

Data Access DE Cloud

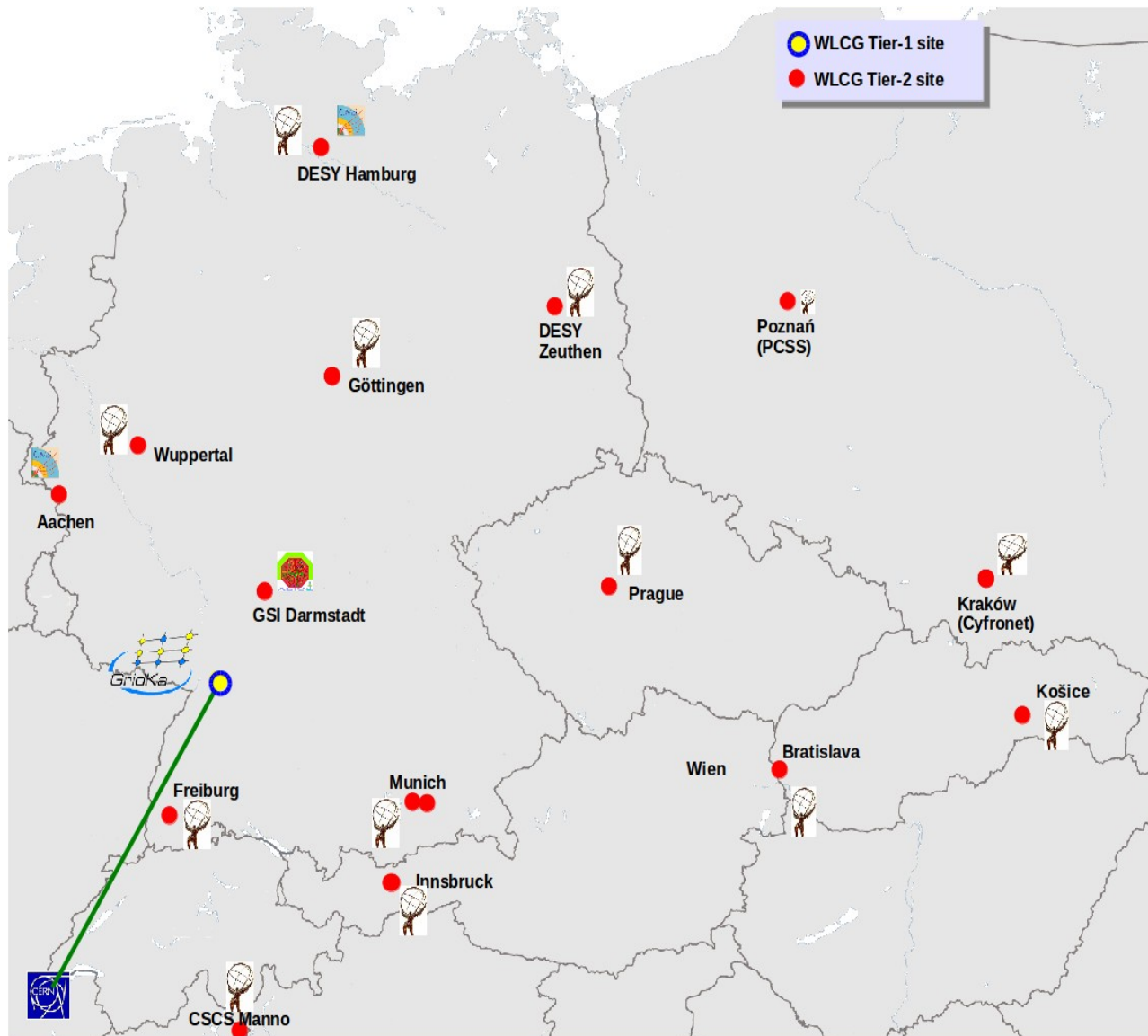
pre-GDB, Cern
May 13, 2014

- Overview
- Local Access
 - Experience with direct IO
- WAN access

Note:

Mainly presenting ATLAS perspective plus a few slides from CMS

ATLAS DE Overview



> 15 Sites in
“ATLAS GridKa cloud”

