

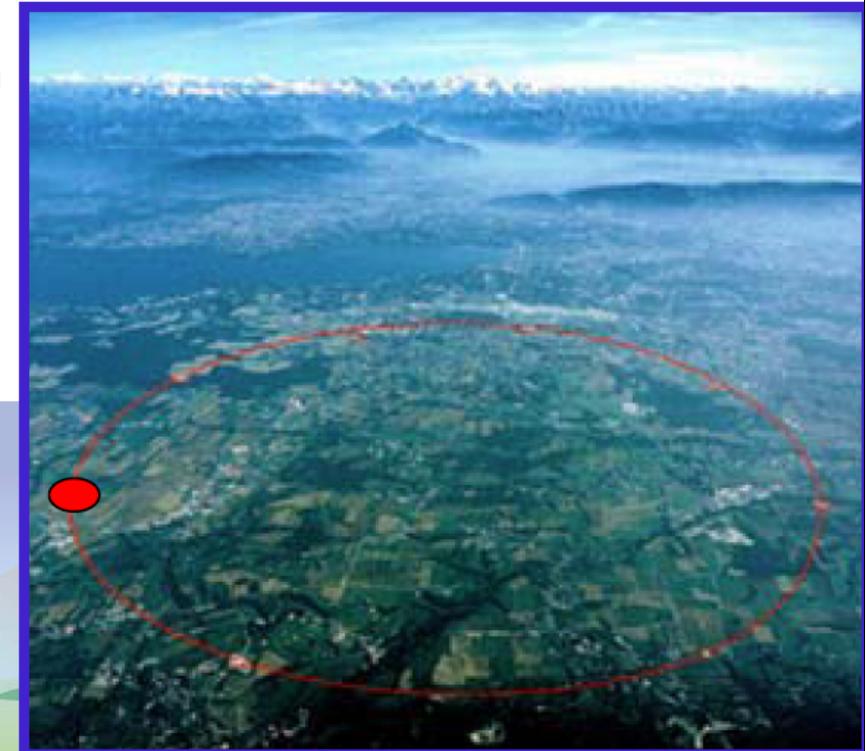
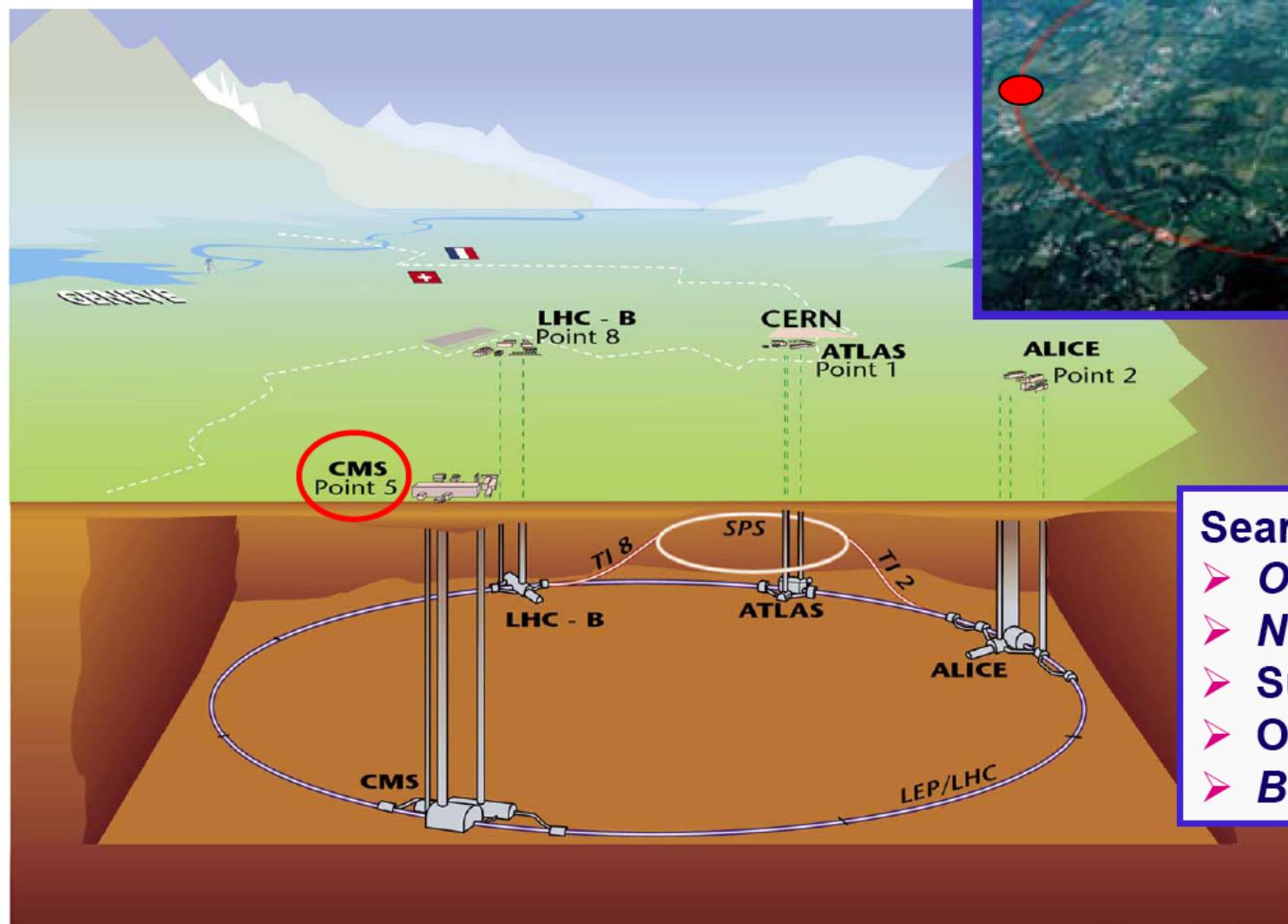


