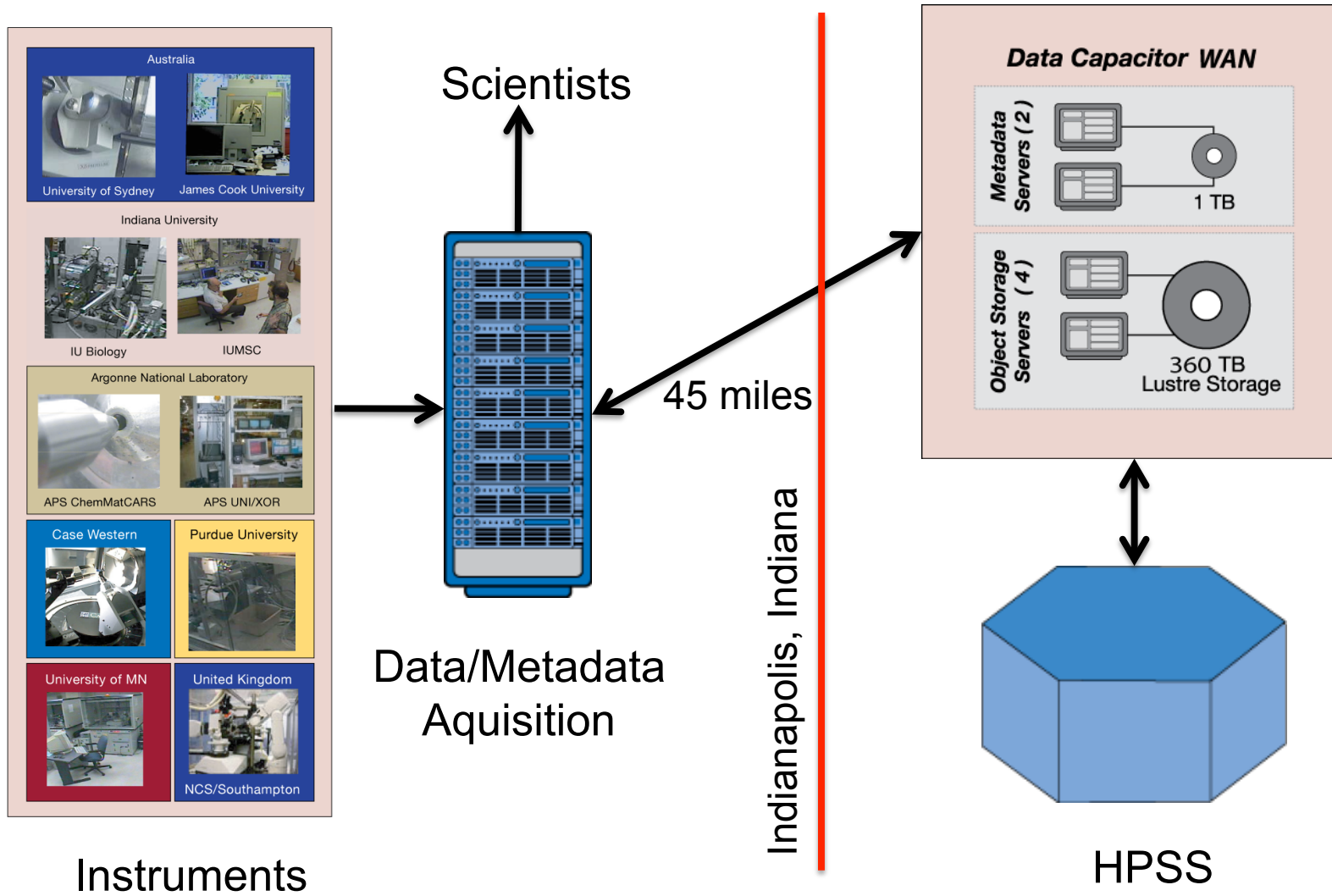


The image displays X-ray diffraction data and instrument status. On the left is a 2D diffraction pattern showing a dark spot on a light background. On the right are two smaller images: the top one shows a single diffraction spot, and the bottom one shows a person operating the instrument in a laboratory setting.

