

# Storage federations, caches & WMS

Rob Gardner

Computation and Enrico Fermi Institutes  
University of Chicago

BigPanDA Workshop at CERN  
October 21, 2013





- Over the past few years the rigid T0/T1/T2/T3 hierarchy changing into a flatter mesh-like infrastructure
- Made possible by faster, reliable and affordable networks
  - & new virtual peering projects such as LHCONE
- Offers opportunity to create new ways of accessing data from production and analysis jobs
  - E.g. remove restriction that CPU and data located at same site
- Federation is a tool to allow jobs to flexibly reach datasets beyond the local site storage





- Explored for past two years as an ATLAS computing R&D project – Federated ATLAS XRootD (FAX)
- XRootD chosen as a mature technology widely used in HEP
- Efficient protocol for file discovery
- Provides a uniform interface to the various backend storage systems used in the WLCG (dCache, DPM are the most common)
- Close collaboration between XRootD and ROOT teams
- Good results in ALICE and CMS





- Create a common ATLAS namespace across all storage sites, accessible from anywhere
- Make easy to use, homogeneous access to data
- Identified initial use cases
  - **Failover** from stage-in problems with local storage
    - Now implemented, in production on several sites
  - Gain access to more CPUs using WAN direct read access
    - **Allow brokering to Tier 2s with partial datasets**
    - Opportunistic resources without local ATLAS storage
  - Use as caching mechanism at sites to reduce local data management tasks
    - Eliminate cataloging, consistency checking, deletion services



# A diversity of storage services



- ATLAS uses 80+ WLCG computing sites organized in Tiers
- Various storage technologies are used
  - dCache, DPM, Lustre, GPFS, Storm, XRootD and EOS
- Single source of deployment documentation: <https://twiki.cern.ch/twiki/bin/view/AtlasComputing/JoiningTheATLASFederation>
- To ease support we have experts for different technologies
- To ease communications we have contact persons per national cloud
- Name lookup is in transition from LFC to Rucio namespace ← expect performance impre

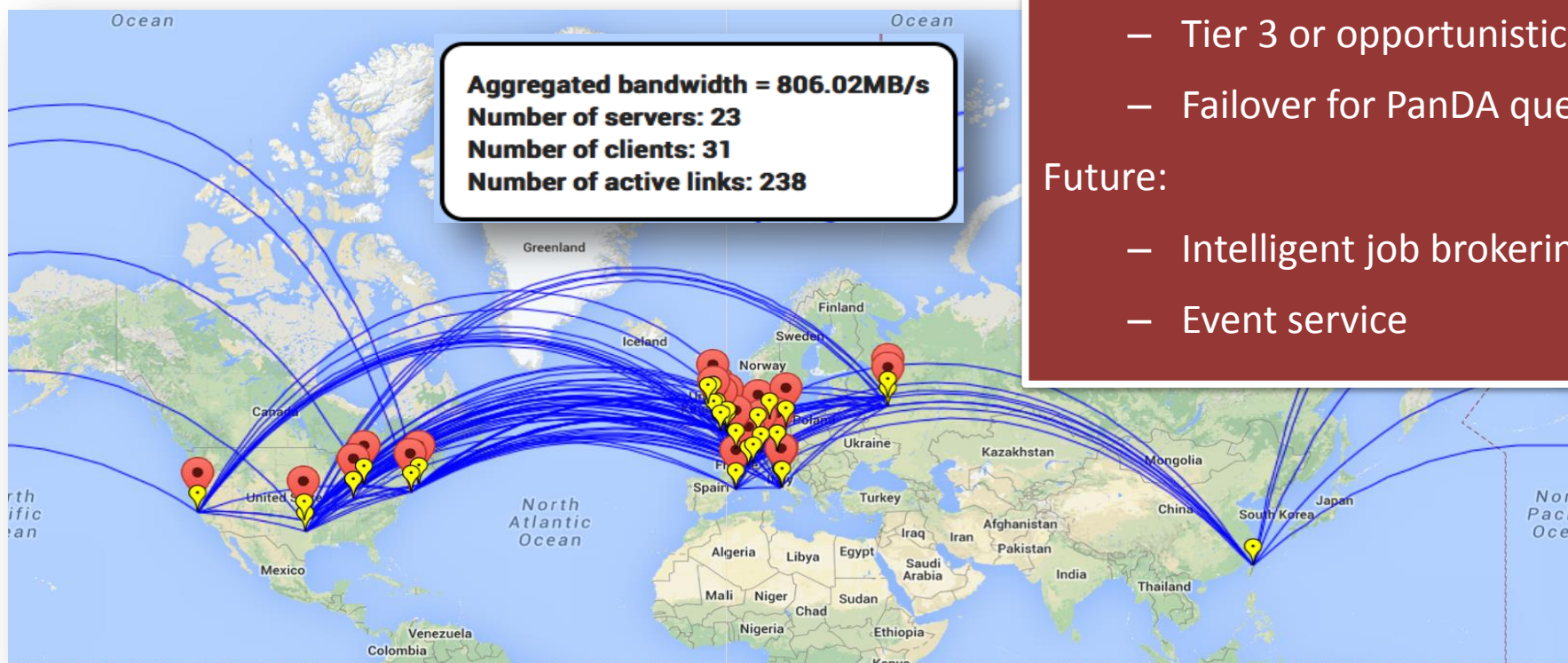


# A deployed federation



FAX (Federated ATLAS Xrootd) is a way to unify direct access to a diversity of storage services used by ATLAS

