# Open Science Grid Storage Workshop

Frank Würthwein (UCSD)

# Outline

- Key drivers of OSG for the next year
- Go through, and try to verbalize what these drivers mean for storage on OSG.

# The top 10: Key Goals for FY10

Top 5:

- Support for LHC user running, support for Tier-1, Tier-2s and focus on Tier-3s.

- Support software, analyses, and data management for LIGO based on requests.

- Easier and more usable incremental upgrades of software for "any" reason – security, faults, functionality.

- Support for other OSG stakeholders requests and increased number of non-physics and campus beneficiaries in OSG.

- Timely and appropriate security and operational response to problems, as well as records of expectations, including SLAs and Policy, across the board.

Next 5:

- Continued improved technical and operational support for data storage use, access and policies.

- Wider adoption of Pilot based workload management, progress in transparency between campus & wide area resources, policies for improved usability & efficiency.

- Articulation and implementation of Security authorization, identification & policy.

- Success of OSG Satellites for moving the infrastructure forward

- Better understanding of role of (separately) MPI, Virtual Machines and Cloud resources and policies in the OSG environment.
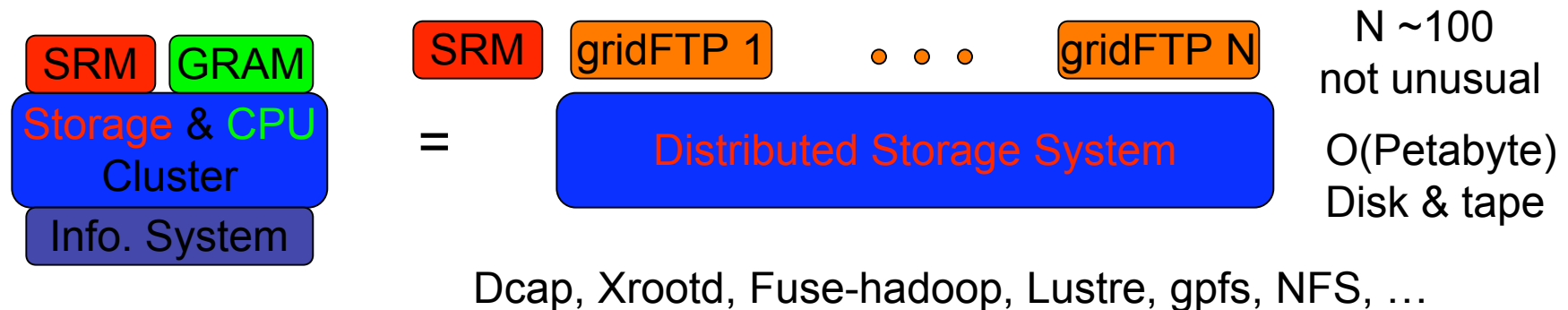
# Storage for LHC

- Basic Model
- Setting the scale
- What works, and what doesn't

# Basic Model for LHC

Uniform WAN access protocols …
… but multiple implementations.

SRM  GRAM

Storage & CPU
Cluster

Info. System

= 

SRM   gridFTP 1   • • •   gridFTP N

Distributed Storage System

N ~100
not unusual

O(Petabyte)
Disk & tape

Dcap, Xrootd, Fuse-hadoop, Lustre, gpfs, NFS, …

Heterogeneous LAN access

Applications designed to transparently access LAN
via multiple protocols. Specific protocol choice is
configured at each site via application software installation.
Generally, only restricted subset of POSIX available.

5

# "Storage Workflow"

- Install application software at sites
- Stage dataset(s) to sites
  - Staging requested by scientific community
  - Requests accepted/denied by site data manager
- Send jobs to site that has the needed input data
  - Output is written to local disk while job is running.
- Stage out output data to "home" Tier-2
  - Register staged out data in filecatalog
    - This produces new dataset that can be either processed locally, or "elevated" to public status, archived to tape, and moved around.

# Setting the Scale in 2009/10

- Tier-2:
  - 200TB to 600TB
  - O(1 Million) files
  - 10-100Hz WAN access to SRM.
  - 1,000 to 10,000 cores accessing storage simultaneously.
  - Want Gbps wirespeed performance on LAN from 100-1000 nodes.
  - Want moderate maintenance cost, < 0.5FTE .

- Tier-3:
  - 20TB to 200TB
  - O(100k) files
  - O(10Hz) WAN access to SRM.
  - 100 - 1000 cores accessing storage simultaneously.
  - Want Gbps wirespeed performance on LAN from 10-100 nodes.
  - Want low maintenance cost, O(1/10th) FTE

*These are stretch goals,*
*maybe not all can be accomplished at once.*

# Setting the Scale (II)

- <span style="color:red">SRM: scale is driven by stage-out. Would want 50-100Hz for both srmls and srmcp by end of 2010.</span>

- WAN IO: want to be able to fill >50% of a 10Gbps network pipe with hundreds of parallel flows.