- in DE:
 - 8 dCache sites (1 T1, 7 T2, 1 T3)
 - several T3 w/ Lustre/GPFS/Sonas
- outside DE:
 - DPM sites Prague, Cracow, Innsbruck
 - dCache in CSCS/Manno

Networking & T2/T3s

- GridKa KIT connected via
 - LHCOPN (multiple 10 Gb)
 - LHCONE (10 Gb) – Desy, Wuppertal
 - X-Win (10 Gb) – Munich, Freiburg, Goettingen
 - dedicated links to Prague (10 Gb) and PL (1 Gb)
- Tier-2/3 sites
 - large T2 sites at Desy-HH, Desy-ZN, Freiburg, Goettingen, Wuppertal, Munich (LRZ-LMU & MPPMU), Prague/CZ, Cracow/PL, Lugano/CH
 - most with substantial T3 add-on (+100% CPU, ~3 PB LOCALGROUPDISK)
 - a few smaller T2s in PL, Austria, Slovakia (since 2012)
 - T3 sites with opportunistic use for ATLAS production
 - Dortmund, Dresden, Mainz

Local Data Access ATLAS DE

- T1/2 Sites in DE all use dCache as SE
 - evaluated early on (~2008) use of direct IO from dCache SE to local WN cluster instead of default copy-to-scratch
 - using native dcap protocol
 - required some development & optimizations in dcap, Root-IO, ATLAS file layout, etc, to reach stable operation and good performance
 - used by default for ATLAS analysis jobs on DE sites since many years

Direct IO vs copy-to-scratch

- Pros:

- no need to download data to /tmp
 - often only fraction of data read (5-20%)
 - data processing volume not limited by /tmp space
- usually higher event processing rate
 - needs efficient caching algorithms (e.g. optimized TTreeCache)