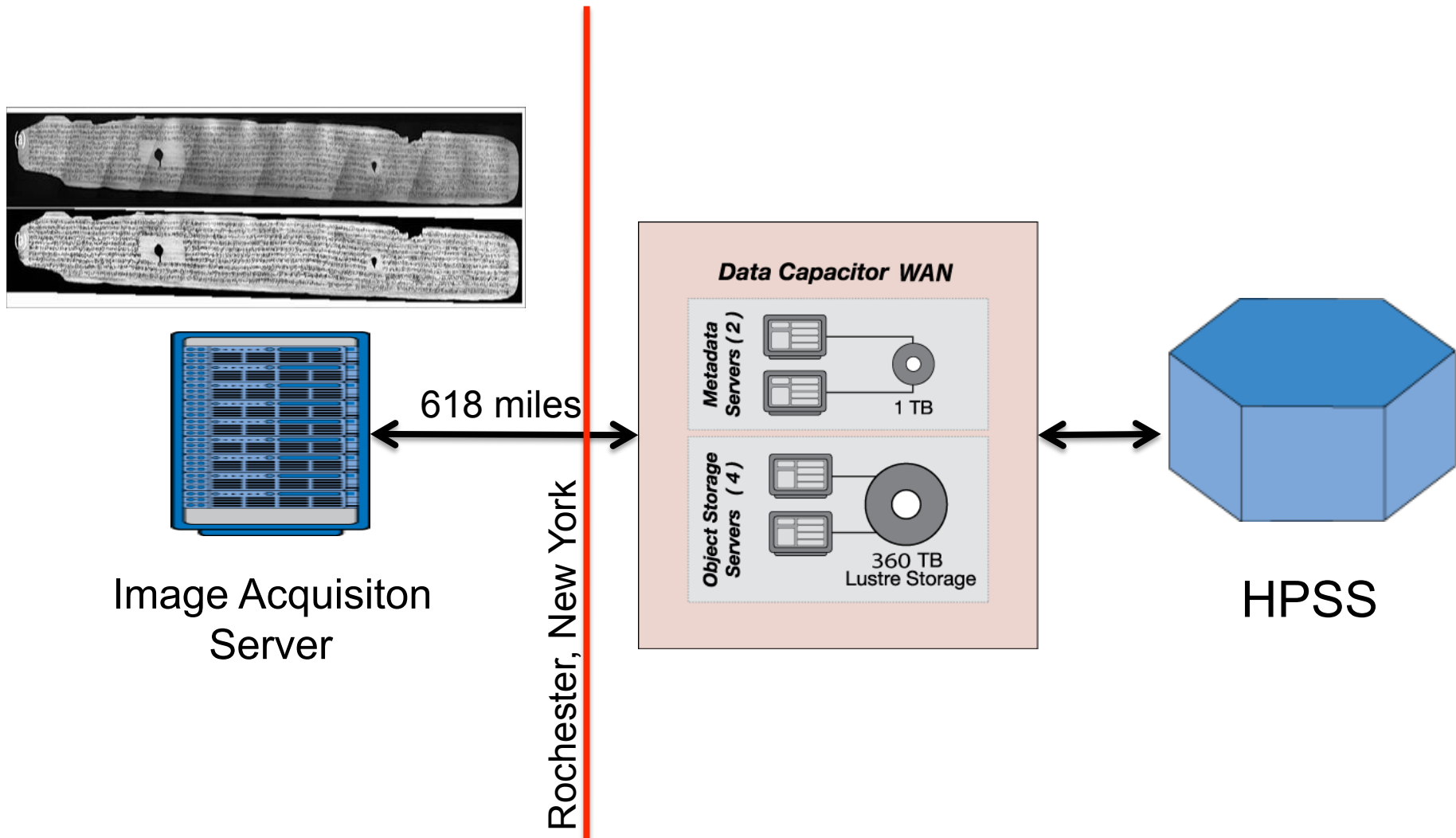
Time	2005-11-04 23:05:19 (UTC)	Instrument bay temperature (C)	19.50
Current data frame	053141.001	X-ray coolant water in (C)	16.40
Crystal temperature (C)	-148.20	X-ray coolant water out (C)	22.30
Instrument enclosure humidity (%)	45.60	CCD Chip Temperature (C)	-55.47
Instrument enclosure temperature (C)	22.70	Frame #	1
Instrument bay humidity	55.60		

X-ray diffractometer output

CIMA cont'd



Preserving Sarvamoola Granthas



One Degree Imager (ODI)

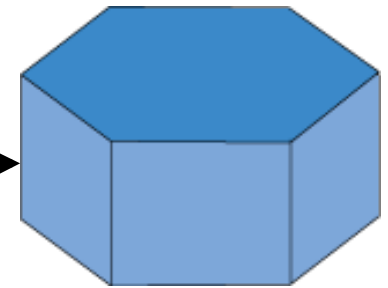
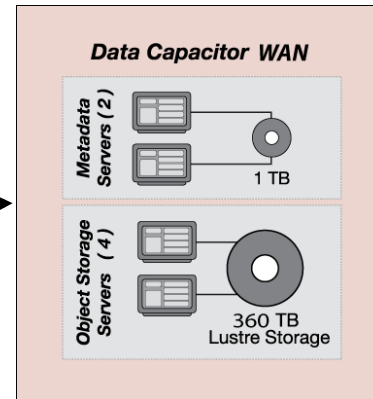
WIYN Telescope