Overview of the USCMS Tier-I at FNAL

Jon Bakken
Virginia Tech Visit to FNAL
April 19, 2010

The Large Hadron Collider (LHC)

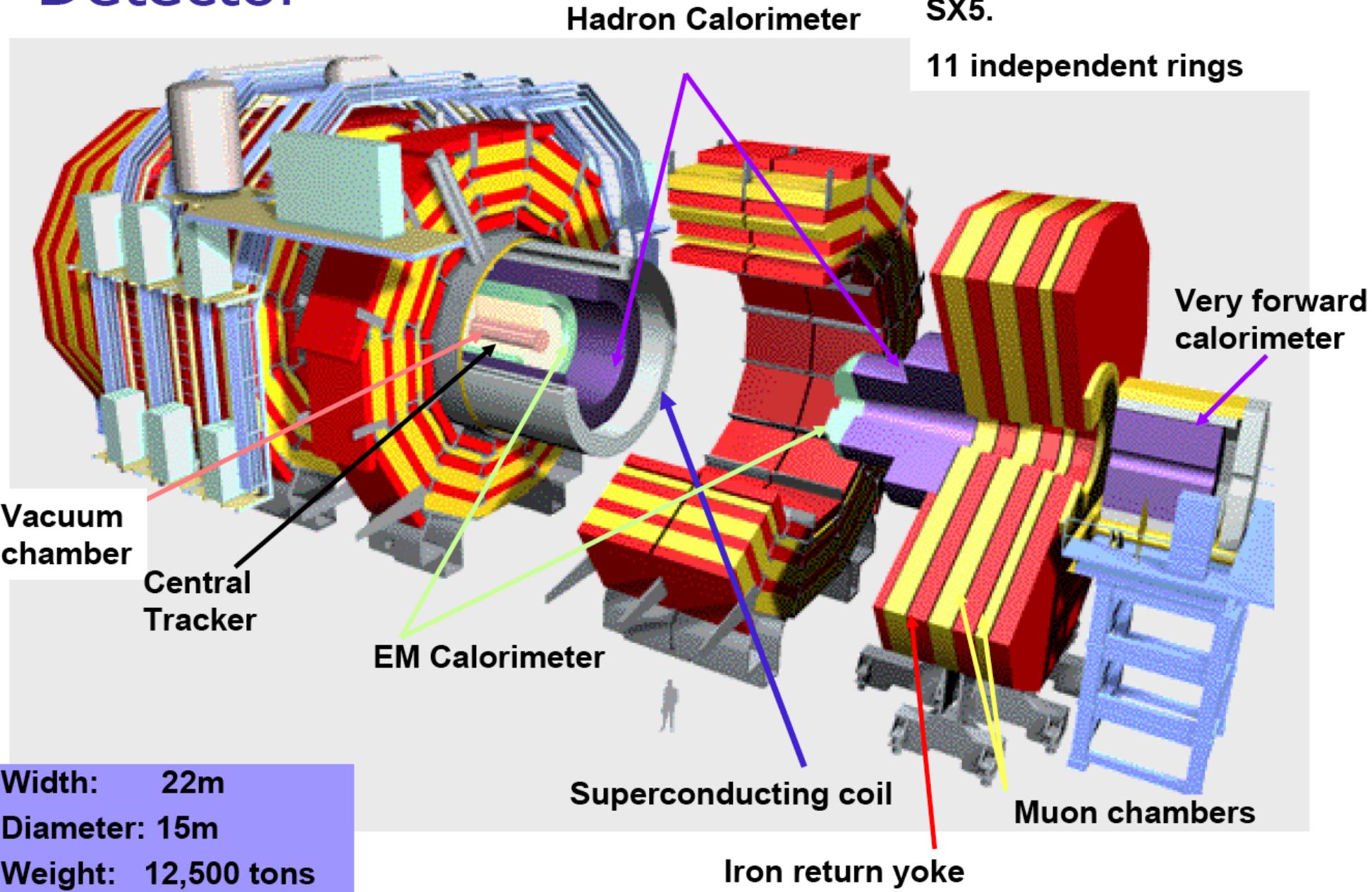
- Proton-proton collider (14 TeV energy)
- 27 km in circumference, 50-175m deep
- between Jura mountains (France) and Lake Geneva (Switzerland)



