

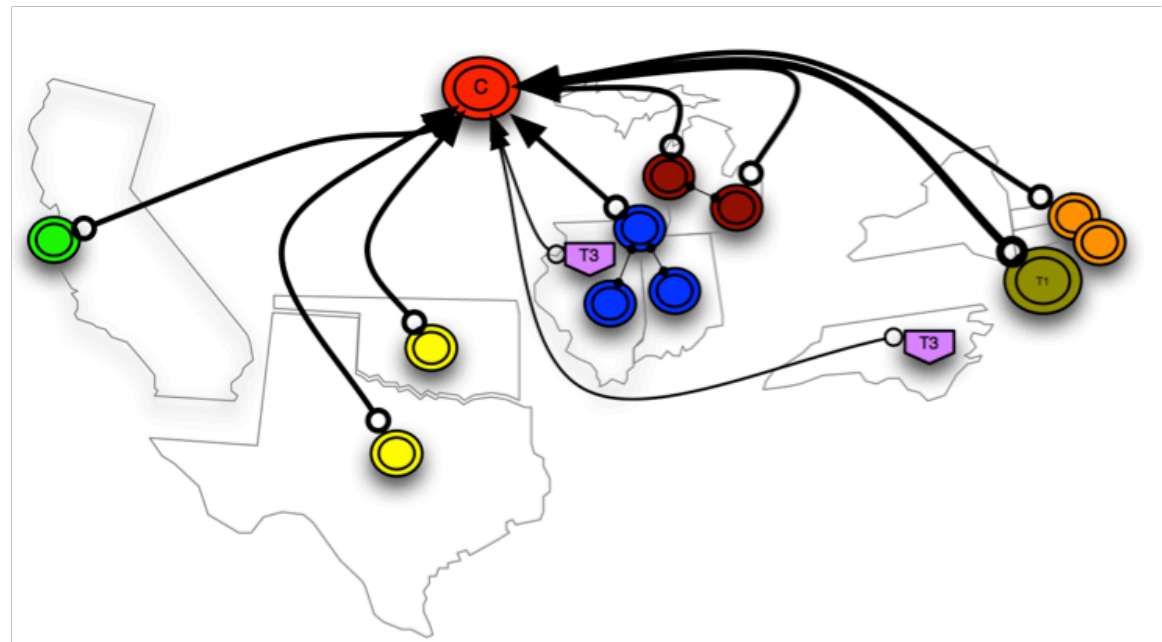
Federating ATLAS Data Stores using Xrootd

Rob Gardner
Lyon Federated Data Meeting
Sep 13, 2012

R&D Activity to Production

- 2011 R&D project FAX (Federating ATLAS data stores using Xrootd) was deployed over US Tier 1, Tier 2s and some Tier3s
- Feasibility testing monitoring, site integrations
- In June 2012 extended effort to European sites as an ATLAS-wide project

BNL Tier 1
AGLT2 (Tier 2)
MWT2 (Tier 2)
SWT2 (Tier 2)
SLAC (Tier 2)
ANL (Tier 3)
BNL (Tier 3)
Chicago (Tier 3)
Duke (Tier 3)
OU (Tier 3)
SLAC (Tier 3)
UTA (Tier 3)
NET (Tier 2)



Federation goals and usage in ATLAS

- Common ATLAS namespace across all storage sites, accessible from anywhere
- Easy to use, homogeneous access to data
- Use as failover for existing systems
- Gain access to more CPUs using WAN direct read access
- Use as caching mechanism at sites to reduce local data management tasks
- WAN data access group formed in ATLAS (see Torre Wenaus' talk later) to determine use cases & requirements on infrastructure

Prerequisites to federation

- Concept of a global name or easy way for analysis users to name input files
- Means to integrate to site's backend storage
 - Name-to-name mapping between LFC name and physical storage path for dCache, Xrootd, GPFS/Lustre(Posix), EOS (in progress)
 - Redirect client to data or proxy server
- System of Xrootd 'redirectors' to open files based network proximity (later cost functions)
- Analysis code that intelligently caches (eg. TTreeCache)

Name translation in ATLAS

- An ATLAS global namespace convention based on unique dataset and filenames has been developed
- While the LFN (and GUID) is unique, the physical path at sites is not
- Use of space tokens as an organizational data management tool means files can live in more than one directory
- Prior to FTS 2.8 and re-writes, DQ2 suffixes appended to subsequent transfer attempts, preventing deterministic lookups
- Therefore an LFC query is required, at least right now
- Rucio (DQ2 replacement) offers a chance to revisit this

Recent Progress (I)

- dcache-xrootd doors now functional
 - new N2N working
 - recipe for adding Tier I sites and others using dCache
 - requires a “helper” xrootd+cmsd pair & billing database publisher for monitoring data
 - New version will have have monitoring plugins for pool services to publish data at source, similar to Xrootd
- Deployment of read-only LFC for unauthenticated LFC lookup (data access still uses GSI security requiring ATLAS voms attribute)






Recent Progress (2)

- Development of general ATLAS monitoring requirements & work with dcache and xrootd teams to standardize (will come out in new xrootd and dcache releases - cf. A.Hanushevsky's talk)
- WLCG dashboard, detailed monitoring services discussions and prototype work (cf. Julia Andreeva, Domenico Giordano talks)
- First European sites federated - from the UK and DE clouds
- Redirectors setup at CERN, good support from
- FAX-enabled pilot available and tested at two sites, implementing “missing file” use-case