NOAO/AURA/NSF

1726 miles

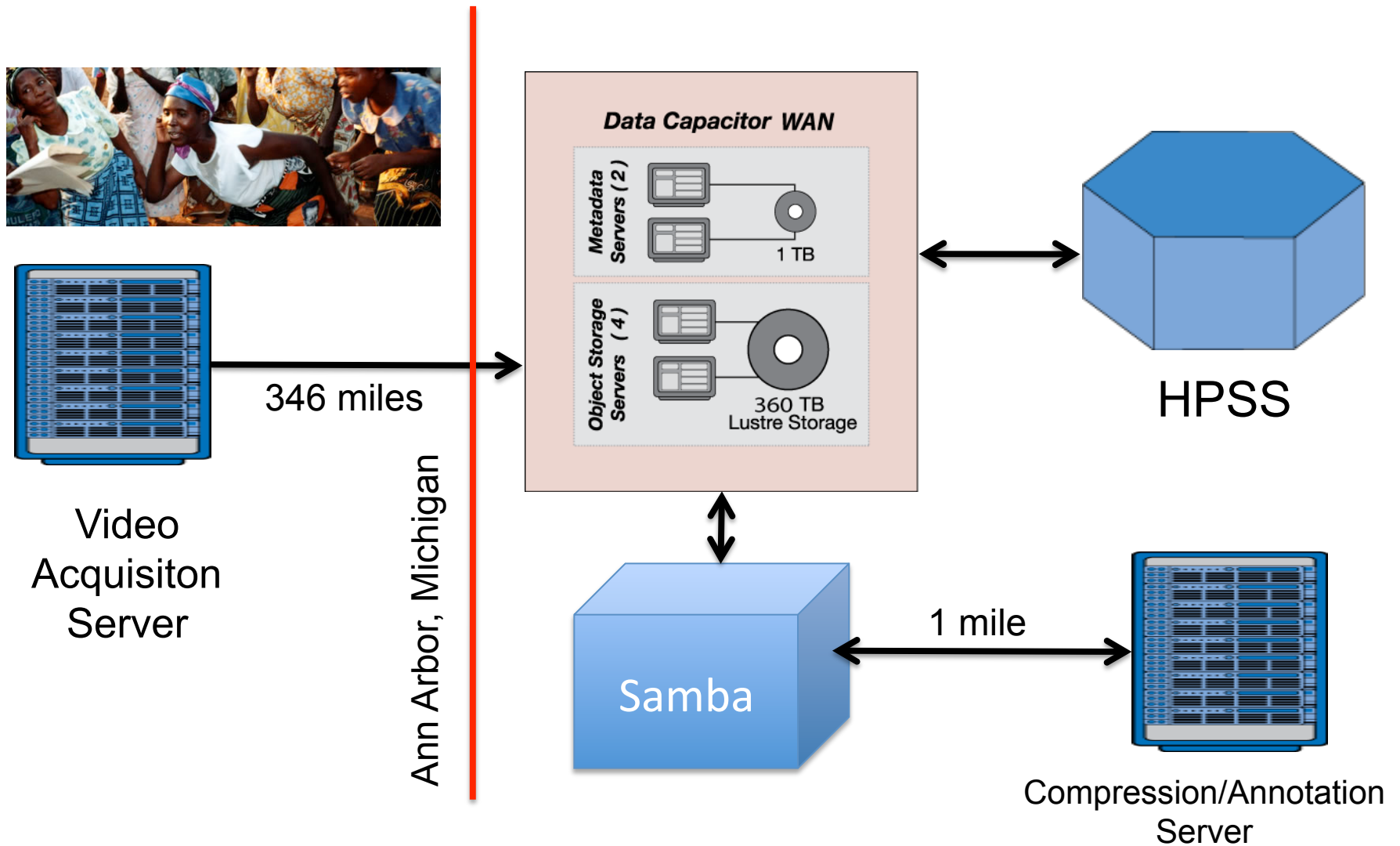
Tucson, Arizona



HPSS

Ethnographic Video for Instruction and Analysis

EVIA