- Read only access to 177 PB
- Global namespace
- Currently 41 federated sites
- Regions covered: US, DE, UK, ES, and CERN



Main use cases today:

- Tier 3 or opportunistic CPU
- Failover for PanDA queues

Future:

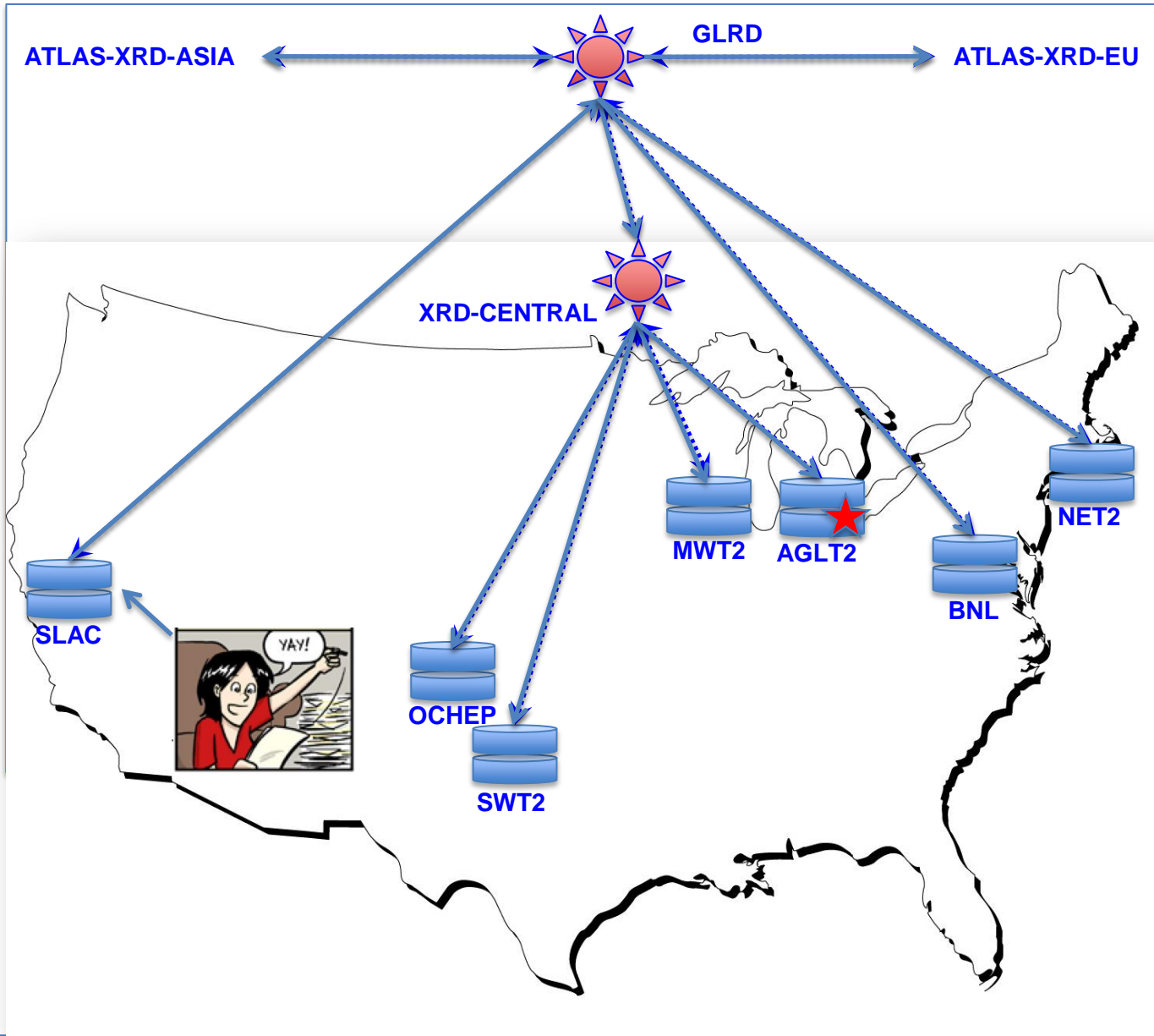
- Intelligent job brokering
- Event service

