



Challenges for the CMS Computing Model in the First Year

Ian Fisk

On behalf of CMS Offline and Computing

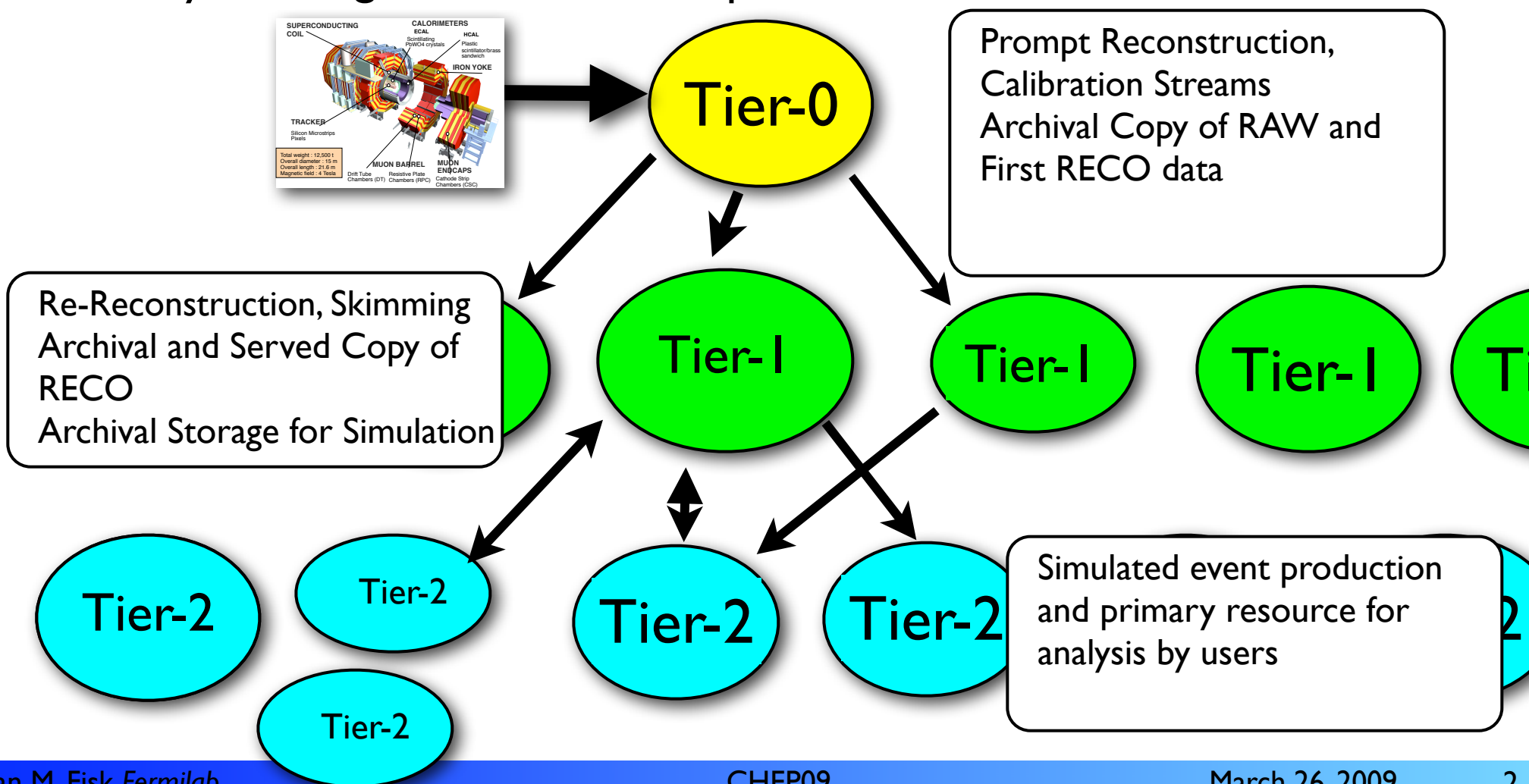
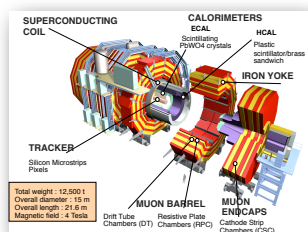
CHEP09