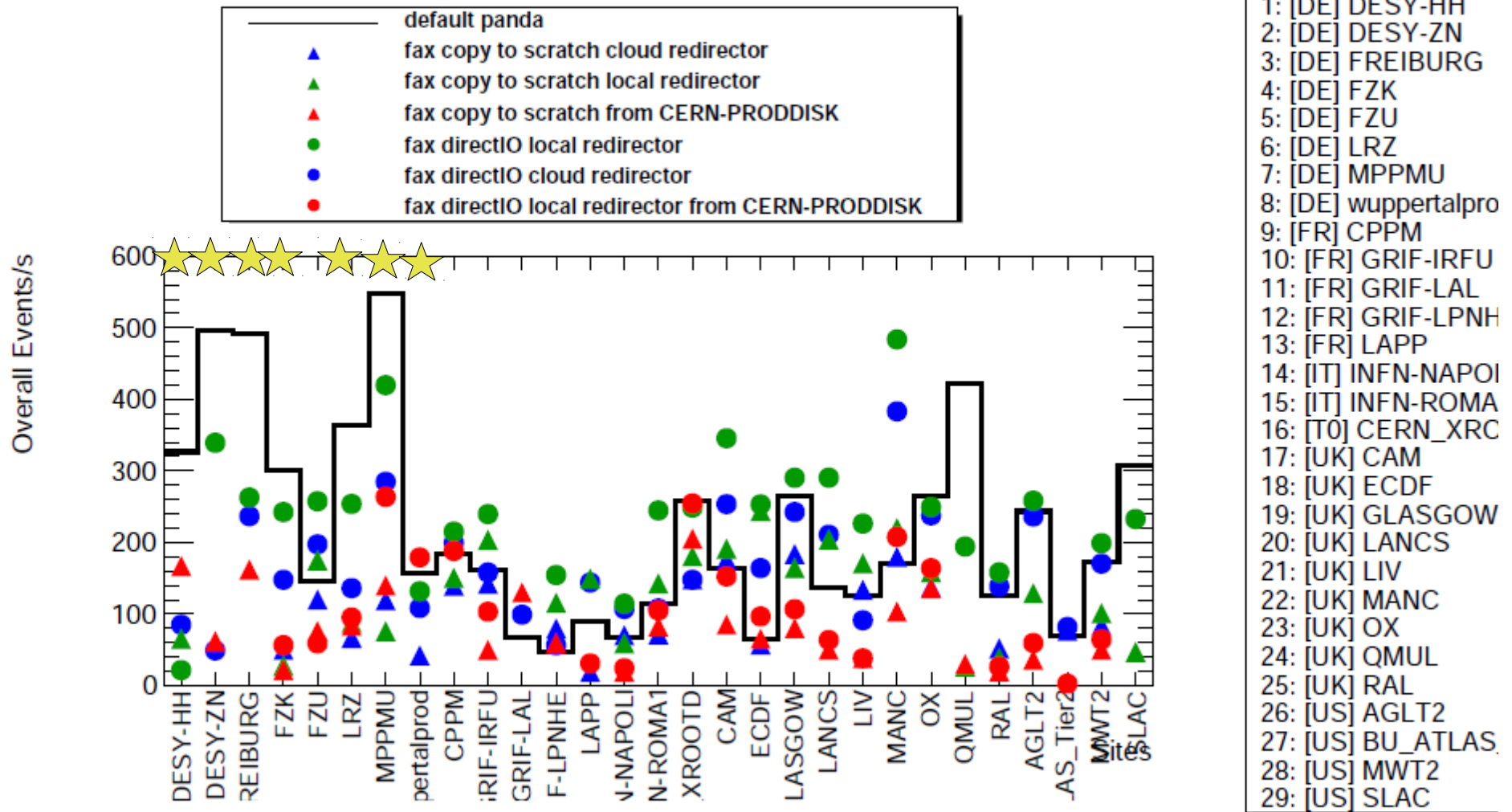
- Cons

- random IO operation can cause high load for large storage servers
- stage-in mode easier to control, check, re-try, fall-back
 - handled in job-script not via Root-IO

IO performance test

- ATLAS is doing systematic tests of IO performance for FAX using the Hammercloud system
 - measures not only remote access but also default local access mode as configured in ATLAS Panda for many different sites
 - Test case:
 - typical ATLAS ntuple analysis H → WW (Root based, no Athena)
 - PandDA default data access as baseline
 - directIO vs copy-to-scratch
 - local fax redirector vs cloud redirector
 - TTreeCache activated for all variables
 - Metrics shown:
 - Event rate relative to whole job (not just payload...)
- Tests done by F. Hoenig (and HC team)