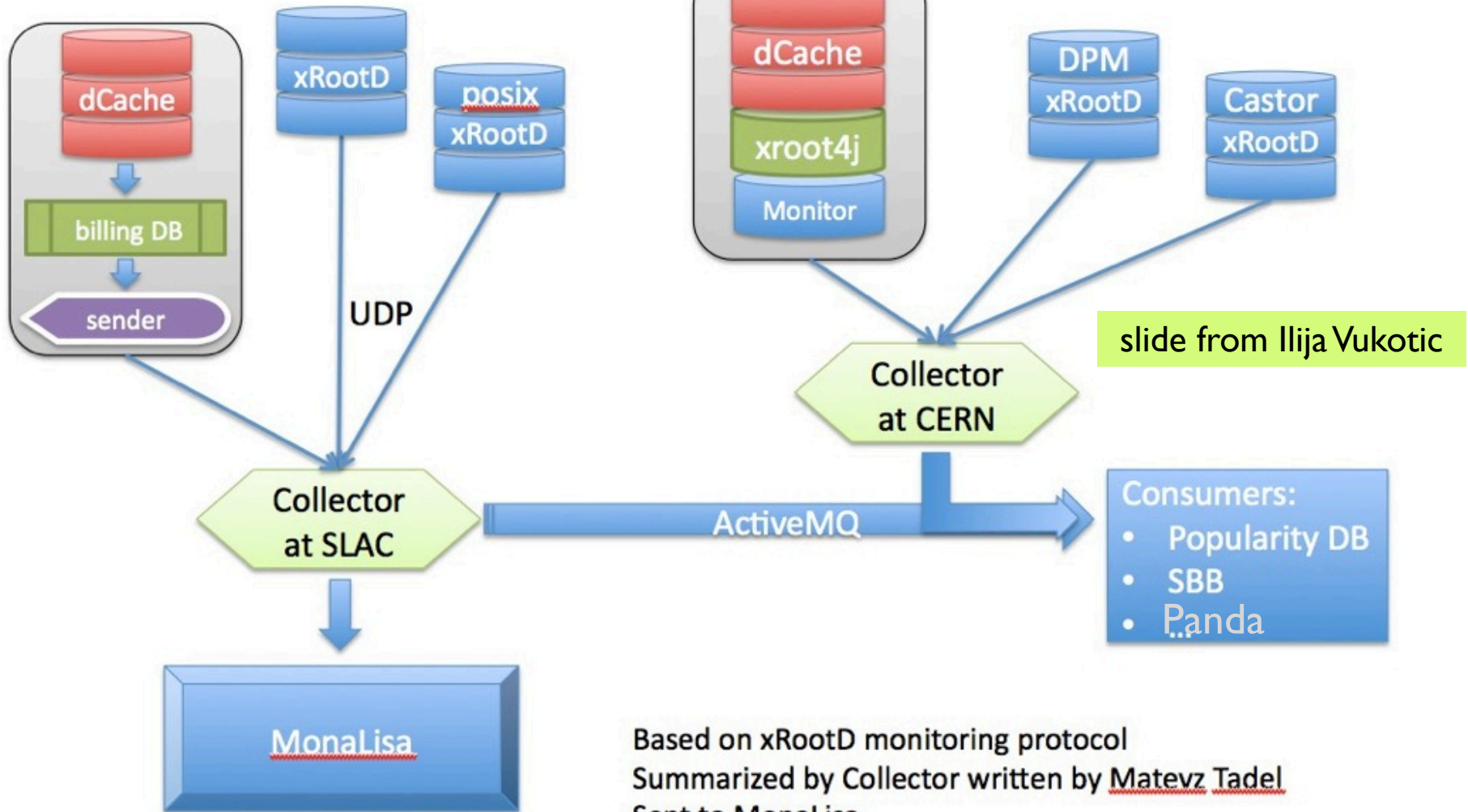
Validating new sites, towards production federation readiness over many regions

Site certification for federation

Notation:  completed  work is in progress  to do  problem

Site	Cloud	SE type-(door)	Regional RD	Federated	X509	GlobalN2N	FAX status mon	UDP Collector	Redir Cloud	Redir Gobal	Fallover	Analy Test
BNL	US	dcache	glrd.usatlas.org									
AGLT2	US	dcache-xrootd	xrd-central to glrd.usatlas.org									
MWT2.org	US	dcache-xrootd	xrd-central to glrd.usatlas.org									
MWT2_UC,IU	US	xrootd	xrd-central to glrd.usatlas.org									
NET2	US	GPFS	glrd.usatlas.org									
SWT2 (UTA)	US	xrootd	glrd.usatlas.org									
SWT2 (OU)	US	Lustre	xrd-central to glrd.usatlas.org									
SLAC	US	xrootd	glrd.usatlas.org									
Wuppertal	DE	dcache-xrootd	atlas-xrd-de.cern.ch									
LRZ-LMU	DE	dcache-xrootd via xrootd proxy	atlas-xrd-de.cern.ch									
Edinburgh	UK	DPM	atlas-xrd-uk.cern.ch									
Glasgow	UK	DPM	atlas-xrd-uk.cern.ch									
Oxford	UK	DPM	atlas-xrd-uk.cern.ch									
QMUL	UK	Storm/Lustre	atlas-xrd-uk.cern.ch									
EOS	EU	EOS	atlas-xrd-eu.cern.ch									
Dubna	RU		atlas-xrd-ru.cern.ch									