# File lookup sequence



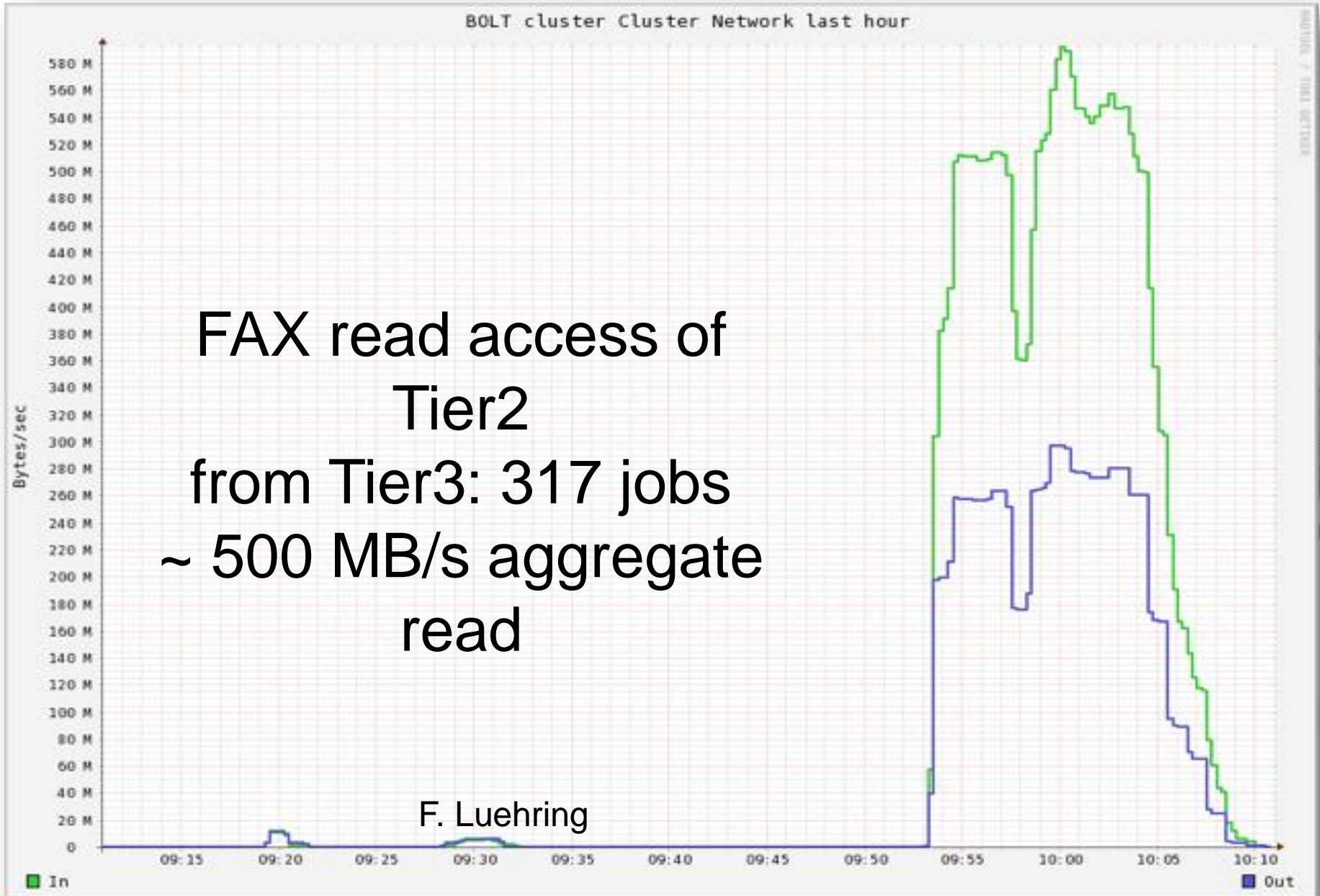
redirector

endpoint



- Data can be asked for from any endpoint and redirector
- Data are transferred directly from server to user
- Searching for the file is fast, delivering is more important
- Ideally one should use the one with best connection
- Usually that's the closest one

# User at Tier 3 reading from “nearby” Tier 2



ATLAS



# PanDA Job Failover using FAX



Currently PanDA only sends jobs to sites having complete datasets.

In case that an input file can not be obtained after 2 tries, the file will be obtained through FAX if it exists on another site.

There are in average 2.8 copies of files from recent reprocessing in FAX so there is a large change of success.