Search for

- *Origin of Mass*
- *New forces/particles*
- *Supersymmetry*
- *Other new particles*
- *Black Holes*

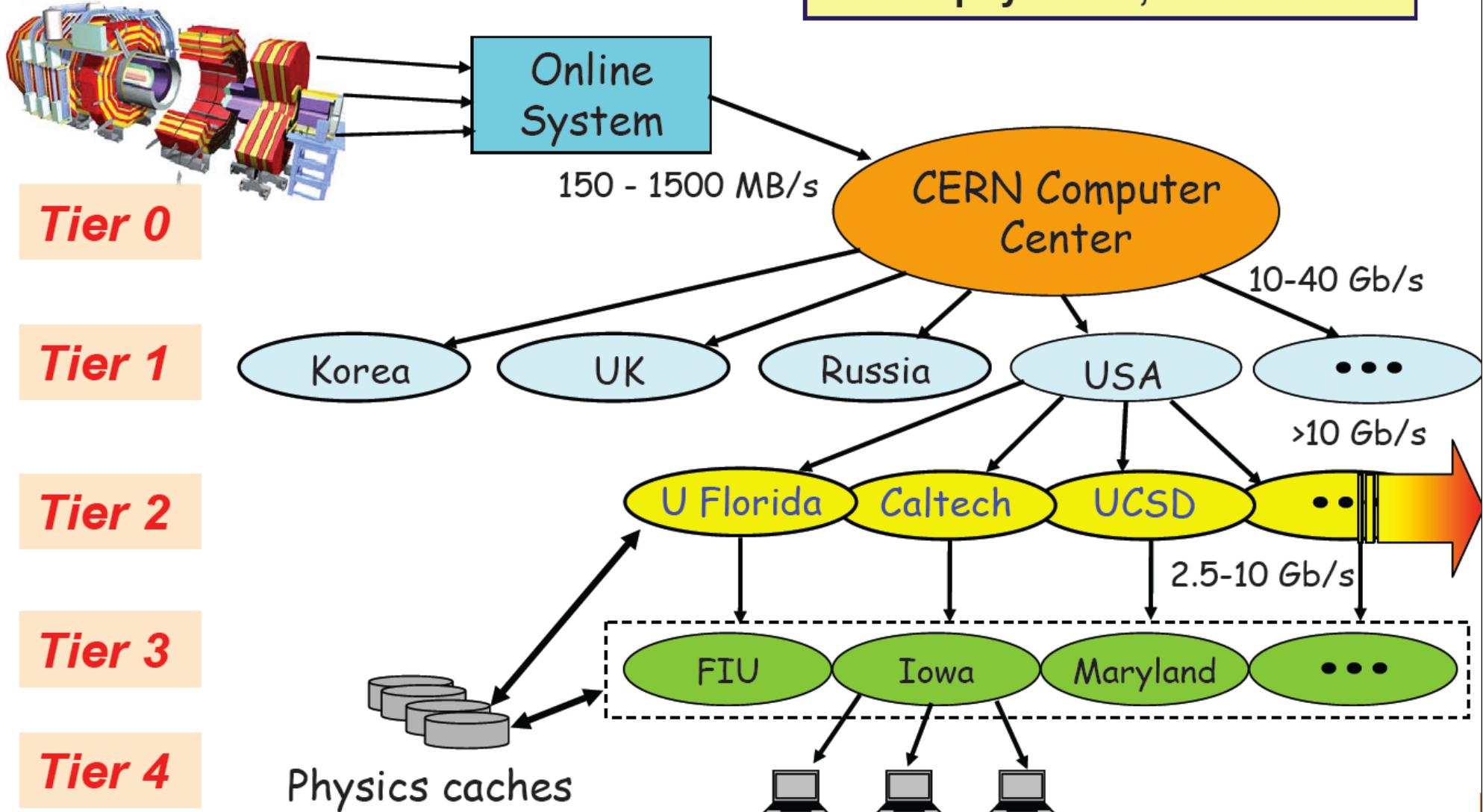
Compact Muon Solenoid (CMS) Detector



LHC Global Data Grid

CMS Experiment

➤ 5000 physicists, 60 countries



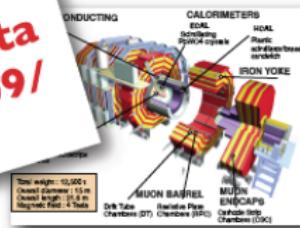
CMS Computing Metrics



♦ Reminder of CMS Computing Model:

- ★ user submit jobs, workload management system submits to where data is
- ★ data movements are triggered by operators, physics organizers, users

**Metrics Fully
Established in Data
Challenges CSA09/
CCRC**



100MB/s between Tier-1 sites
25k-50k jobs/day at Tier-1s

Approximately equal ingest rate from Tier-2s as Tier-0

75k-150k jobs/day at Tier-2 sites between analysis and simulation
40 users supported on average