Monitoring infrastructure



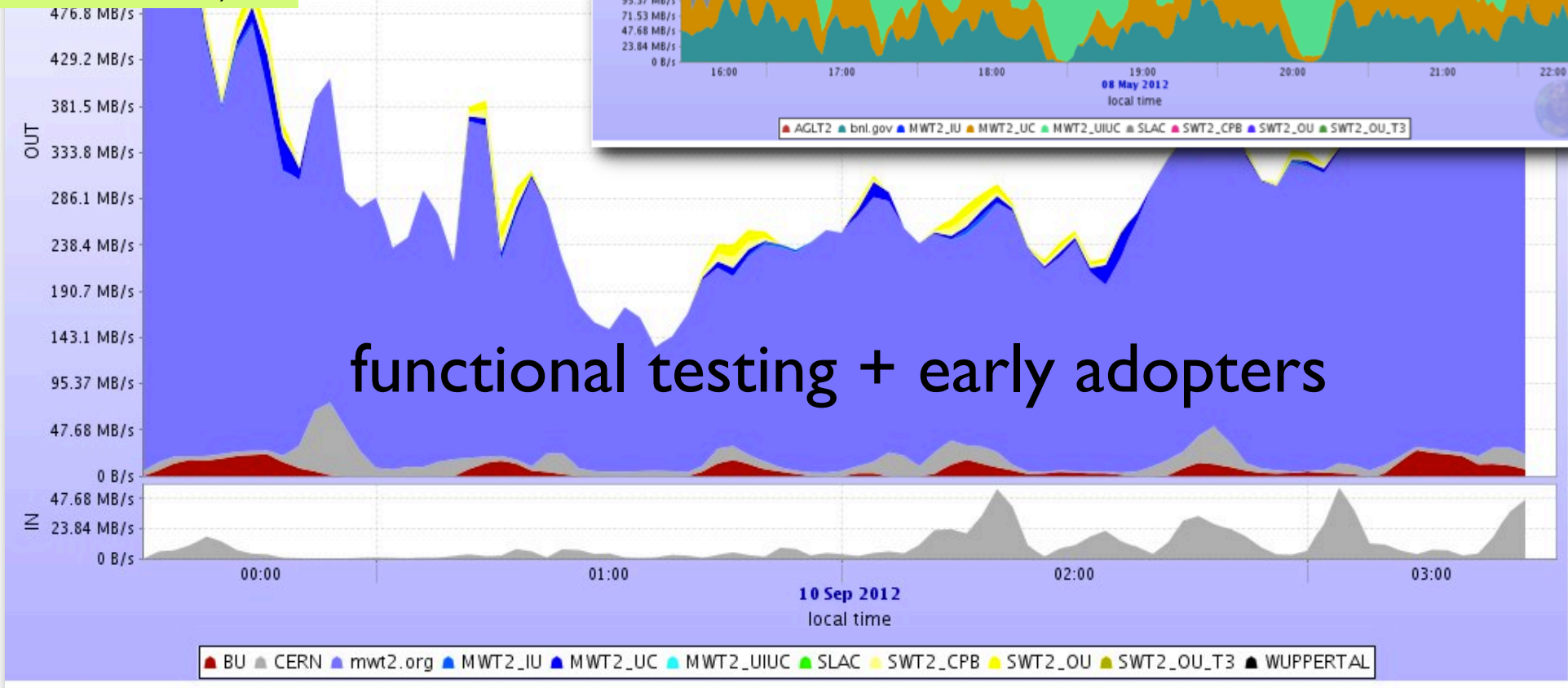
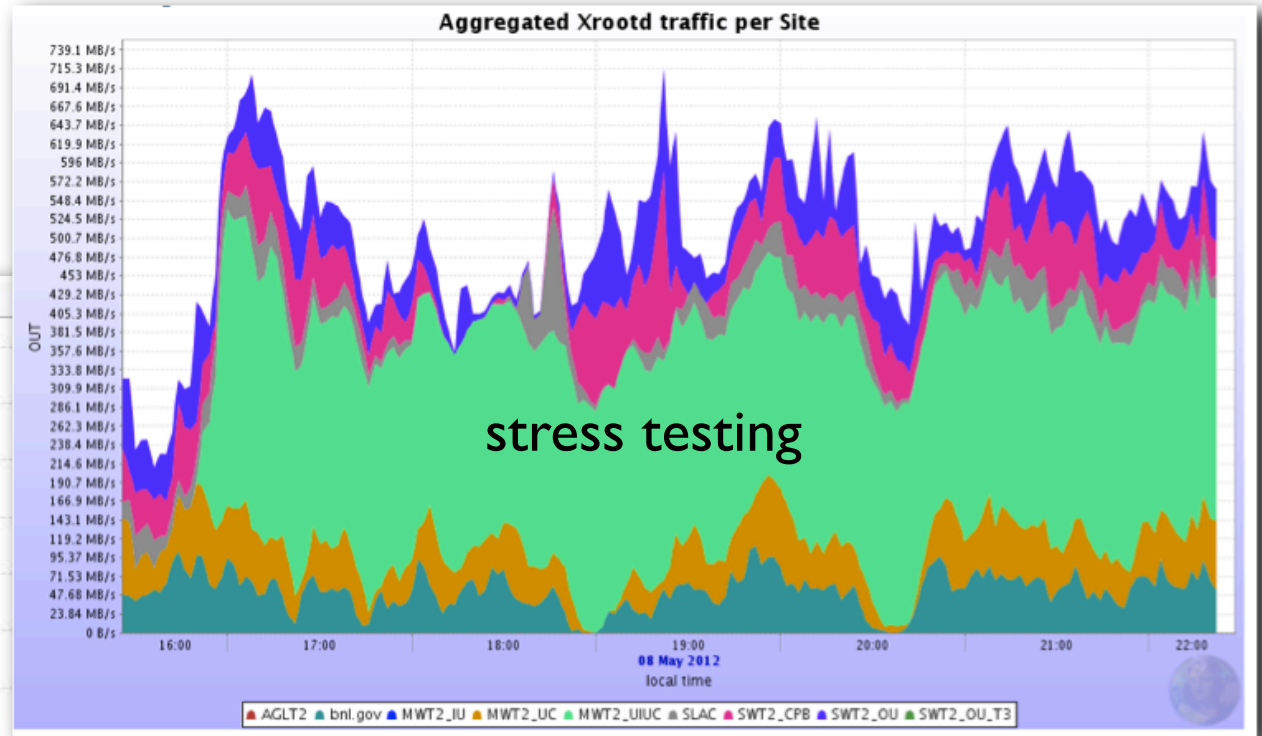
Based on xRootD monitoring protocol
 Summarized by Collector written by Matevz Tadel
 Sent to MonaLisa
 Sent to ActiveMQ

Collaborative work w/WLCG, AAA, Xrootd; cf Julia Andreeva et. al.
<https://twiki.cern.ch/twiki/bin/view/LCG/WLCGTransferMonitoring>

More in later talks!

Summarized monitoring

MonaLisa visual fed from UDP Collector at SLAC (much help from Matevz Tadel)



Xrootd DATA POPULARITY

CERN IT Experiment Support



PROOF OF CONCEPT ON ATLAS EOS

Home xrd monitor Plots Tables

FAX REAL-TIME MONITORING

Provide the real time snapshot of the last 10 mins of collected Xrootd monitoring data. UTC Time is reported.

Show 10 entries

starttime	endtime	server domain	client domain	read bytes	write bytes	file size	filename
2012-09-11 15:24:06	2012-09-11 15:25:50	slac.stanford.edu	uchicago.edu	797152257	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:25:17	2012-09-11 15:25:24	atlas-sw2.org	atlas-sw2.org	797152257	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:24:14	2012-09-11 15:25:09	atlas-sw2.org	uchicago.edu	797152257	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:24:12	2012-09-11 15:25:07	bu.edu	uchicago.edu	0	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:24:14	2012-09-11 15:25:04	uchicago.edu	mwt2.org	0	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:23:28	2012-09-11 15:24:38	uchicago.edu	oceph.ou.edu	0	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:22:39	2012-09-11 15:24:13	slac.stanford.edu	uchicago.edu	797152257	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:23:27	2012-09-11 15:24:07	bu.edu	oceph.ou.edu	0	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:23:24	2012-09-11 15:23:58	slac.stanford.edu	oceph.ou.edu	797152257	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root
2012-09-11 15:23:29	2012-09-11 15:23:54	atlas-sw2.org	oceph.ou.edu	797152257	0	797152257	/atlas/dq2/user/ilijav/HCTest/user.ilijav.HCTest.1/group.test.hc.NTUP_SMWZ.root

Showing 1 to 10 of 76 entries

Get source JSON || Resource URL: [/eosat/xrdrealmonfax?](https://eosat/xrdrealmonfax?)