- LAN IO to worker nodes: ~ 10Mbps per core. Gbps LAN is sufficient for some time to come.
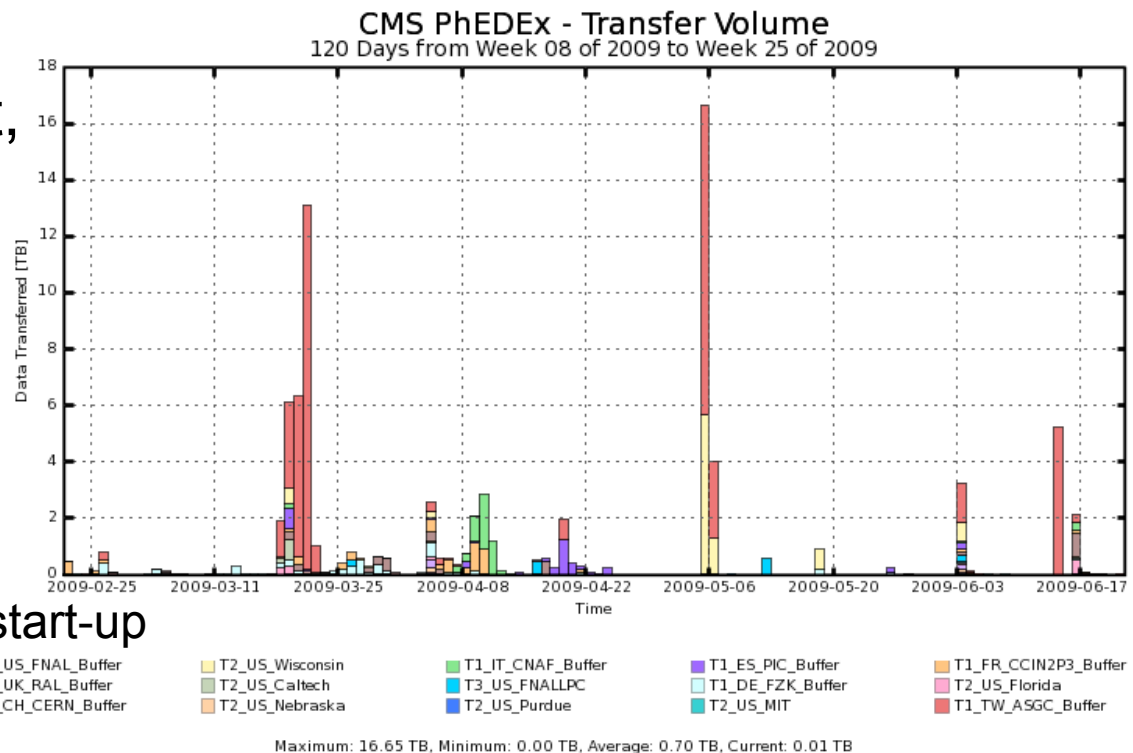
# What works:

- Organized and scheduled WAN data movement is a solved problem for 09/10.

- Larger WAN IO capabilities needed by ~11/12.

We want a 10TB dataset, and we get it within a couple days or so.

It's not always perfect, but it's quite useable:

10/pb of e/mu RECO ~ 20TB at start-up



CMS PhEDEx - Transfer Volume
120 Days from Week 08 of 2009 to Week 25 of 2009

Maximum: 16.65 TB, Minimum: 0.00 TB, Average: 0.70 TB, Current: 0.01 TB

# What doesn't work:

- ## Reliability of read access:
  - E.g. CMS ran 0.5 Million few hours test jobs during a 2 week period in June 2009 at 80% success rate.
    - Test jobs had no stage-out.
    - Roughly 90% of all failures of jobs were due to LAN read errors.
- ## Safety of storage on disk:
  - Disks go bad routinely, and T2s have no backup.
  - Without replication, user files get lost.
- ## Scalability of SRM:
  - CMS stages out all user generated files via the WAN to the Tier-2/3 that hosts the user's disk space.
    - The rate SRM is required to support thus scales with the amount of resources a single user can get across the grid.

# LIGO

- For next year, LIGO wants to move beyond "Einstein@Home" use of OSG.

⇒Operate a storage Element @ home for stage-out of results, and as source for stage-in at opportunistic sites.

⇒Incorporate opportunistic storage into LIGO workflows.

⇒Understand how to best use the mix of local storage and SE storage.

- And possibly other goals I'm unaware of.

# OSG Community in general

- Modest tools for space reservation
- Soon tools for storage discovery
- <span style="color:red">No tools in sight for managed WAN data movements in/out of sites.</span>
- Not clear if OSG support of storage is good enough for people outside LHC/LIGO
- Most of storage on OSG is not fully posix compliant. Workload tools need to be aware of the limitations.

# Continued Support from OSG

- Work with Storage Software providers on:
    - Reducing cost of operations
        - Increasing reliability
        - Reducing complexity
    - Increasing scalability
    - Improving ease-of-use

*How might OSG achieve this?*

*Testing, packaging, deployment & operations support, …*

Probably more of the same we've been doing so far.

# Summary

- I hope this brief presentation was useful to put things in perspective.

- In the end, our goal for the next 1.5 days is to provide a forum where site admins and developers can exchange information among and between each other.

- And if this leads to us rethinking the priorities for the next year, then this is not a bad thing.