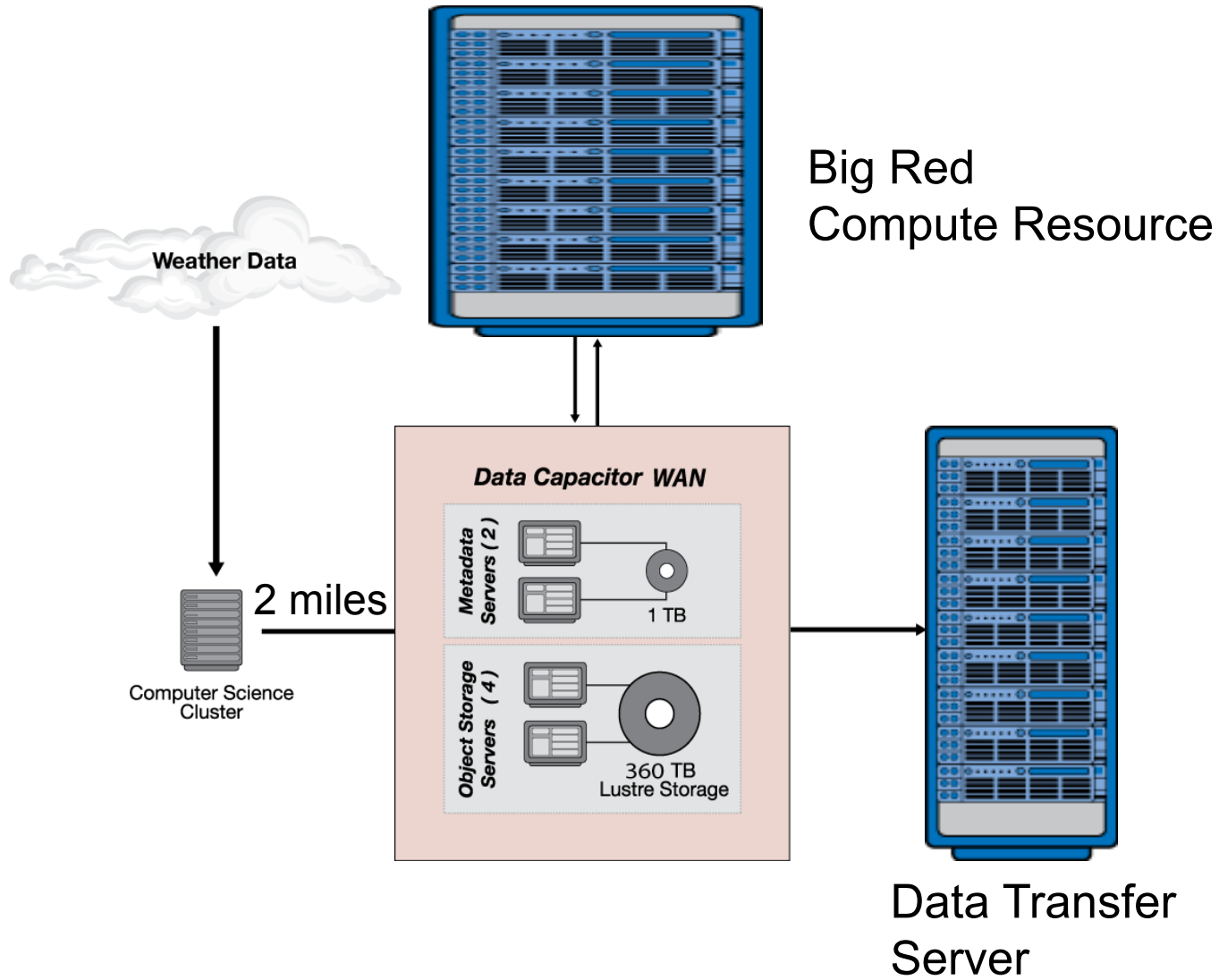


One Stop Resource Shopping

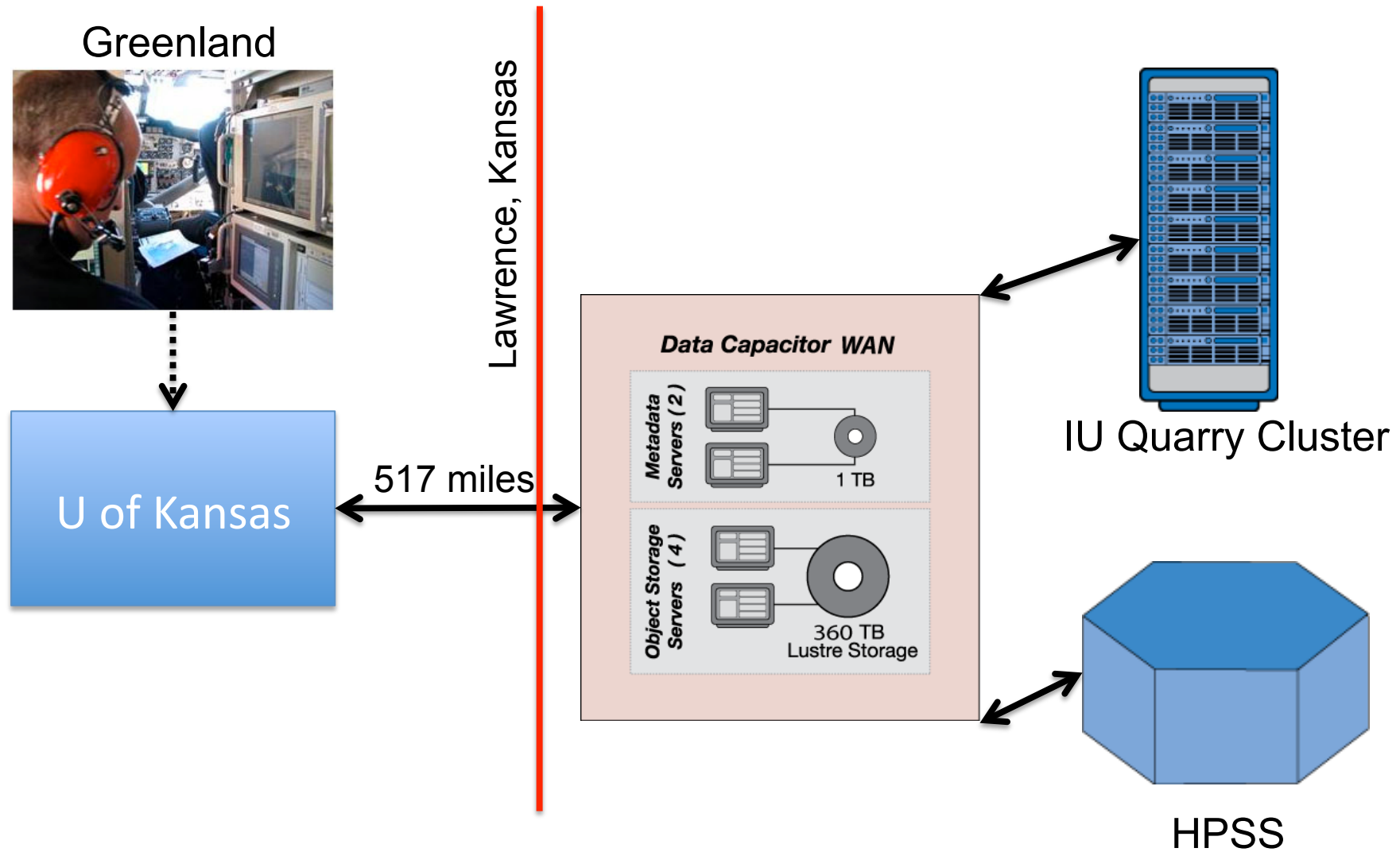


Linked Environments for Atmospheric Discovery