Complete Overview sorted by cloud



★ dCache direct IO

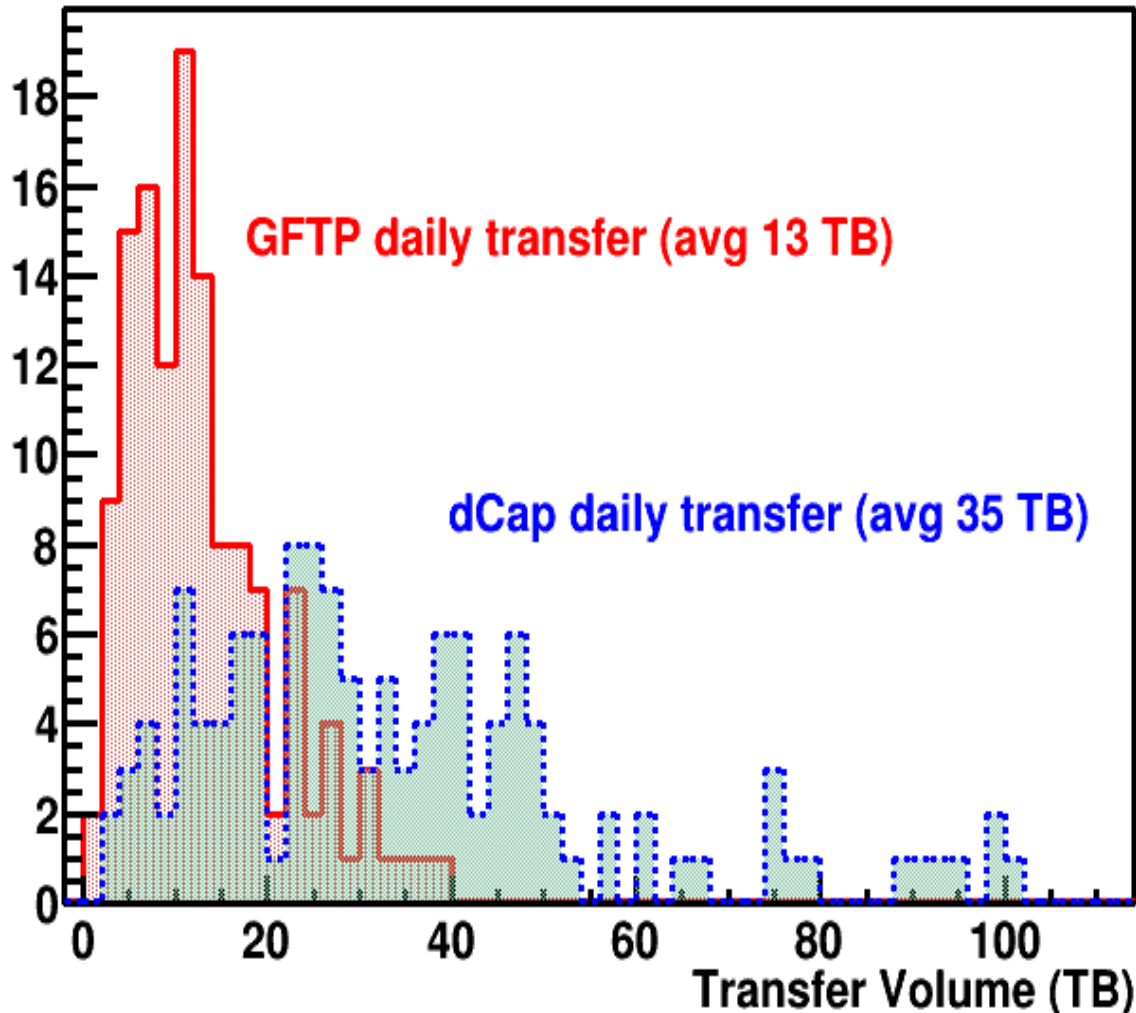
Aside: Direct IO protocol

- dcap :
 - Default and mostly used at ATLAS-DE sites
 - supports tunable caching schemes & vector-read
 - dCache-only, limited support
- xroot :
 - HEP std, supported by xrootd, dpm, EOS, dCache
 - considered move from dcap → xroot protocol for ATLAS-DE
 - tested on 1 site (LRZ-LMU) for ~2 month
 - suffered from low-rate (5-10%) random job-failures, not reproducible
 - ticket open w/ dCache, involved xrootd experts, still unresolved
- NFS 4.1
 - Used in production at Desy for some non-LHC VOs
 - Issues to get NFS4.1 supported on client side
- [http/webdav/davix](http://webdav/davix):

ATLAS DE federated data access

- In general decent network connectivity between DE sites via GPN X-WIN, distance moderate ($O(\text{few } 100 \text{ km})$)
- Special cases for federations:
 - Desy-HH \leftrightarrow Desy-ZN, dedicated 10/20 Gb link
 - LRZ-LMU \leftrightarrow MPPMU, only 500 m apart, dedicated 10 Gb
 - Sites (Desy-*, LRZ, Wuppertal, Bonn) with fast local file system (Lustre, GPFS, Sonas, ...), investigate use as cache
- [xrootd/FAX](#)
 - for dCache setup rather straightforward \rightarrow details see next slide
- [http/webdav](#)
 - used at some sites for prod input download (R.Walker, aria2c & Meta-links)
 - investigating direct IO w/ Davix