Tier-0

600-800MB/s to 7 Tier-1 centers

Tier-1

Tier-1

Tier-1

Tier-1

Tier-2

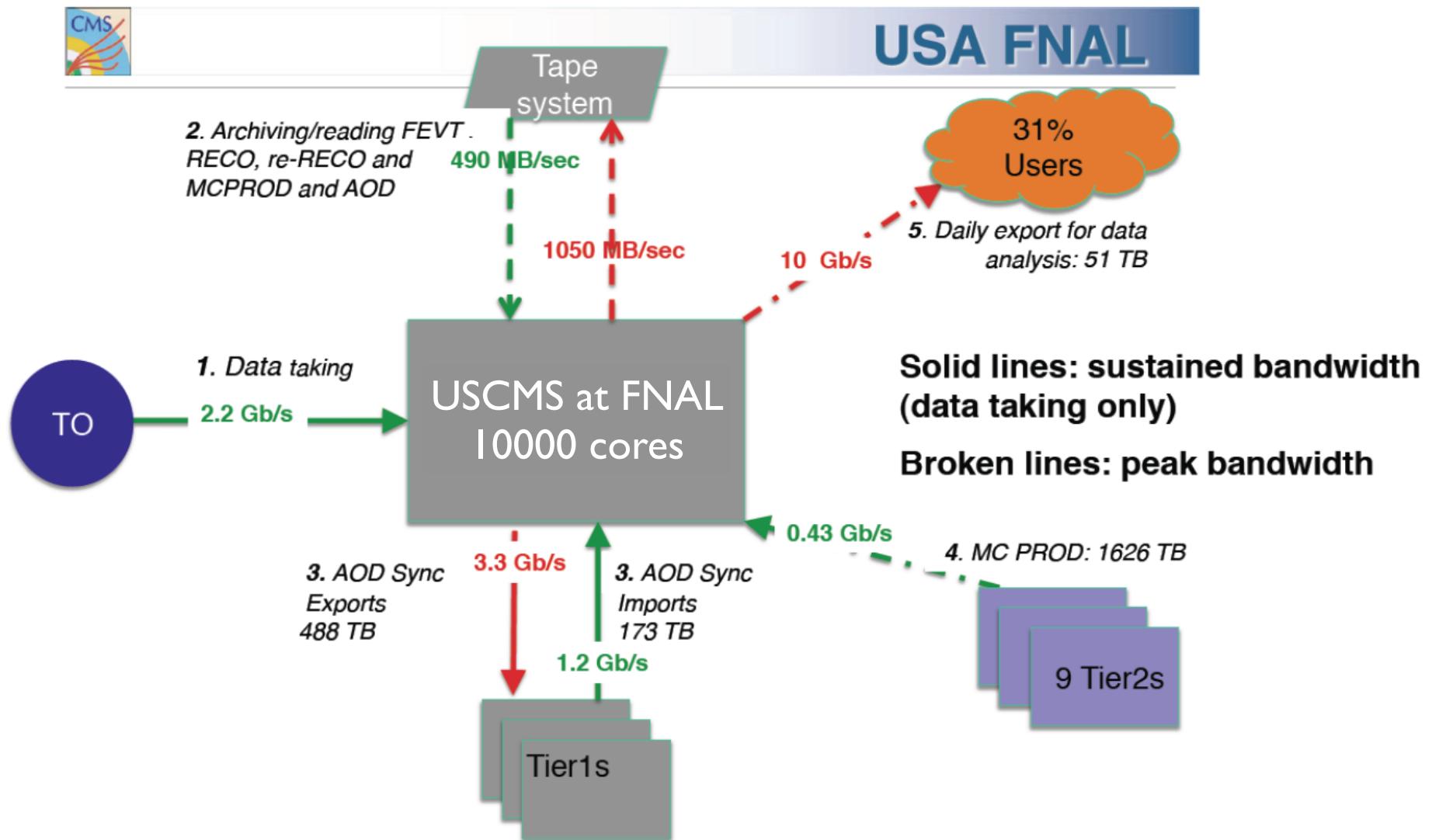
Tier-2

Tier-2

Tier-2

Burst Transfers driven by user needs
50MB/s-500MB/s
Full mesh of permutations

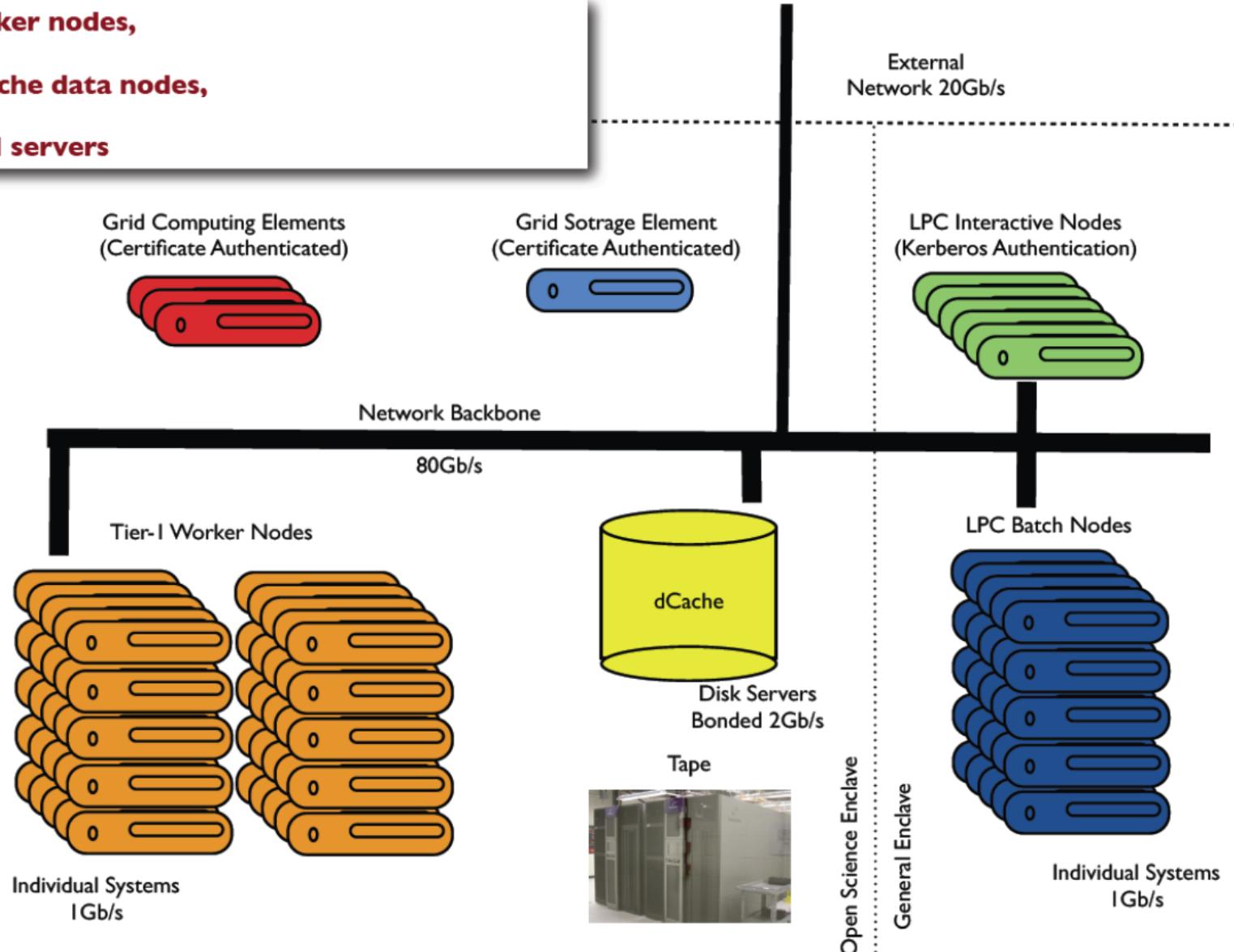
CMS Bandwidth Needs



Site Networking



- 1300 worker nodes,
- 200 dCache data nodes,
- 75 critical servers



After setting the goals for the USCMS facility, we went through 3 phases:

- We can code something that solves our exact problem ourselves and avoid all the issues we don't like
 - One moment of thought makes you realize very smart people have been working for years and this approach is foolhardy
- We should use standard industry tools and buy products from professionals.
 - This doesn't work for everything since our environment has special needs that are not addressed adequately by industry tools
 - You can't afford it, and it takes FTE effort to run industry tools, too
- You figure out who has solved problems like yours **and continues to support them and work with them** - **OSG is a very good choice.**
 - Keeping production services running is HARD!
 - In-house expertise is mandatory, outside groups can't do it alone.