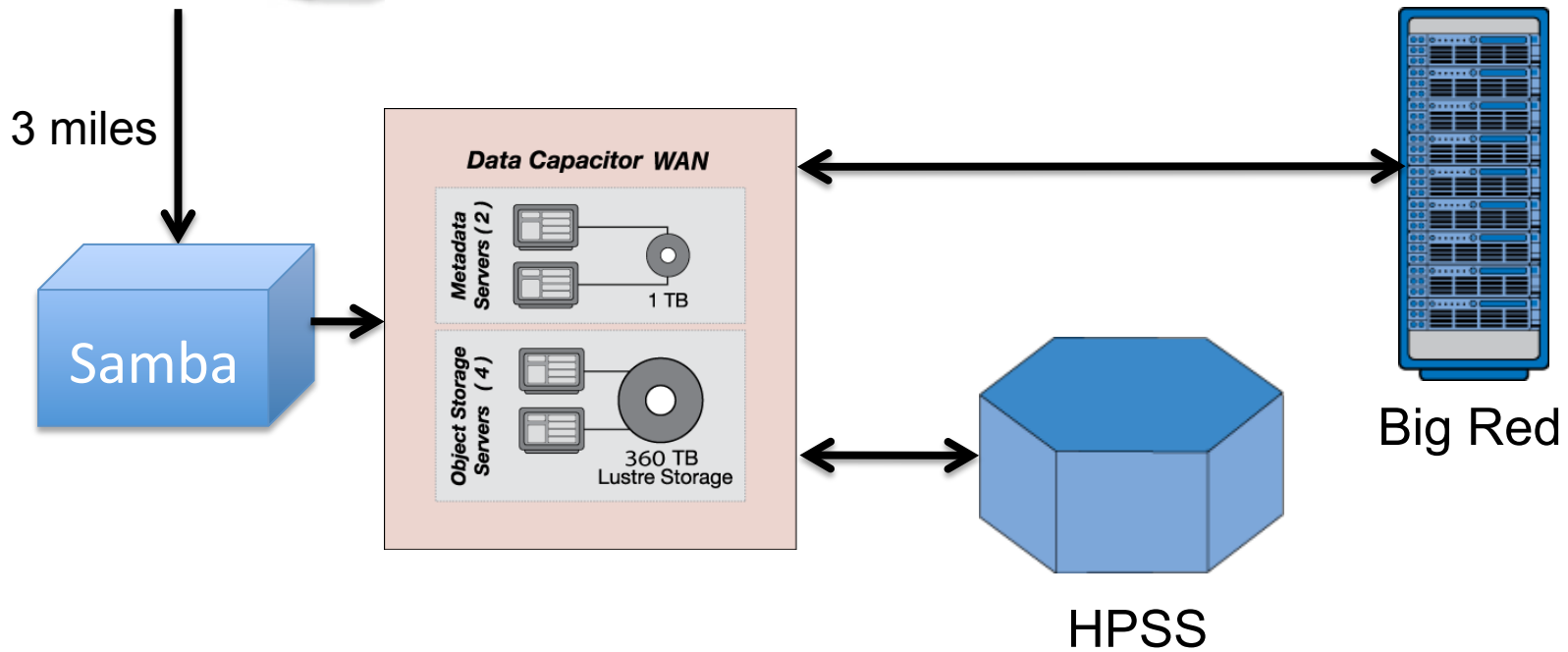
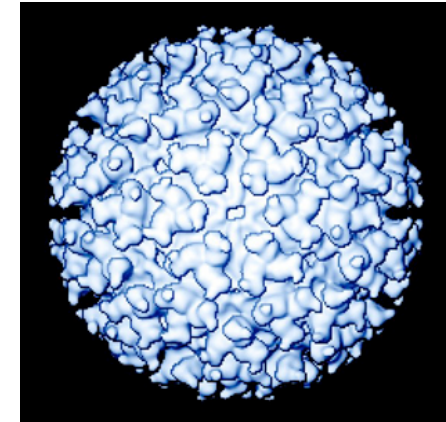
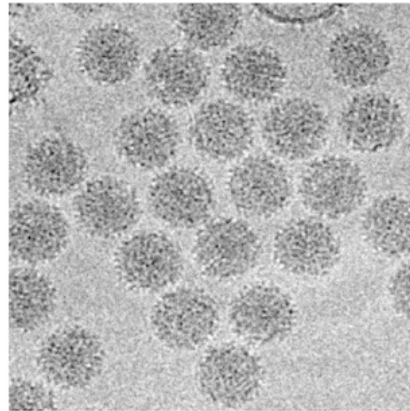
LEAD