Local vs WAN IO



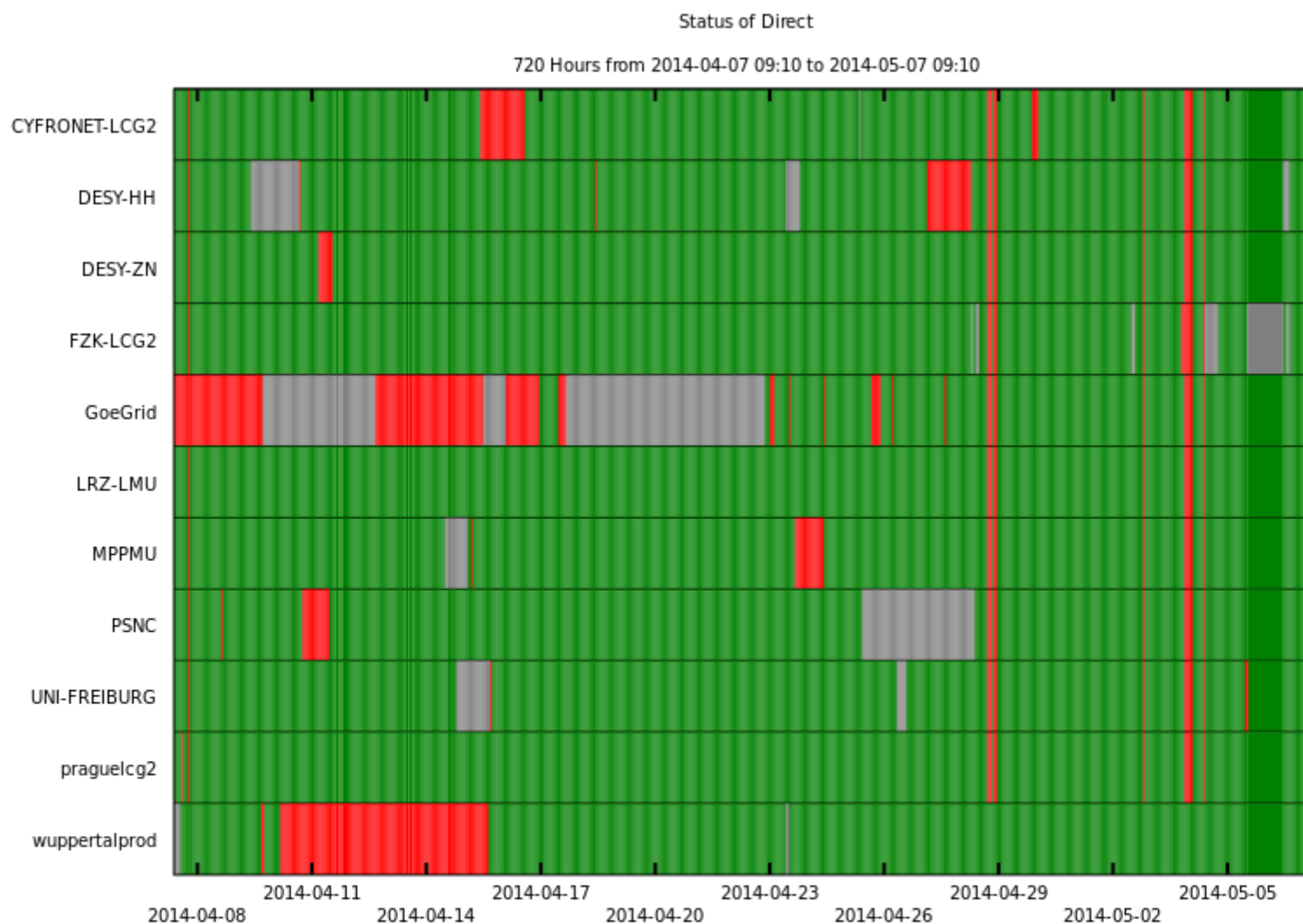
- Example of LRZ-LMU T2:
 - daily transfer volume Jan-Apr 2014
 - local transfer in avg factor 3 larger
 - local peak local volume >150 TB/day
 - spikes > 6 GB/s
 - cf WAN bandwidth 6 Gbit/s
- Room for addtl WAN direct IO at level of 10-30% wrt LAN IO

FAX deployment in DE cloud

- two main options how to integrate existing dCache site
 - Plugin for dCache xrootd door (re-direction and N2N service)
 - dependencies on details of dCache version
 - rather invasive setup (billing DB access for monitoring)
 - in use at GridKa (and US T2s)
 - redirects connection to pool → good scaling
 - Proxy xrootd setup on extra node
 - can talk to dCache xrootd door or Posix-NFS mounted storage
 - non-invasive setup, just talks to dCache xrootd door (as any client job)
 - all traffic between outside world and storage via this node (no redirection to pools)
 - implicitly limits & protects WAN connection of SE site
 - a really crucial requirement for sites with shared WAN access
 - used in DE at all sites except GridKa T1