The USCMS facility has the following major components:

- **Networking**
 - Cisco Nexus 7000 & 6509. 40-60G between switches. (Industry)
- **Tape storage** for archiving data
 - We use FNAL's **Enstore**. About 6 PB on tape right now. (Lab)
- **Data disk** for serving data to production and users
 - We use DESY/FNAL **dCache**. About 6 PB of disk right now. (OSG)
- **Data transfer** - We use **dCache/SRM** for this as well (OSG)
- **User disk** for a global file system, home, data and scratch areas
 - We use **BlueArc** - About 300 TB right now. (Lab)
- **Interactive** nodes - we have load-balanced cluster & direct connections
 - We use **Scientific Linux**, **Cisco** load balancer (Lab/Industry)
- **Batch** worker nodes
 - We use Wisconsin's **Condor** - 10,000 batch slots (OSG)

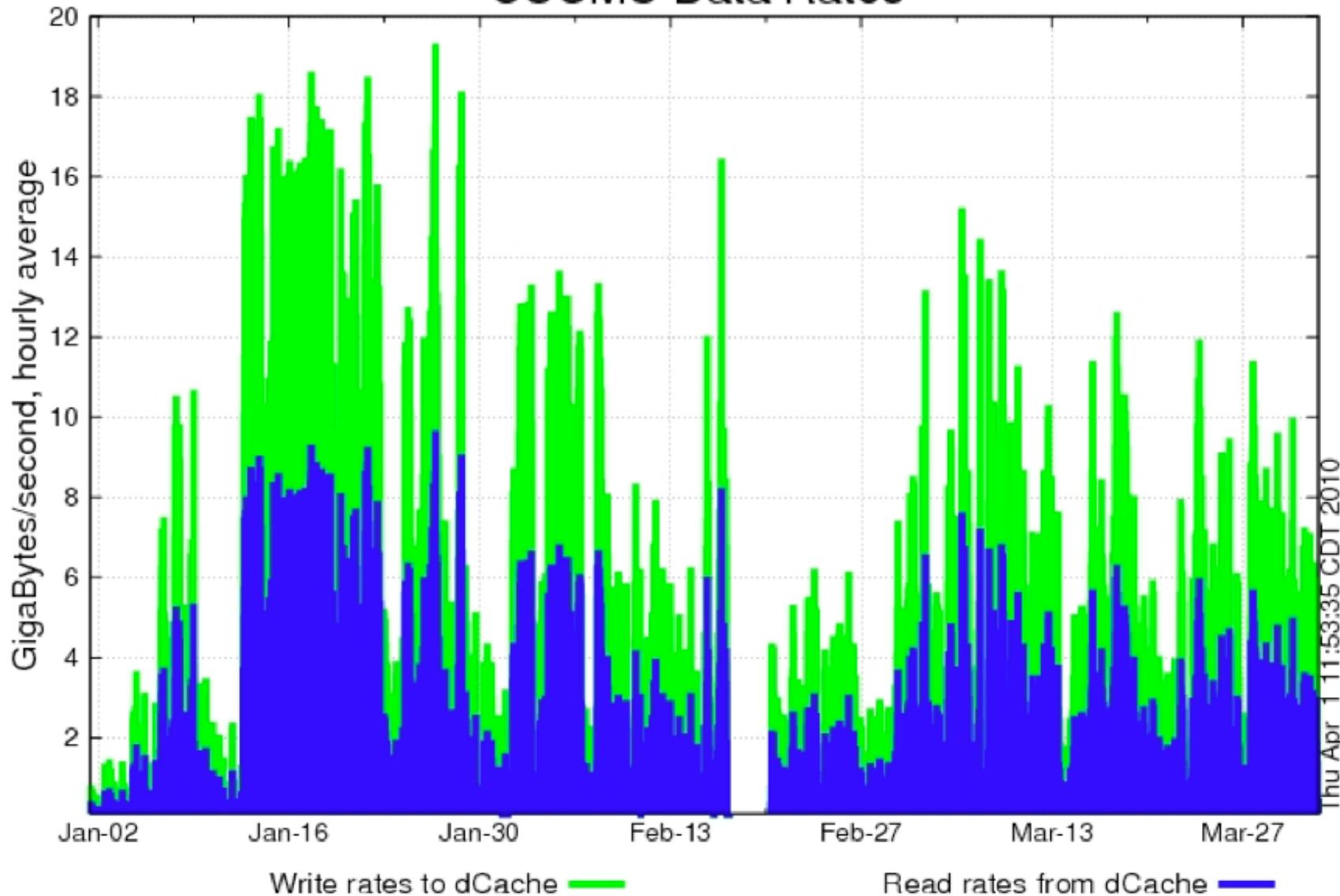


Data disk and transfers:

- We chose **dCache/SRM** since FNAL has long history of developing this with DESY and it was working in our environment.
- Tape, local transfers and WAN transfers all integrated into 1 package.
- dCache is complex to run, and if you do not have tape requirements you should consider Hadoop and Bestman. Many USCMS Tier-2s have used it and are happy - the system is reliable and “less work” to maintain. **OSG also supports Hadoop and Bestman**