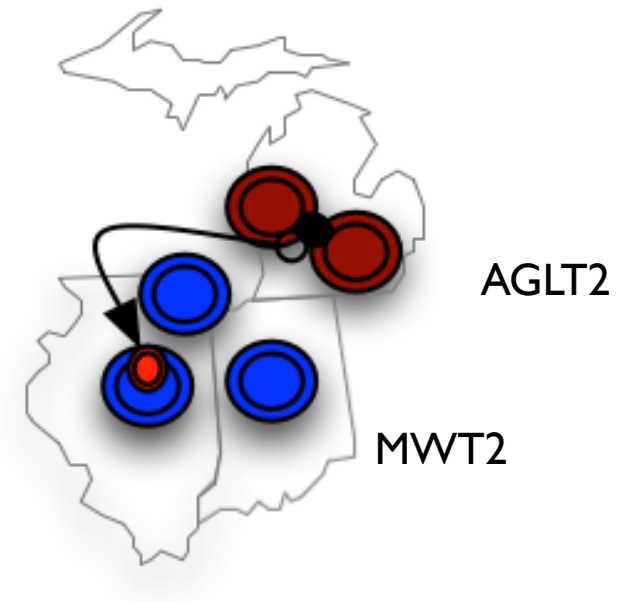
publishing detailed Xrootd stream from
SLAC to ActiveMQ at CERN

cf. Domenico Giordano's talk

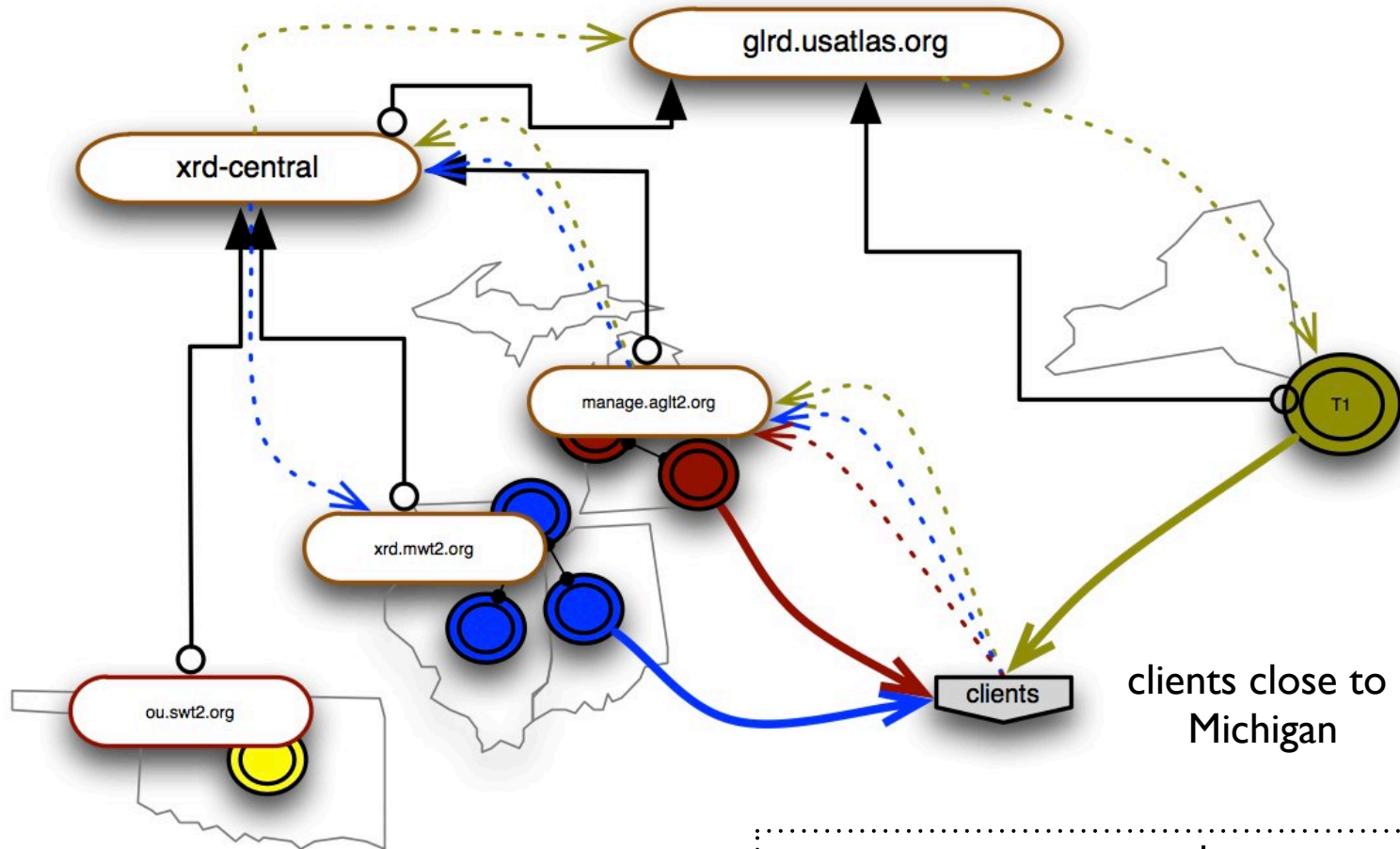
Capabilities: regional redirection

```
root://xrd-central.usatlasfacility.org/global-name
```

- MWT2 and AGLT2 internally are multi-site federations using dCache (5 total)
- Combined storage ~ 5 PB and about 10K job slots
- ANALY_AGLT2 & ANALY_MWT2 within 7 ms RTT
- Added SWT2 Oklahoma (< 10ms)
- Not found here? Fails over to US cloud level redirector gird.usatlas.org



Regional redirection



case 2: control data
data

client redirected to file on site
in region

Regional WAN analysis: jobs running on one site reading data from itself and two others

PandaID, Owner, Working group	Job	Status	Created	Time to start	Duration	Ended/ Modified	Cloud/Site, Type
1594538009 Robert W. Gardner Jr.	jobsetID=2511 runGen-00-00-02	finished	2012-09-06 07:46	2 days, 1:27:14	0:07:16	09-08 09:20	US/US.ANALY OU_OCHEP_SWT2, analysis-run
Out: user.rwg.SMExample OU_OCHEP_SWT2_xrd-central.001/							

Associated build job: [11141357](#)

Job ran on Oklahoma analysis site
read data from two nearby federations

Panda monitor

Run jobs in this job set: [11141358](#)