**PanDA Monitor**  
Times are in UTC

### Panda report on jobs failovers to FAX over last 24 hours

Record count: 138

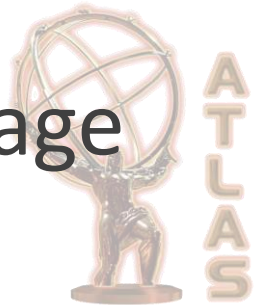
Show 50 entries Search:

	Site	Jobs	WithFAX [files]	WithoutFAX [files]	WithFAX [GB]	WithoutFAX [GB]		
+	DE: GoeGrid	1	1	19	0.17	2.15		
-	FR: ANALY_LPSC	1	1	1	0.15	0.06		
	<b>PandaID</b>	<b>Time</b>	<b>WithFAX</b>	<b>WithoutFAX</b>	<b>WithFAX [GB]</b>	<b>WithoutFAX [GB]</b>	<b>Status</b>	<b>User</b>
	<a href="#">1951899183</a>	2013-10-10 13:57:36	1	1	163463017	60679111	finished	mark hodgkinson
+	US: ANALY_MWT2_SL6	127	136	6428	52.68	1089.72		
+	US: OU_OCHEP_SWT2	9	9	99	5.38	38.39		

Showing 1 to 4 of 4 entries



- In production in US and UK for a few months
- Very low rates of usage (production storage elements already very efficient) and no major problems seen
- However, what happens if an unstable storage service fails all together?
  - Not yet seen, do a region's Tier 1 or neighboring Tier 2 become stressed?
  - Controlled testing planned



# FAX usage in Job Brokering



One can broker jobs to sites that don't have all or part of input data and use FAX to read them directly. Beneficial in cases where

- A site has very often full queues as some specific data exist only there
- A site has free CPUs, good connectivity but not enough storage/data
- One can use external sources of CPU cycles (OSG, Amazon/Google cloud based queues,...)

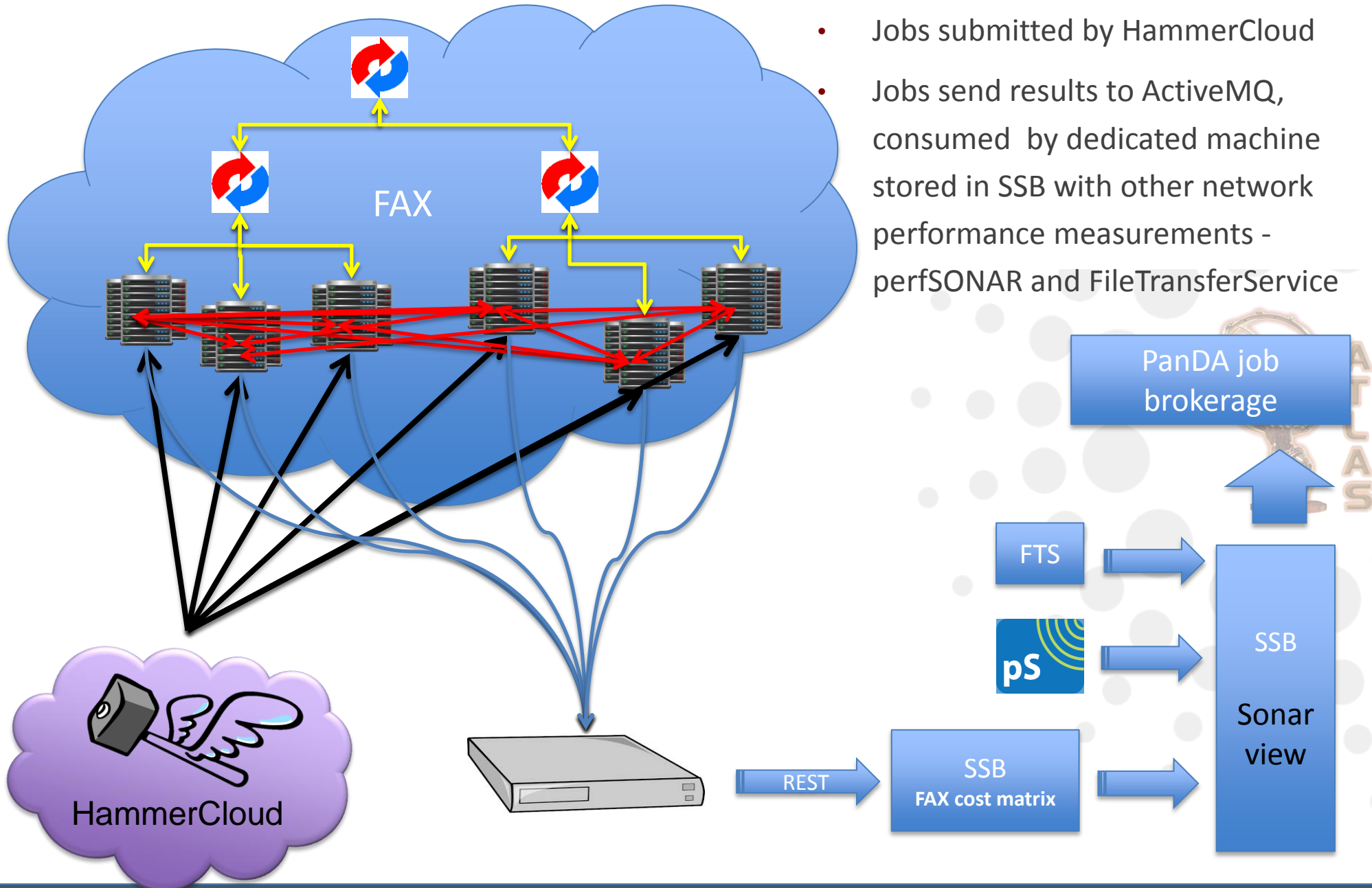


For this approach to be efficient, system has to “know” expected available bandwidth between a queue and all of the FAX endpoints