- Lustre is also a choice, but it is just as complex as dCache to run and has other considerations

- Our dCache system performs well - ~30K simultaneous transfers, with aggregate IO rates ~18 GB/sec.

USCMS Data Rates



Batch workers:

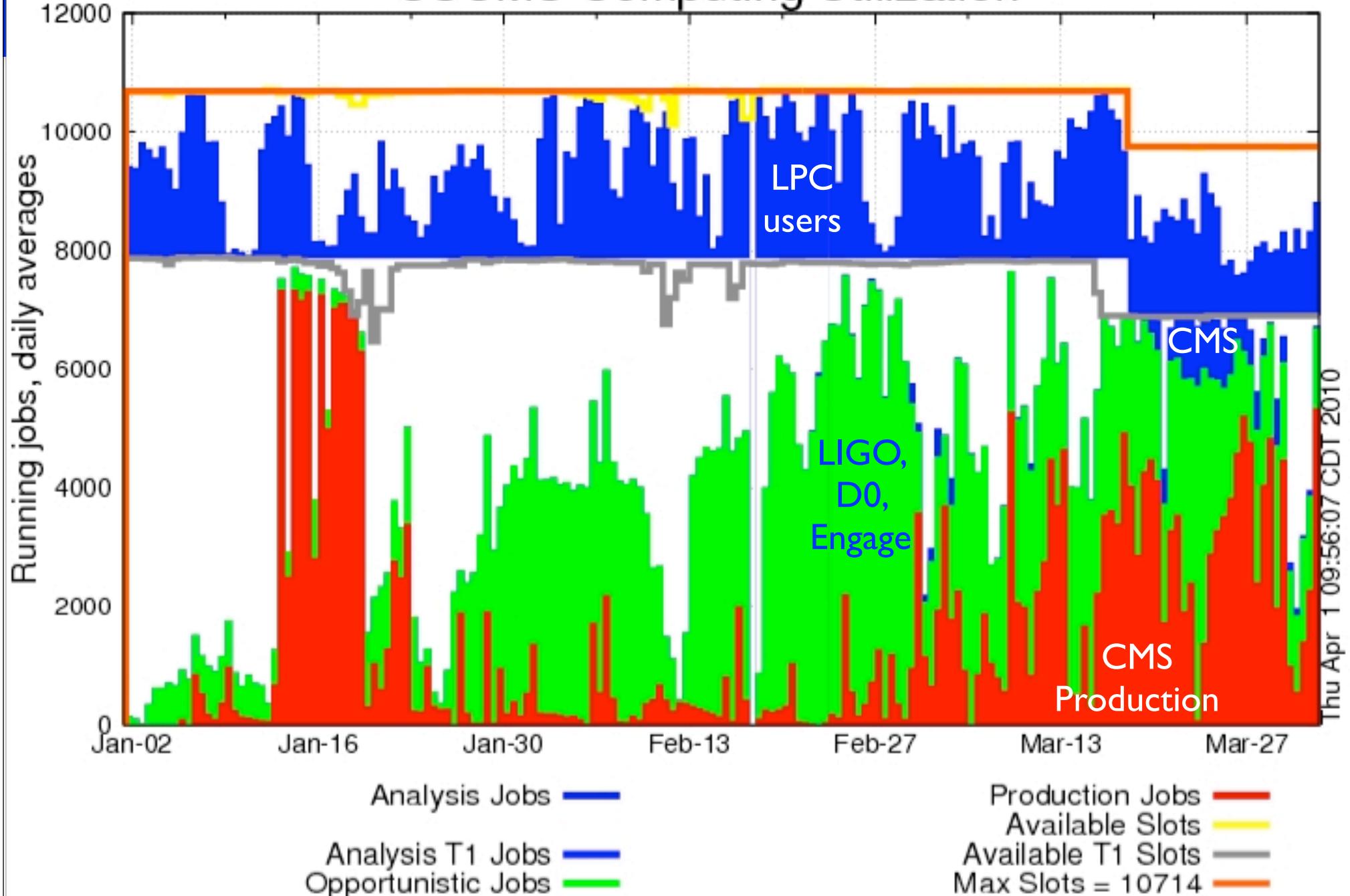
- We use Condor for 3 reasons:
 - It works, it is very reliable and it is highly configurable.
 - priorities, fair-share, extensive logs, lots of documentation, easy to use
 - Support from Wisconsin is fantastic
 - We've pushed the scaling envelope several time in recent years and the Condor team has always provide solutions
 - Security concerns are always addressed in a timely fashion
 - Requests for new features or enhancements are satisfied quickly, essentially making the Condor team part of our team here.
 - Send jobs to fastest cores preferentially, multiple cores for same job, glideinWMS, etc..
- We've investigated other batch systems, and they do not satisfy these 3 items as well as Condor.
 - And, the other batch systems are typically not free and have per core license fees. Multiplying by 10000 cores is a big number.