Job 1594538009 details

3 files for job 1594538009:

Filename	Type	Status	Dataset
user.rwg.0906074600.297152.lib_002511.lib.tgz guid=0dd239c9-309a-4df4-ad35-2f3552f8ae1b	input	ready	user.rwg.0906074600.297152.lib_002511
user.rwg.002511_1594538009.log.tgz guid=e42d4df2-c6d0-499a-b731-4479938d7351 Space token OU_OCHEP_SWT2_USERDISK	log	ready	user.rwg.SMExample OU_OCHEP_SWT2_xrd-central.001/ (destination block: user.rwg.SMExample OU_OCHEP_SWT2_xrd-central.001.1209060246)
user.rwg.002511_00025.SMWZExampleXYZ.log.tgz guid=eac1209e-747e-4efe-ab1d-20c7d05bf664 Space token OU_OCHEP_SWT2_USERDISK	output	ready	user.rwg.SMExample OU_OCHEP_SWT2_xrd-central.001/ (destination block: user.rwg.SMExample OU_OCHEP_SWT2_xrd-central.001.1209060246)

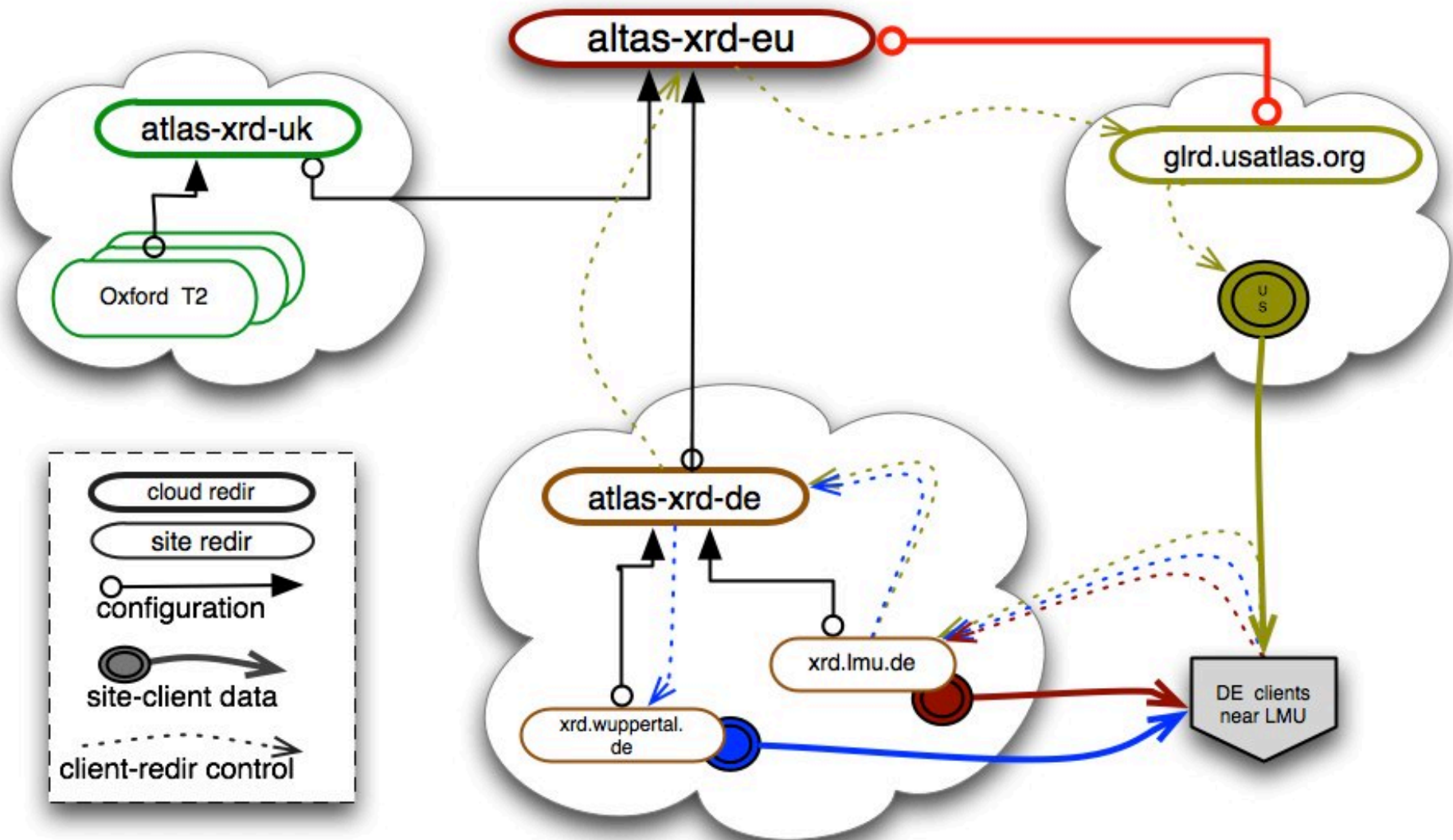
root://xrd-
central.usatlasfacility.org//
filename

read data from:
MWT2 (128.135.x.x)
AGLT2 (192.41.x.x)
OU (129.15.x.x)

```

2012-09-08 09:14:14.486698 : INFO connect: 128.135.158.237:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 128.135.158.190:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 192.41.230.187:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 192.41.230.21:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 192.41.236.55:23197 cmd: SMWZd3pdExample
2012 : INFO connect: 128.135.158.180:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 192.41.236.63:23982 cmd: SMWZd3pdExample
2012 : INFO connect: 128.135.158.252:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 128.135.158.186:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 192.41.236.60:24377 cmd: SMWZd3pdExample
2012 : INFO connect: 192.41.230.23:20398 cmd: SMWZd3pdExample
2012 : INFO connect: 128.135.158.221:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 128.135.158.189:1094 cmd: SMWZd3pdExample
2012 : INFO connect: 192.41.236.60:23501 cmd: SMWZd3pdExample
    
```