March 26, 2009

CMS Distributed Computing Model

CMS has been developing a distributed computing model from early in the experiment

- ➔ Variety of motivating factors (infrastructure, funding, leverage)
- ➔ Many challenges still face the experiment

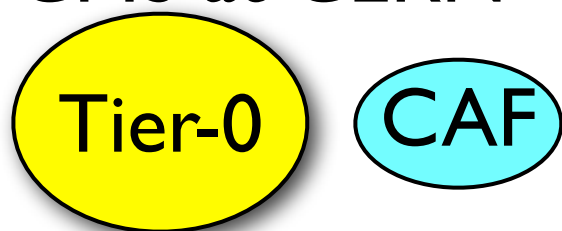


Commissioning the CMS Model

CMS will be commissioning a distributed computing model and a detector simultaneously

- ➔ There are not enough resources at any single location to perform all the analysis. (Run2 comparison)

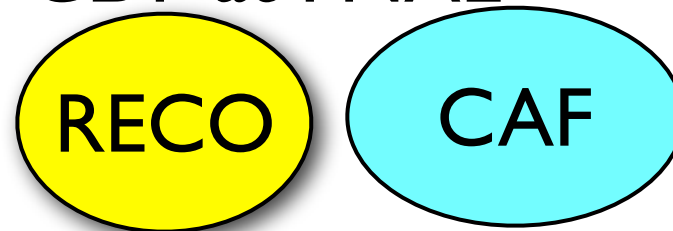
CMS at CERN



Host Lab

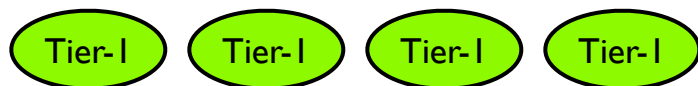
Roughly a Factor of 4 more analysis resources centrally located for CDF

CDF at FNAL

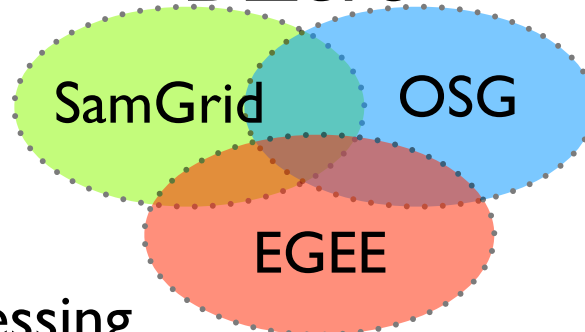


- ➔ During running all the Reprocessing resources are located at remote facilities (Run2 comparison)