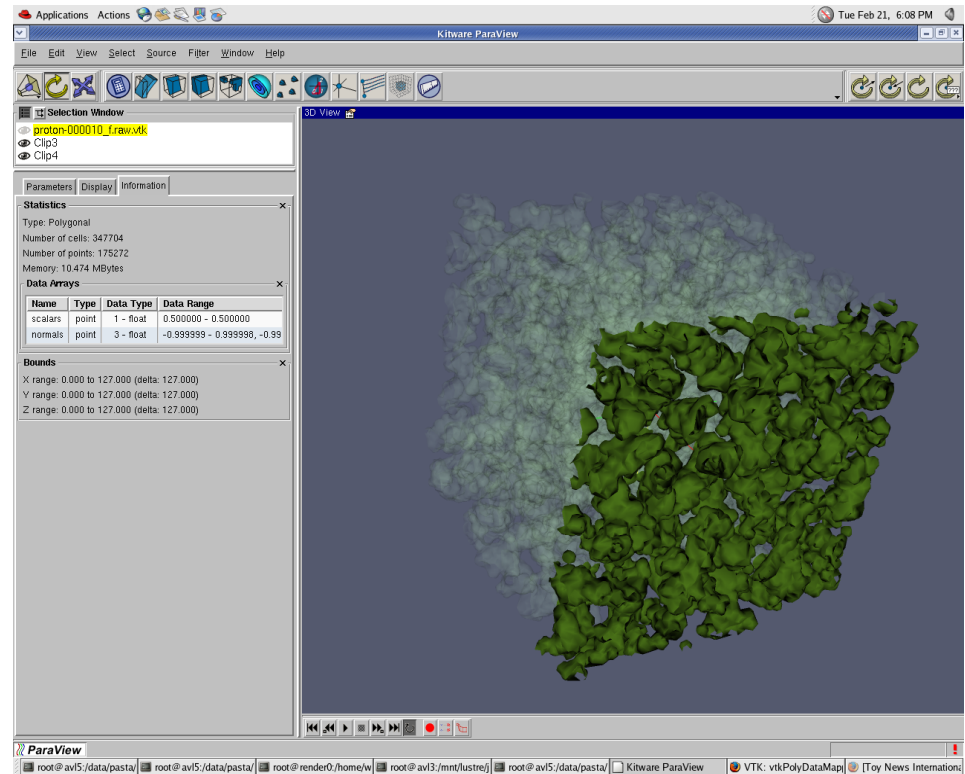
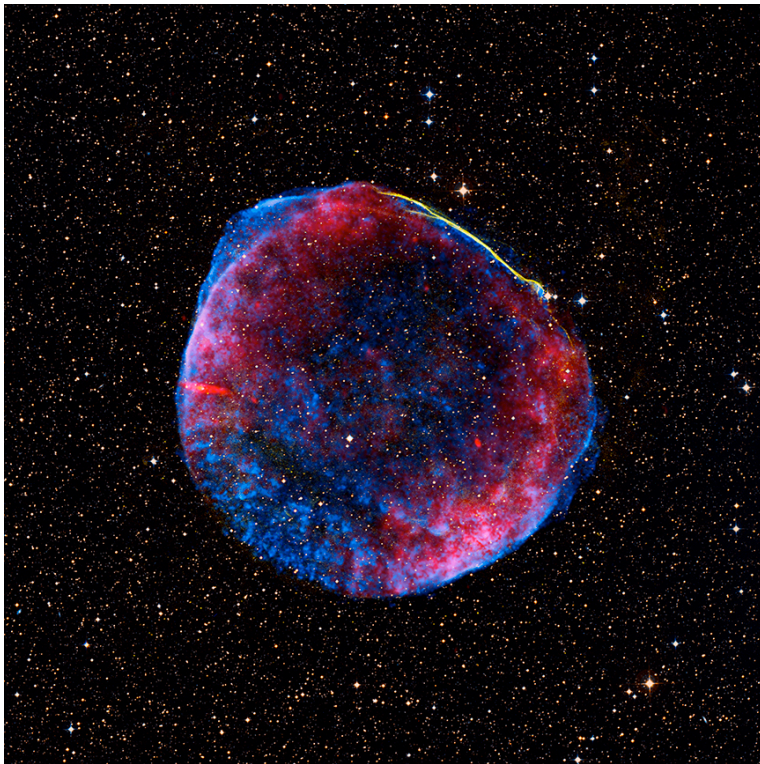
Center for the Remote Sensing of Ice Sheets (CReSIS) Workflow