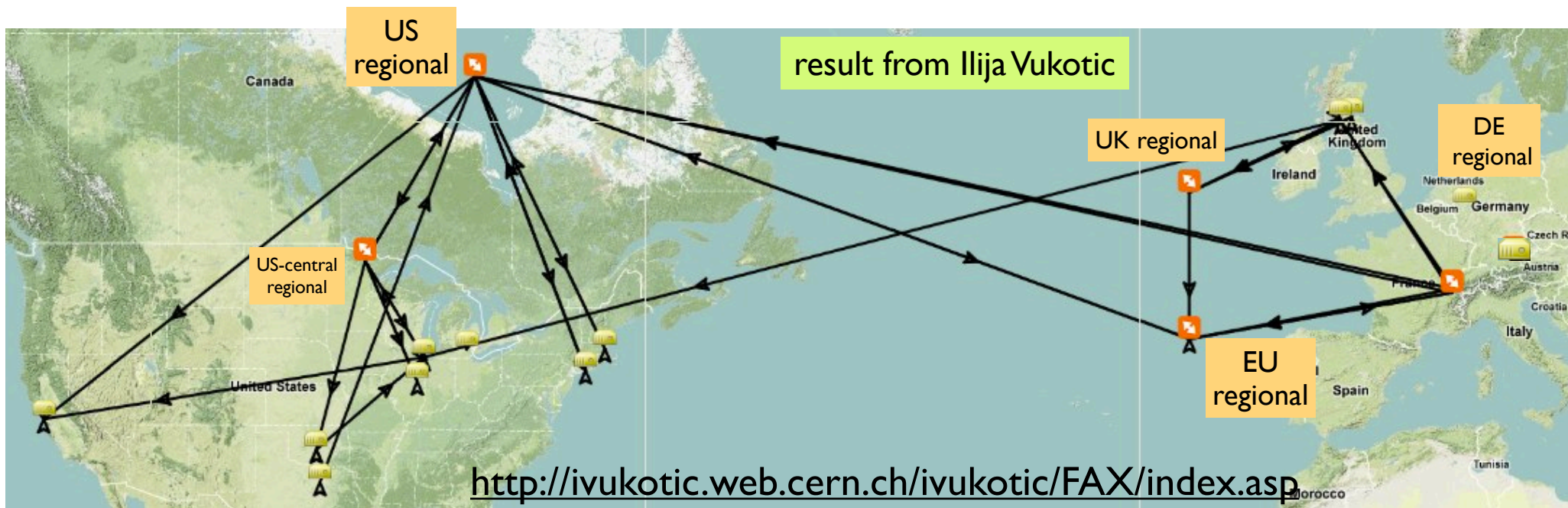
Four levels of redirection: site-cloud-zone-global



Start locally - expand search as needed

Topology validation

- Launch jobs to every site, test reading of site-specific files at every other site
- Parse client logs to infer resulting redirection

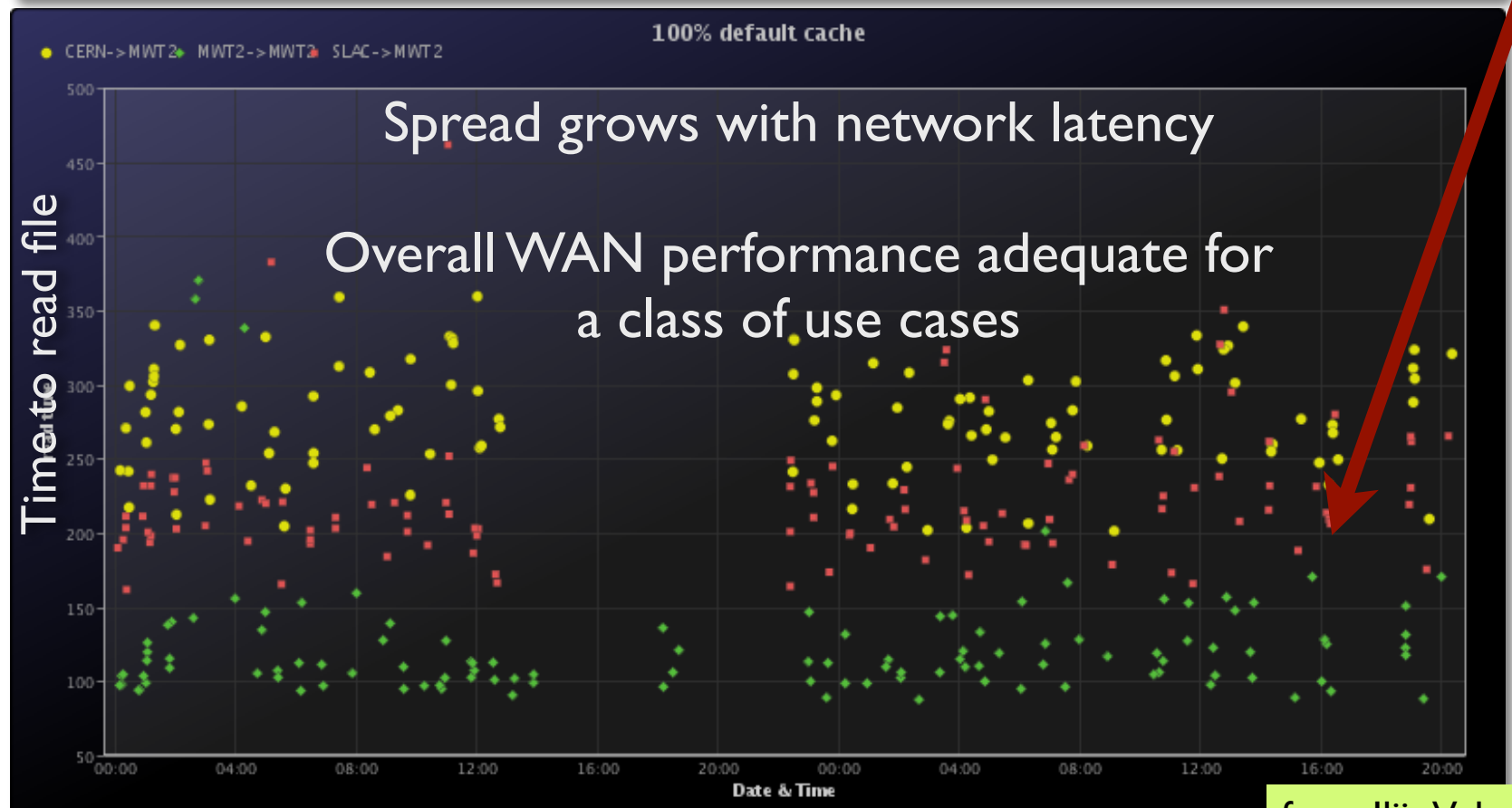
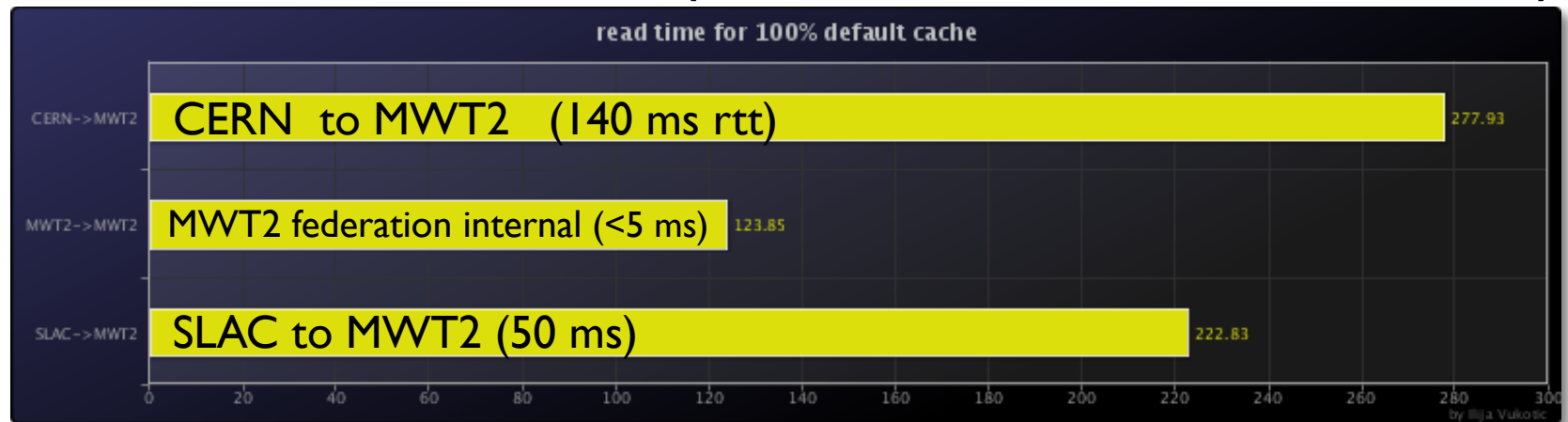