FAX status in DE

- FAX status – mostly ok
 - all larger sites deployed service except CSCS/CH

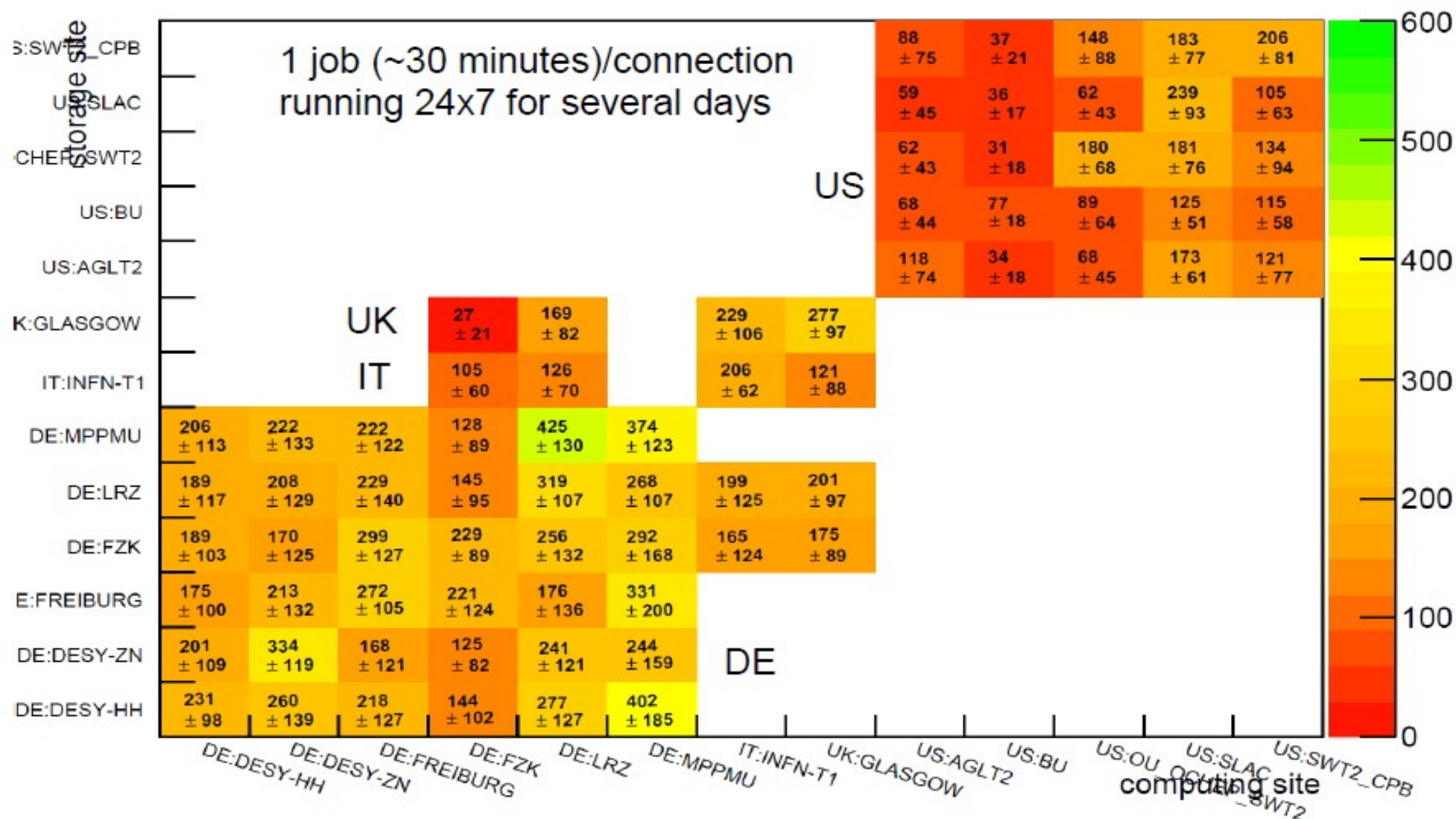


FAX tests

- detailed Hammercloud testing of FAX going on (world-wide)
 - replicated dedicated unique datasets per sites to evaluate remote access
 - execute typical ATLAS ntuple analysis
- move away from LFC based file-lookup to deterministic Rucio name scheme gave substantial improvement in performance and stability
- next two slides show examples of tests with low load
 - only few jobs running at a time per site
 - very decent performance and stability
- plan to do further load tests with increasing # parallel jobs to determine limitations