- **And control loads accordingly**
- Continuously collecting and storing data needed to do this (**Cost Matrix**).

# Cost matrix

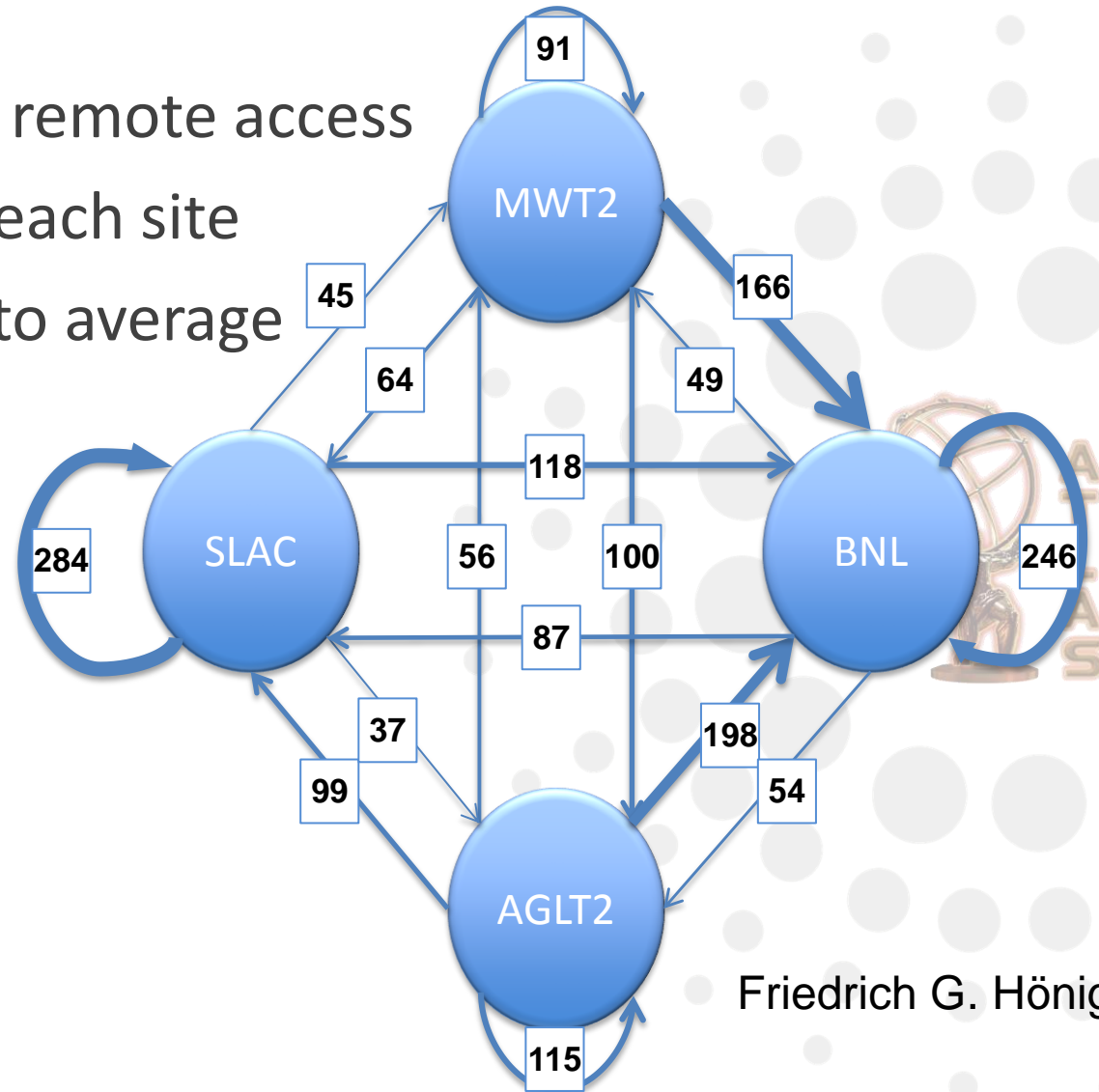
- Measures transfer rates (memory-to-memory) between 42 ANALY queues and each FAX endpoint
- Jobs submitted by HammerCloud
- Jobs send results to ActiveMQ, consumed by dedicated machine stored in SSB with other network performance measurements - perfSONAR and FileTransferService