XRD redirector



federating site 16

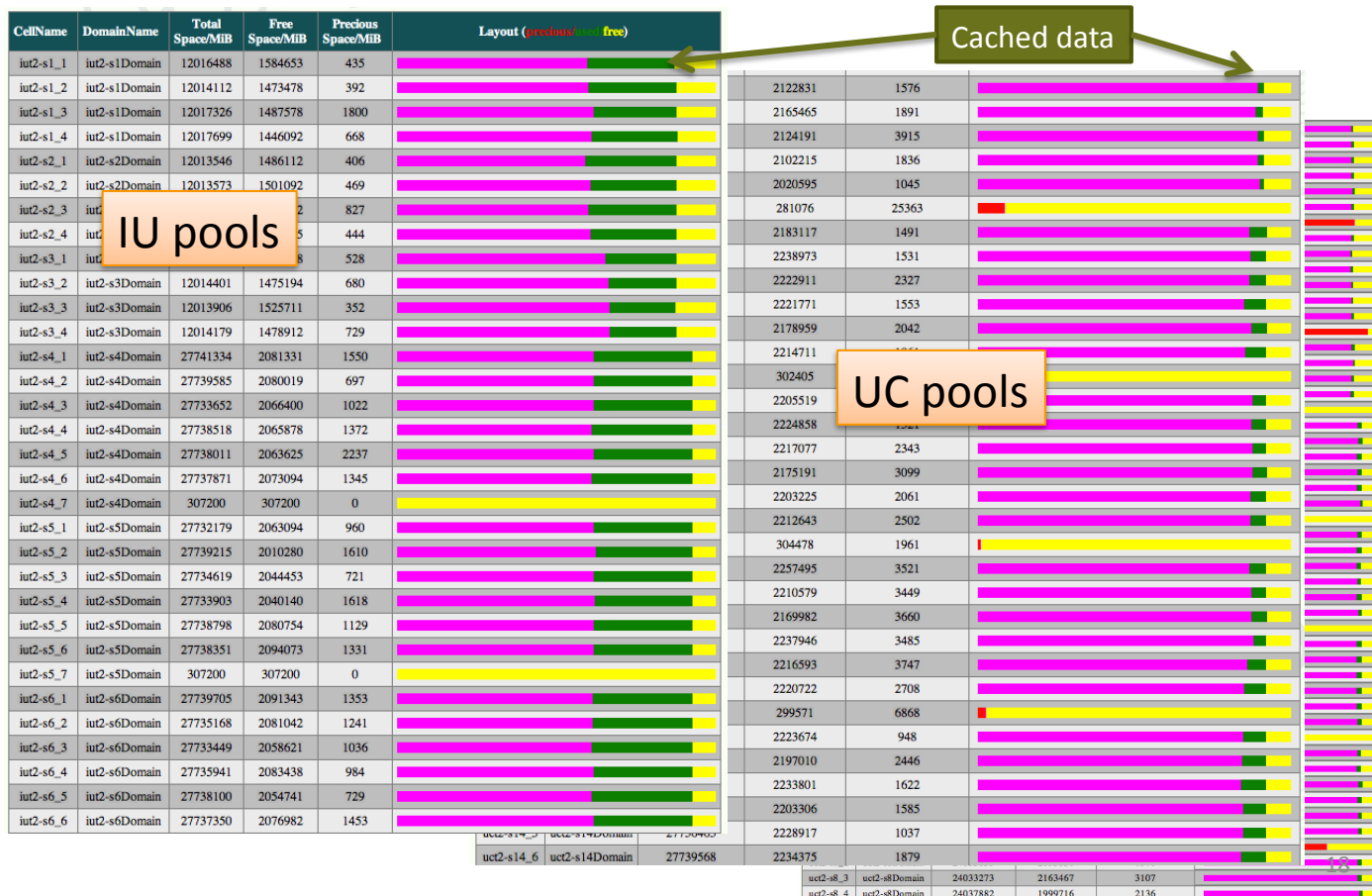
WAN Read Tests (basis for “cost matrix”)