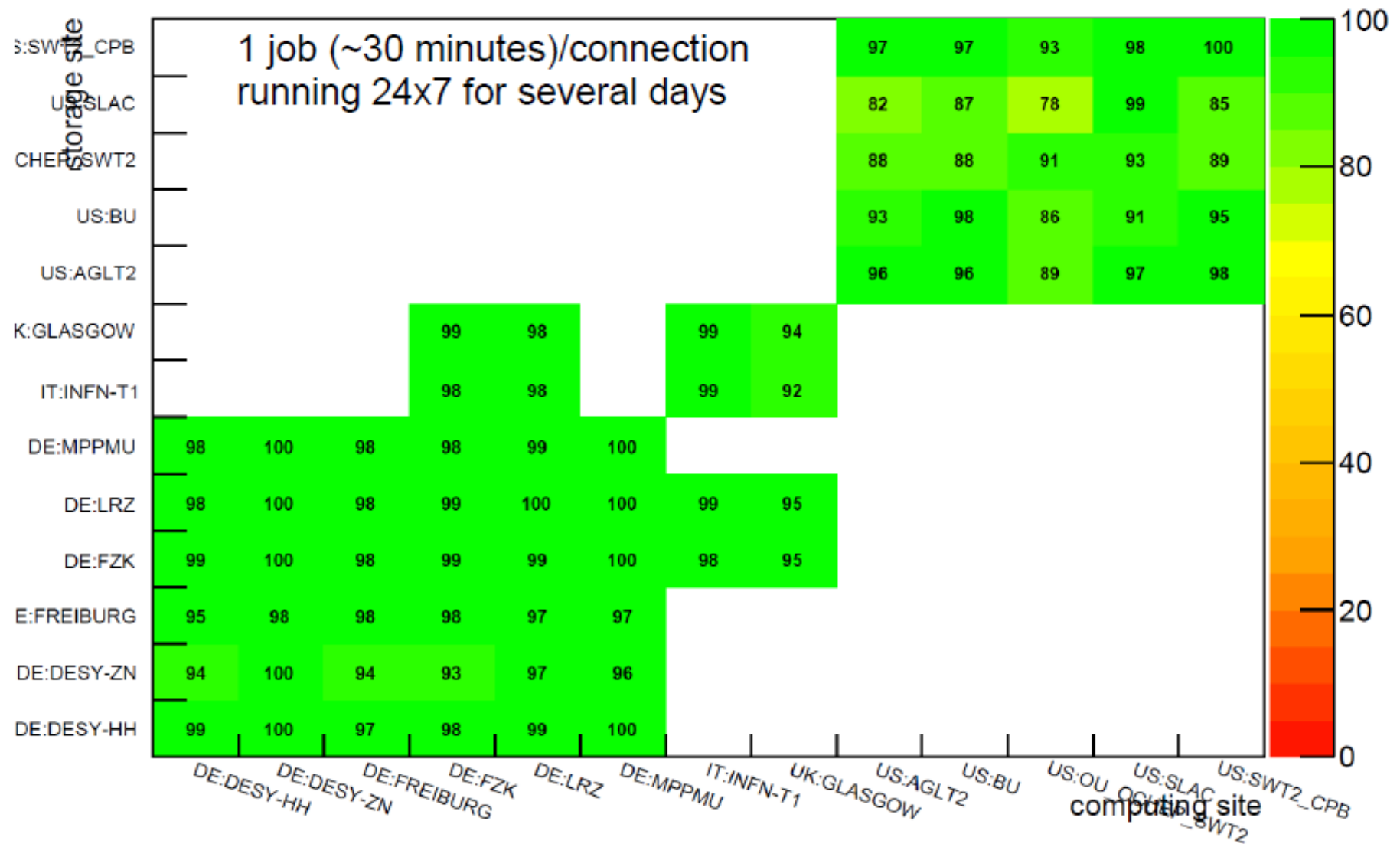
HC FAX stress test results: far-away tests

mean overall eventrate per job [1/s]