# Simulate with WAN accesses



- HammerCloud based test
- Running real analysis code remote access
- **100 concurrent jobs** from each site
- Arrow width proportional to average event rate



WAN performance can be as good as LAN

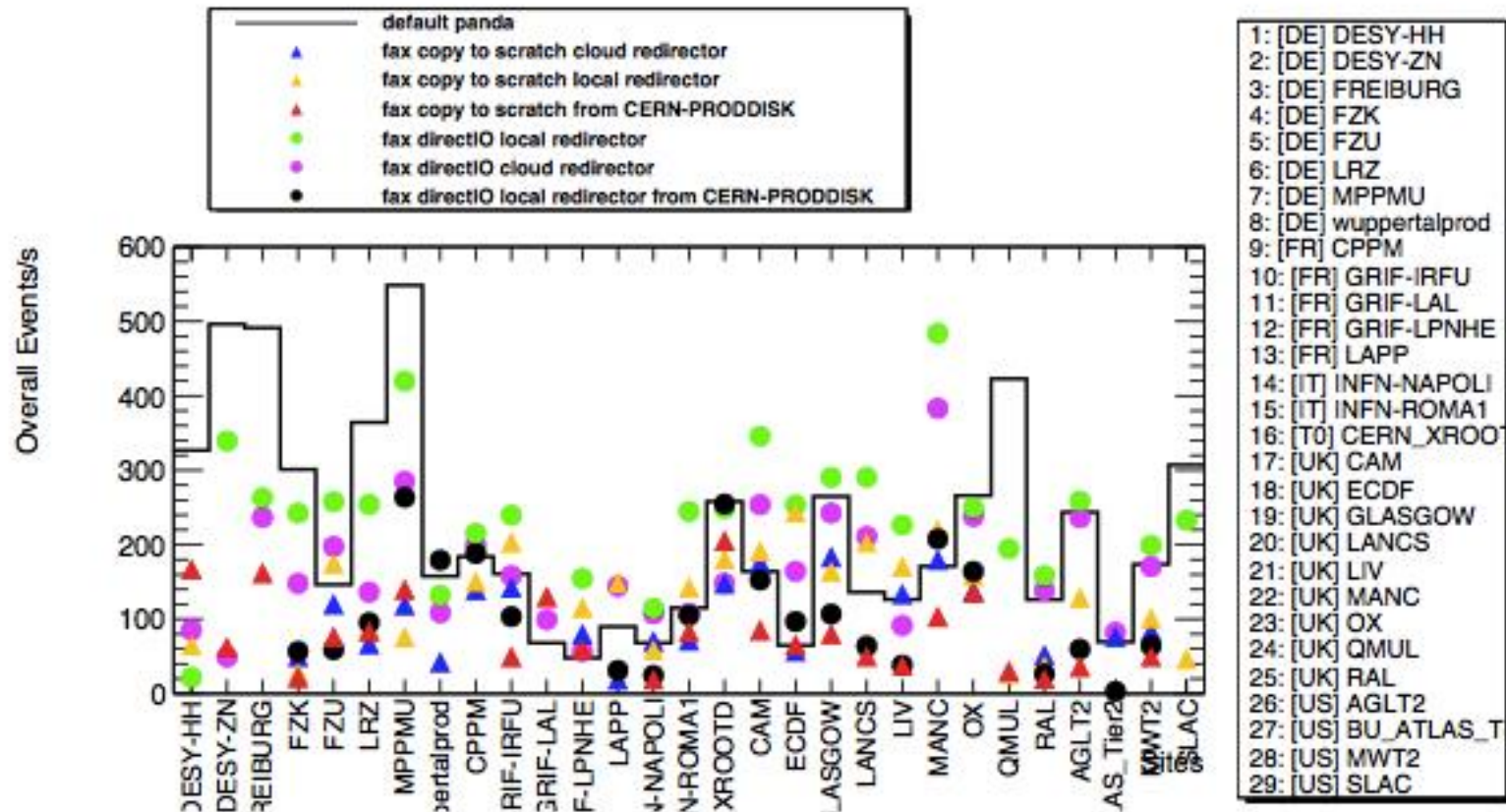
At this scale not unreasonable.