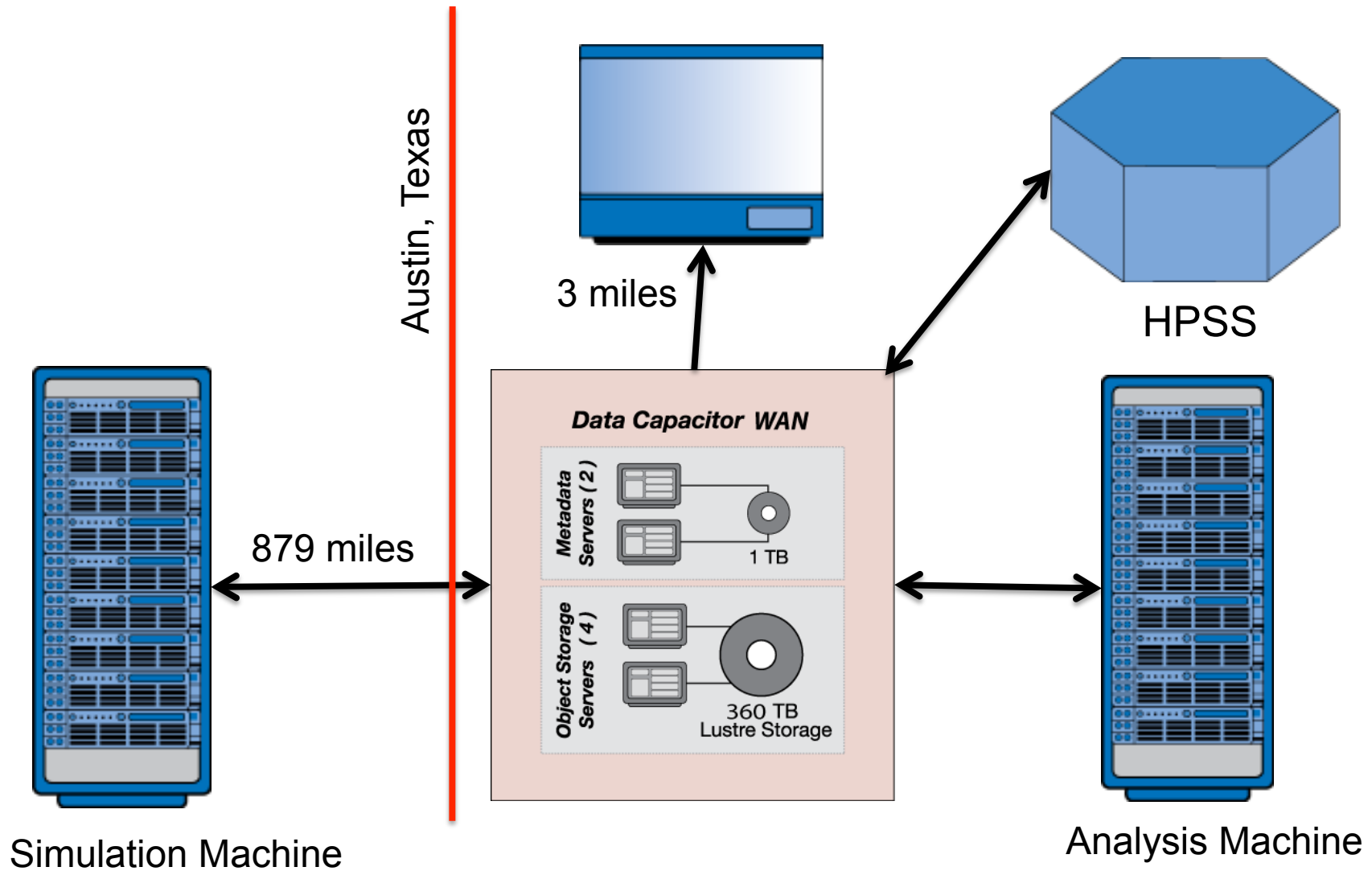
CRYO Electron Microscopy



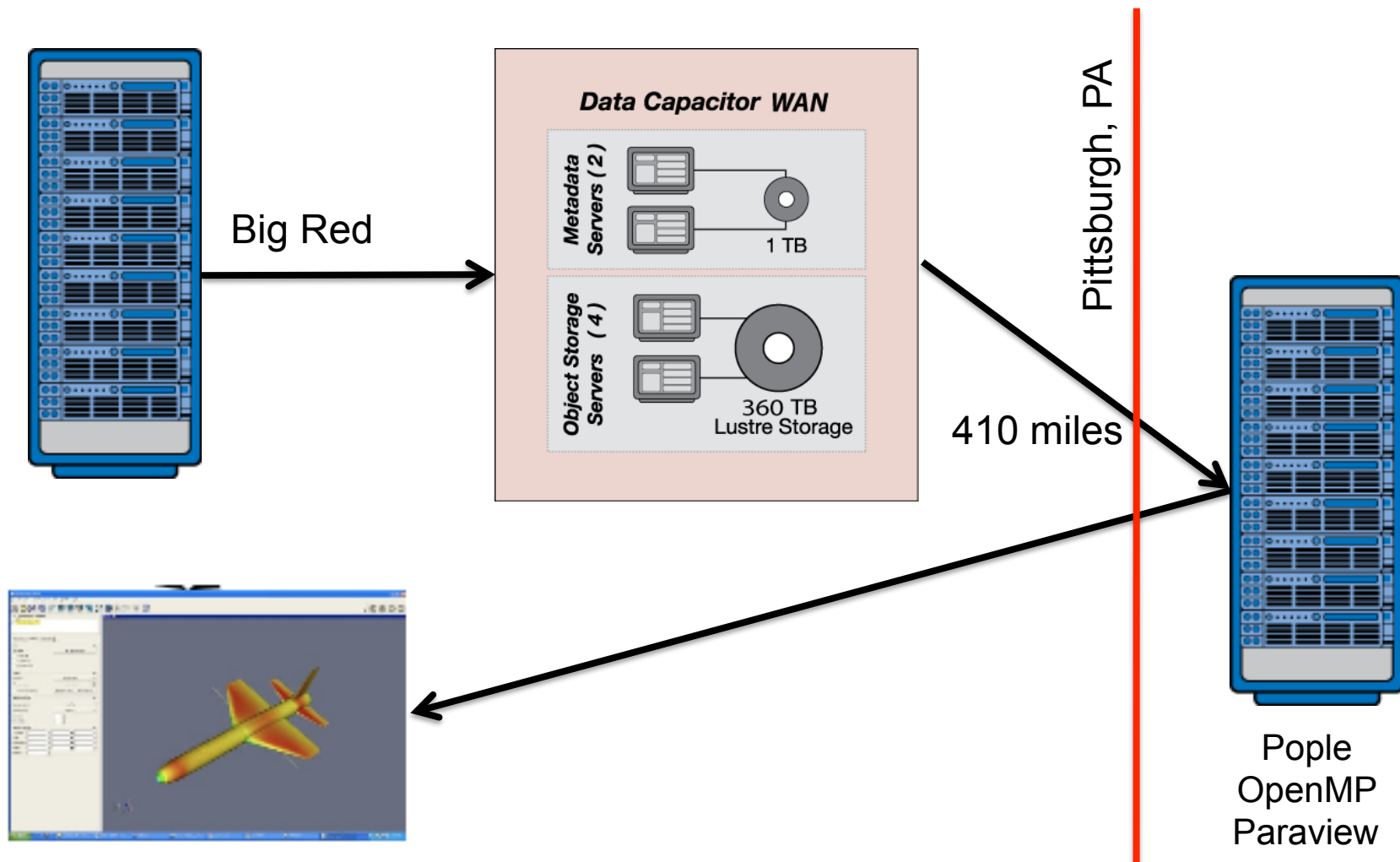
Equation of State Simulations and Plasma Pasta