HC FAX stress test results: far-away tests

percent job success



Remote IO – wishlist

- load tests
 - need to understand how many jobs we can run per site/connection
- monitoring to quickly identify FAX activity in case of problems
 - several cases of WAN/firewall overload reported
 - caused by other VOs (AFAIK)
 - took days/weeks to investigate
- clarify use-cases
 - e.g. for direct IO fall-back not yet possible in practice
 - combine remote IO and local caching
 - many potential sites with O(100 TB) Lustre/GPFS
 - too small as T2 storage but promising as cache
- evaluate full-featured http/davix direct IO

- > Two use-cases
 - > 1.) “Fallback”
 - Try to open file locally
 - In case of local open failure ask regional xrootd-redirector
 - > 2.) Join site storage into federation
 - Publish local files into regional redirector
 - Deliver files to remote clients
- > Both deployed at German Tier-1 and Tier-2 sites
 - T1_DE_KIT
 - T2_DE_DESY
 - T2_DE_RWTH
 - All use dCache based storage systems

> Basically routine operation

- Actively used
 - > Some to many TB per week [DESY example]
 - > cf: LAN IO >1000 TB/wk
Phedex several 10 TB/wk

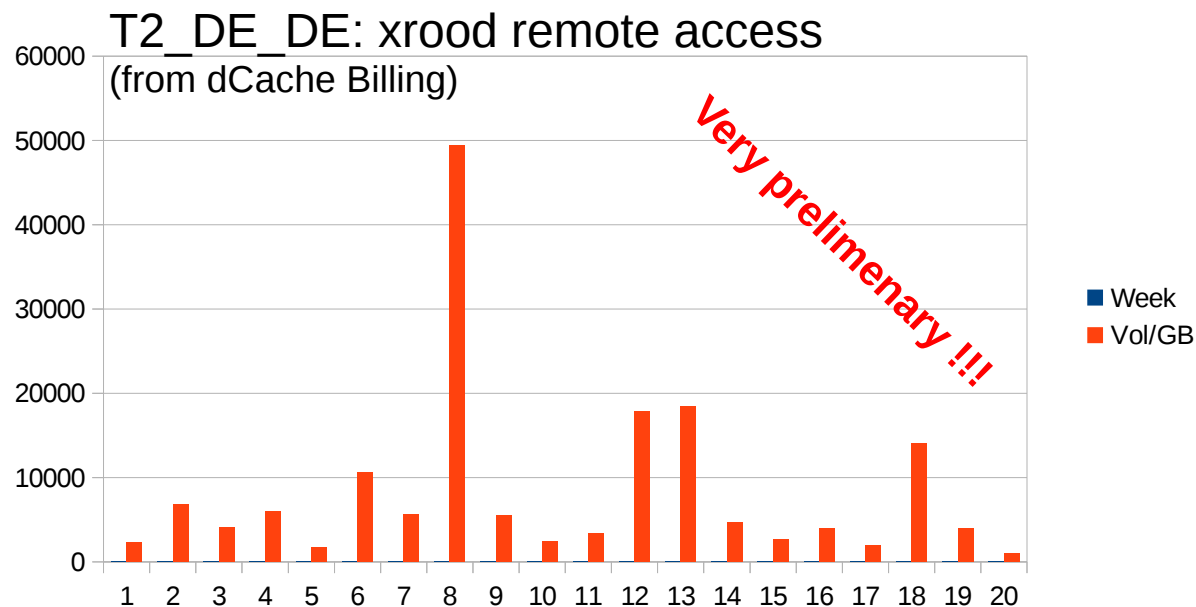
- So far not much trouble

> Documentation by CMS community

- Quite useful

> Issues and concerns

- Protection against overloading the WAN by
 - > Too many Fallback file openings
 - > Too many remote access – need advice for dCache mover tuning
- Support of “detailed xrootd monitoring plugin”
 - > Not part of dCache
 - > Required on every storage pool



Slide C. Wissing
Desy

Summary

- Good experience with LAN direct IO at ATLAS-De sites since many years
- reducing protocol zoo and/or use of common std desirable
 - not easy to achieve in practice
- WAN/remote IO
 - FAX/xrootd largely deployed at DE sites
 - performance and stability looks promising
 - http/Webdav/Davix
 - in use for simulation input download (aria2c) at few sites
 - still testing for analysis direct IO
- CMS:
 - AAA in routine use at CMS DE sites