More testing on-going

Friedrich G. Hönig



## Complete Overview sorted by cloud

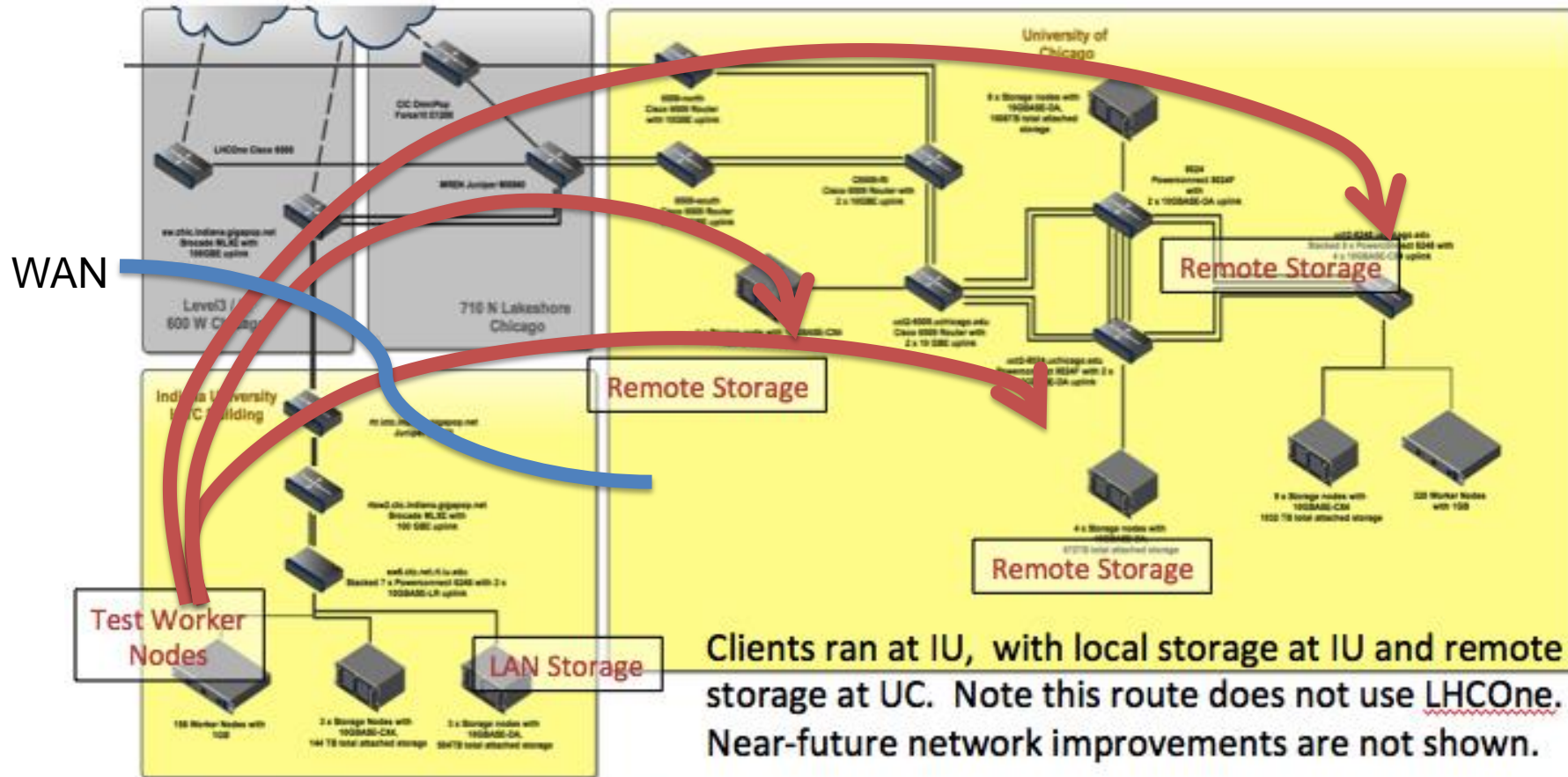


Friedrich G. Hönig

# Caching studies over WAN



Compare caching or direct access read over WAN with RTT ~ 3-5 ms



Clients ran at IU, with local storage at IU and remote storage at UC. Note this route does not use LHCOne. Near-future network improvements are not shown.

```
[root@iut2-c200 ~]# traceroute uct2-s20.uchicago.edu
traceroute to uct2-s20.uchicago.edu (128.135.158.170), 30 hops max, 60 byte packets
 1 149.165.225.254 (149.165.225.254) 20.328 ms 20.325 ms 9.372 ms
 2 et-10-0-0.2012.rtr.ictc.indiana.gigapop.net (149.165.254.249) 0.448 ms 0.491 ms 0.484 ms
 3 149.165.227.22 (149.165.227.22) 5.087 ms 5.196 ms 5.273 ms
 4 10.4.247.230 (10.4.247.230) 5.229 ms 5.396 ms 5.547 ms
 5 uct2-s20.uchicago.edu (128.135.158.170) 5.323 ms 5.347 ms 5.220 ms
```