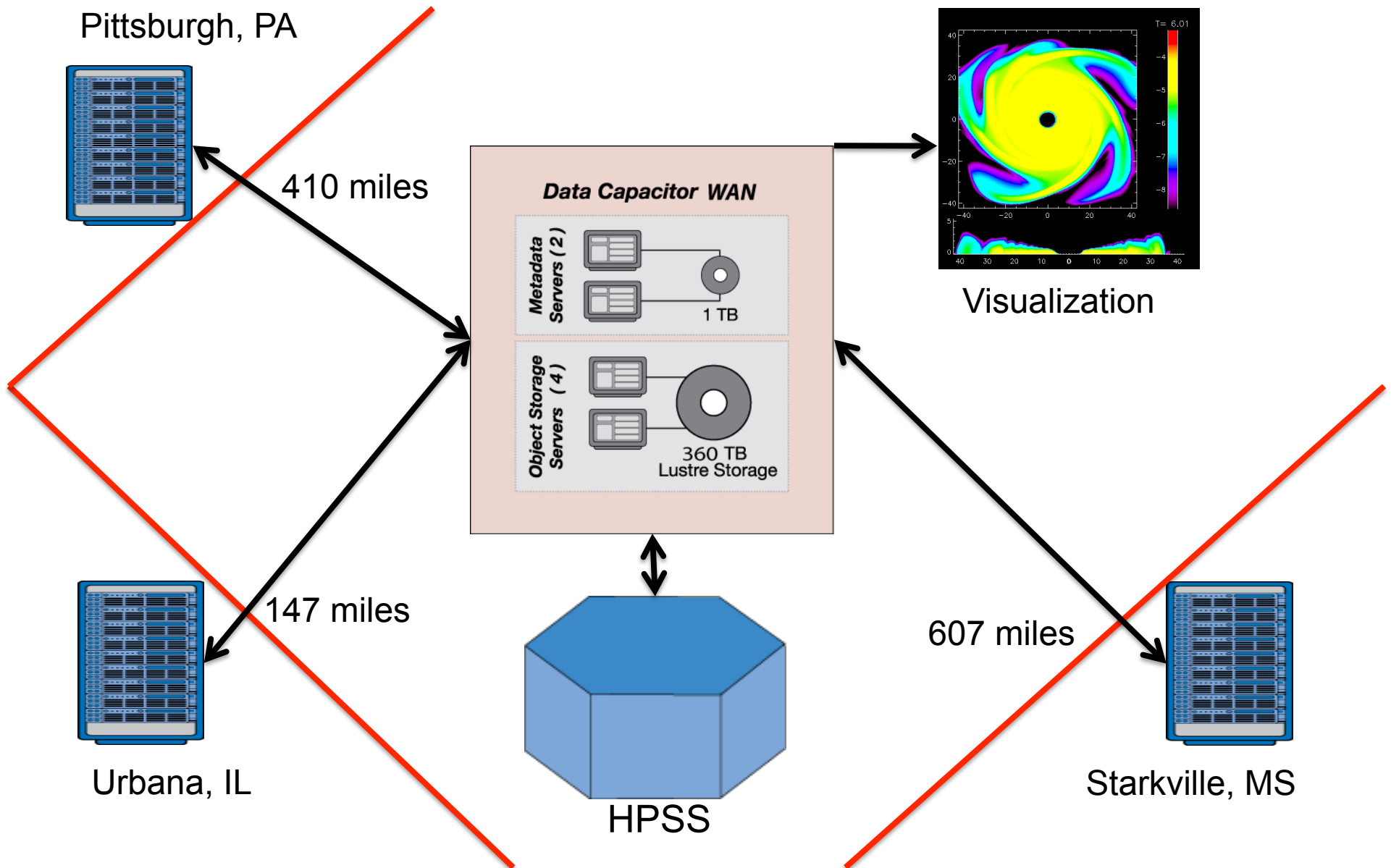
EOS and Plasma Pasta



Computational Fluid Dynamics



Gas Giant Planet Research



TeraGrid Future

- TeraGrid distributed Lustre WAN filesystem
 - 1.8.x
 - Distributed OSSs
 - NSF has funded servers to be deployed at five sites
 - IU, NCSA, NICS, PSC, TACC
 - IU's UID mapping code
 - PSC will run MDS
 - Will move to 2.0 code in the future

Non TeraGrid Exploration

- Dresden
 - ZIH (Technische Universität Dresden)
- Denmark
 - Risø – National Laboratory for Sustainable Energy
- Finland
 - Metsähovi Radio Observatory

Many Thanks

- Josh Walgenbach, Justin Miller, Nathan Heald, James McGookey, Resat Payli, Suresh Marru, Robert Henschel, Scott Michael, Tom Johnson, Chuck Horowitz, Don Berry, Scott, Teige, David Morgan, Matt Link (IU)
- Kit Westneat (DDN)
- Oracle support and engineering
- Michael Kluge, Guido Juckeland, Matthias Mueller (ZIH,Dresden)
- Thorbjorn Axellson (CReSIS)
- Greg Pike and ORNL
- Doug Balog, Josephine Palencia, and PSC
- Trey Breckenridge, Roger Smith, Joey Jones (Mississippi State University)

Support for this work provided by the National Science Foundation is gratefully acknowledged and appreciated (CNS-0521433)

Thank you!



Questions?

ssimms@indiana.edu

dc-team-1@indiana.edu

<http://datacapacitor.iu.edu>