from Ilija Vukotic

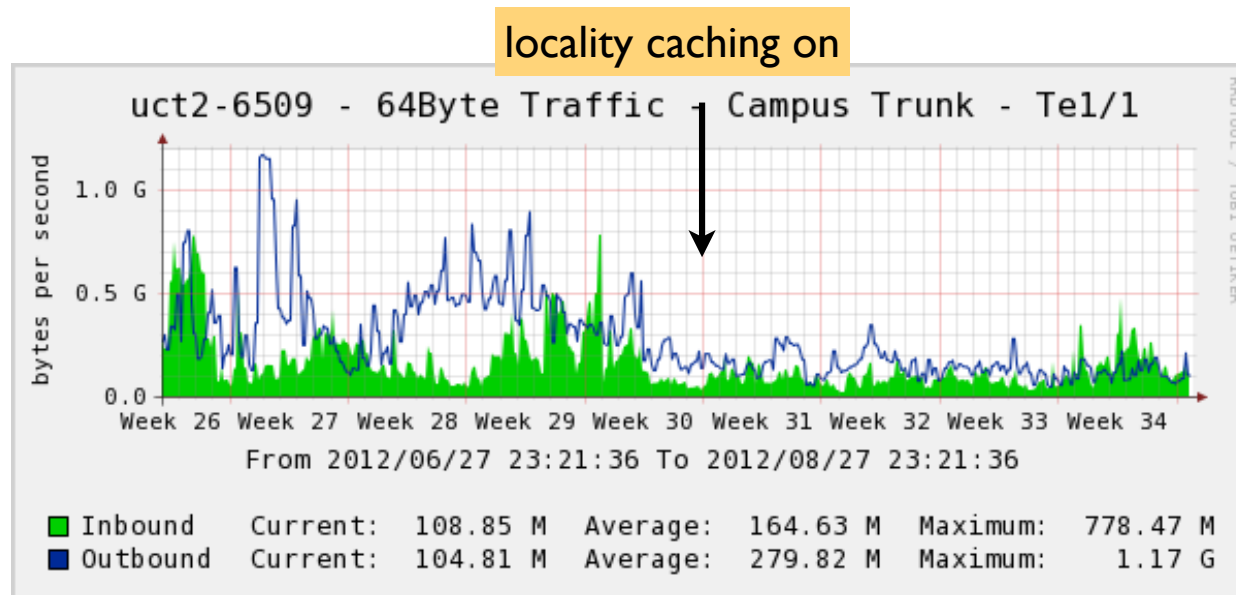
WAN direct access versus caching

- MWT2 federation is split between three sites separated by 5 ms (storage & cpu); storage at two of the sites.
- Recently configured dCache's 'locality' caching



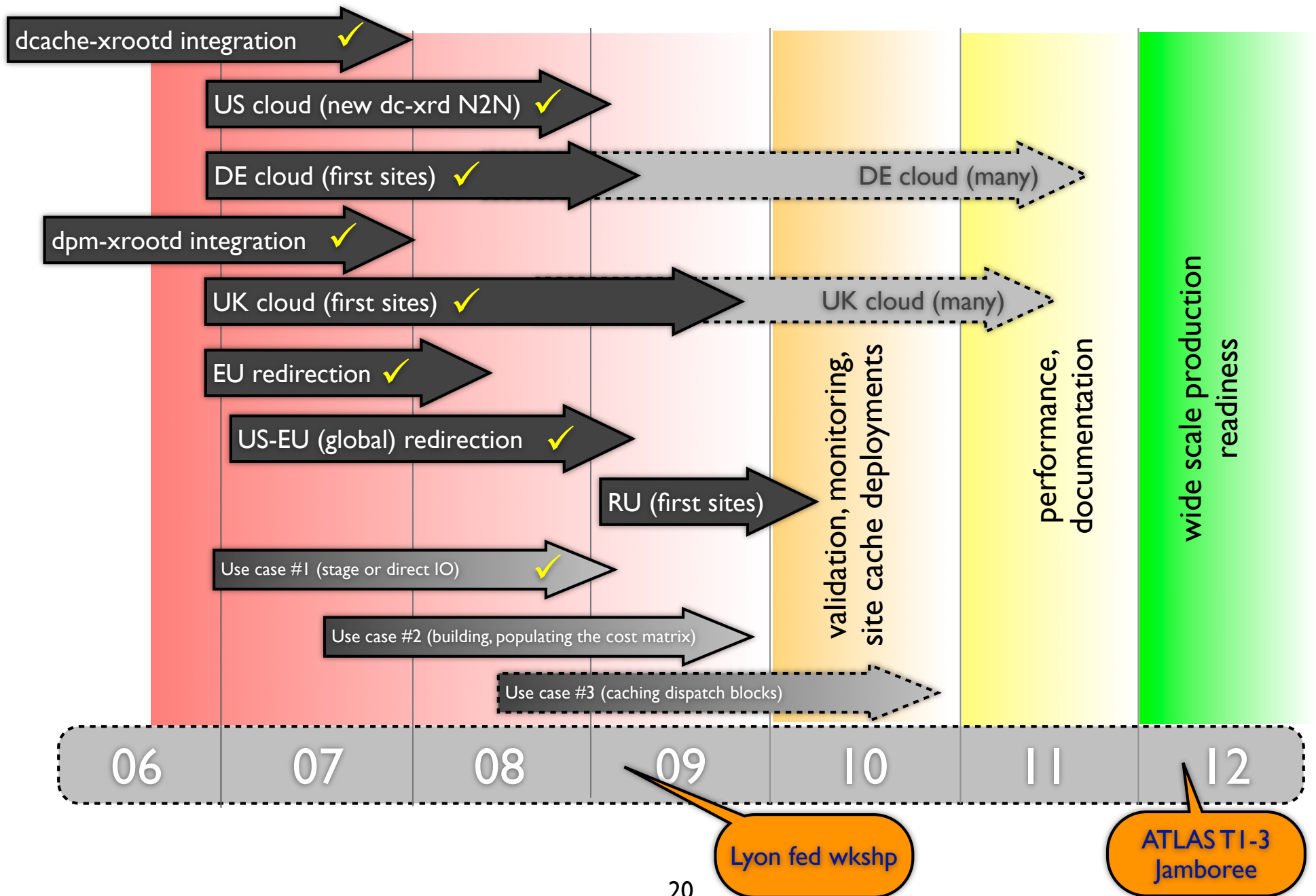
WAN direct access versus caching

- Expect WAN traffic reduced provided there is re-use and non-trivial fraction of file read
- Evidence of reduced WAN traffic
- Suggests equipping sites with federated caches, alongside managed storage



WAN
IO

ATLAS timeline



Summary

- Federated xrootd infrastructure for ATLAS now an ATLAS-wide project - presence in three clouds plus EOS
- Regional to global topology testing (up & down the trees)
- Accessible from Panda analysis jobs, either direct access or stage-in from pilot
- In parallel much work in ATLAS to optimize analysis code for wide area reads
- While **not yet in production** ATLAS is on a good course to provide new capabilities to users using federating Xrootd technologies