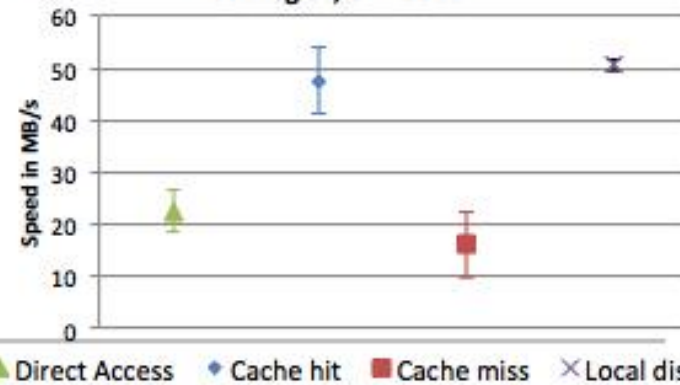




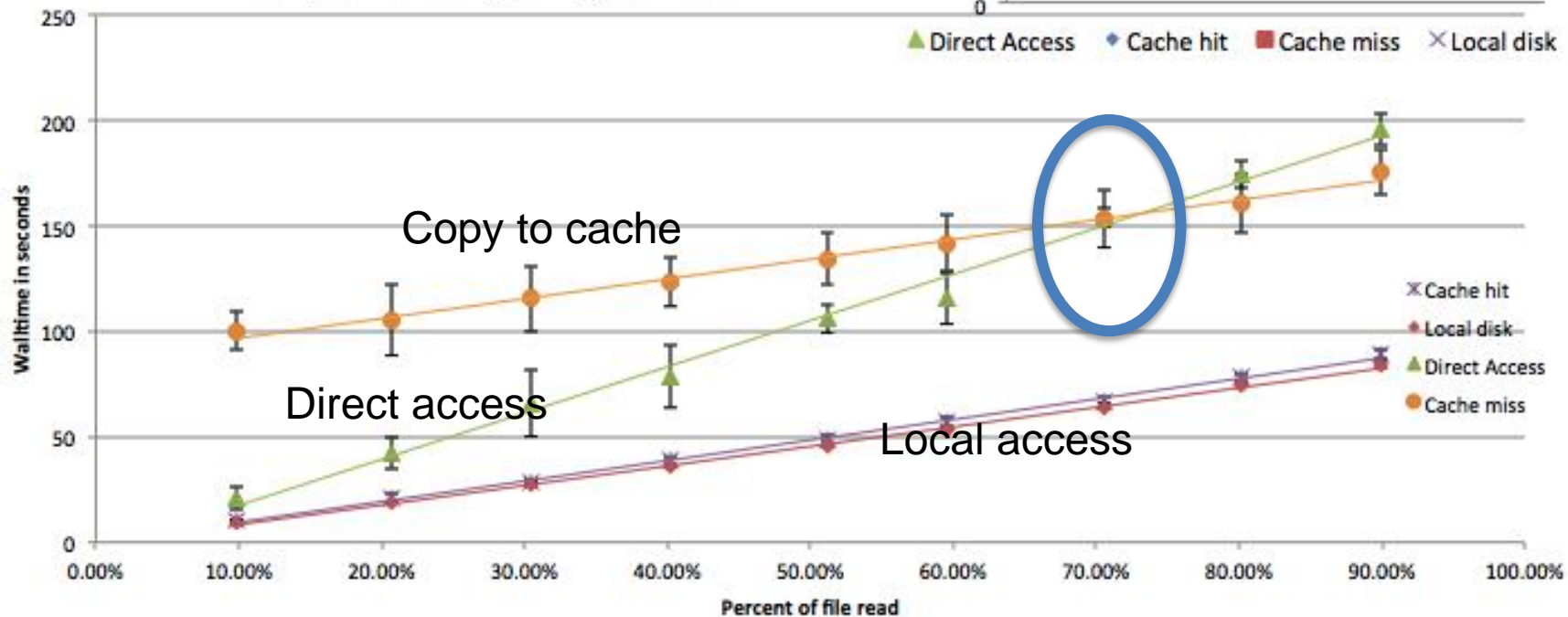
## Direct comparison of caching strategies

No caching is the optimum strategy when less than 75% of the file is read. When more than 75% is read, caching becomes optimal.

Comparison of Download Speed of Caching Strategies, on X5660



Comparison of Caching Strategies on X5660





# The other meaning of “Cache” (in Rucio)



From Vincent Garrone

- A cache is storage service which keeps additional copies of files to reduce response time and bandwidth usage.
- In Rucio, a cache is an RSE (Rucio Storage Element), tagged as volatile.
- The control of the cache content is usually handled by an external process or applications (e.g. Panda) and not by Rucio.
- The information about replica location on volatile RSEs can have a lifetime.
- Replicas registered on volatile RSEs **are excluded from the Rucio replica management system** (replication rules, quota, replication locks).



*There will be much celebration by site admins when this is realized!*



- Broadly deployed, functional Xrootd federation covering vast portions of ATLAS data
- First use cases are implemented – Tier 3 access, PanDA failover from production queues
- Inputs for intelligent re-brokering of jobs being collected
- Caching studies, and development of caching services – still to do