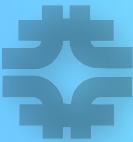


2 Condor Batch Systems within USCMS:

- Production Cluster - This cluster is in the “Open Science Enclave” at FNAL, meaning it accepts grid jobs from any recognized VO, typically through the FermiGrid submission host. (7000 cores)
 - CMS production work is given highest priority, but some idle cycles
 - We expect this cluster to be 100% busy with CMS work soon.
 - Some special rules - for example, disks can't be executable on these workers and on nodes inside the FNAL general enclave.
- LPCCAF Cluster - This cluster is in the “General Computing Enclave” and its slots are reserved for USCMS LPC users only and all submissions are local submissions from the LPC interactive cluster. These nodes were purchased with money specifically ear-marked for the LPC at FNAL. (3000 cores)
 - We still tightly control access, but no restrictions on disk mounting and it is generally easier for LPC users to use these resources.

USCMS Computing Utilization





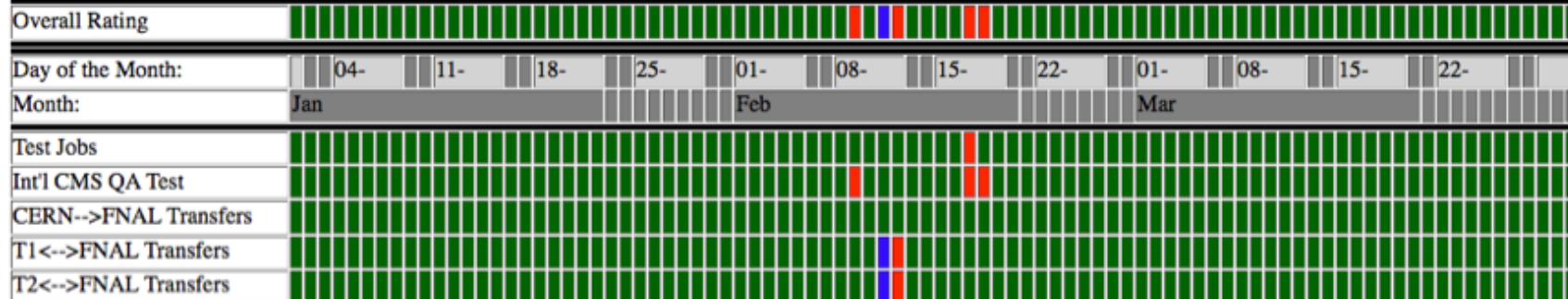
Measuring performance - Metrics

- It's very important to make sure you are achieving your goals.
- Within the OSG, a extensive suite of tests are performed & these are very useful for understanding how well you are doing.
- For the Tier-I, we also have CMS, EGEE, and our our performance tests

When you do not meet your goals - life can get complicated and lots of discussions about why it happened. We are required to report our metrics to FNAL, CERN, and the DOE.

- If you are not careful, nonproductive finger pointing. (The network failed - how could the test pass?)
- Very important to remember that the metrics are management tools, and your real goal is to provide a reliable service to your customers.

T1_US_FNAL Quality Metrics



USCMS Tier1 Combined Metric Assesment



Availability Ranks for the last 90 days [2010-01-02 -> 2010-03-31]