CMS



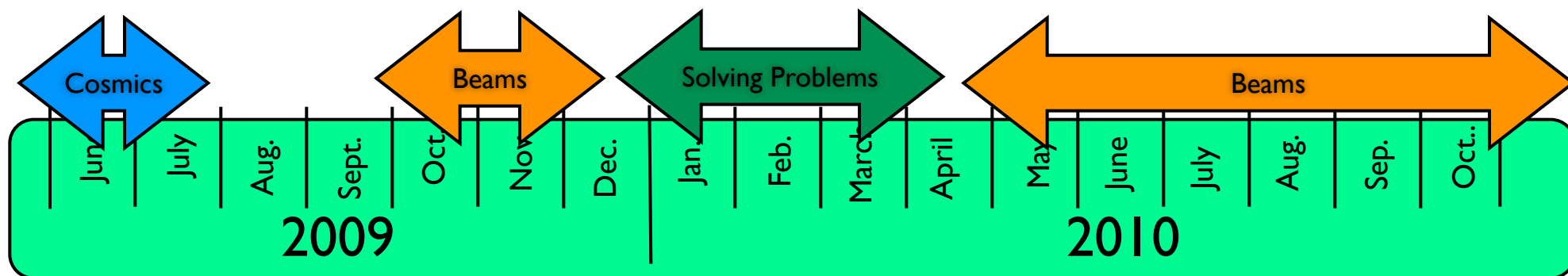
DZero



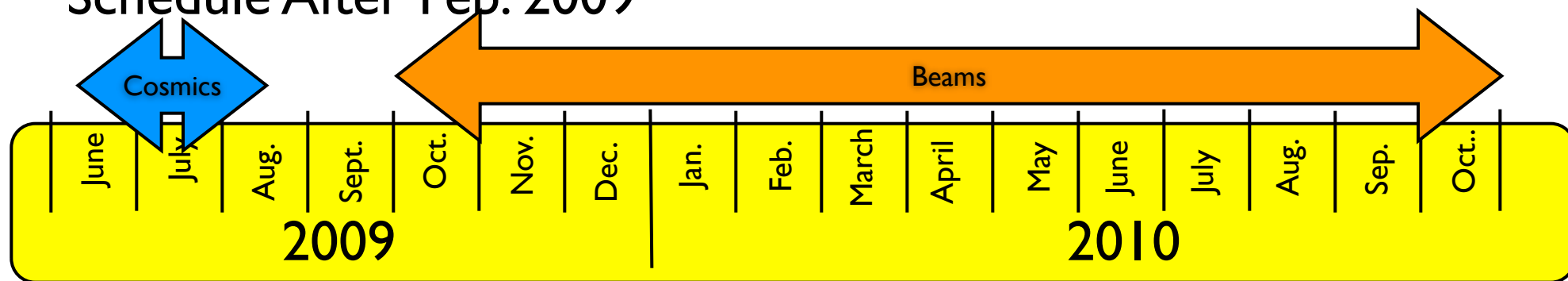
- DZero did successful global reprocessing
 - Well after other elements were commissioned

Schedule Updates

Schedule Before Feb. 2009



Schedule After Feb. 2009



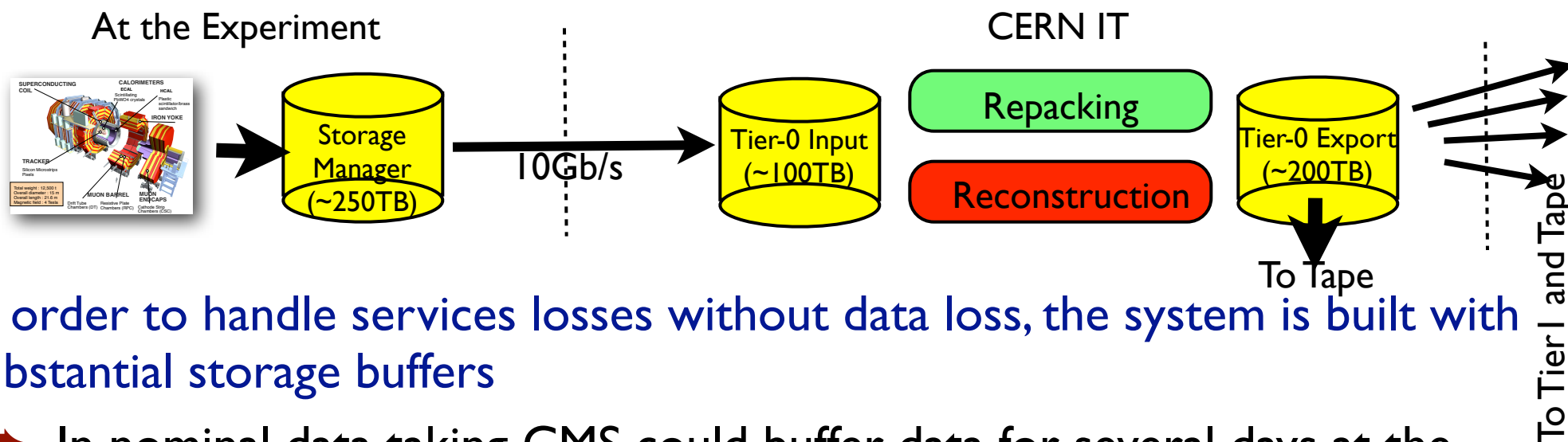
The new LHC schedule has a long initial run

- ➡ If the summer cosmics are counted, CMS will be in operations for ~16 months
- Increases needs for commissioning and development now

Initial Running Conditions

At the nominal trigger rate the CMS Data Processing Infrastructure has computing capacity to keep up with data.

- ➔ Selected events are reconstructed in an hour. Remaining data is reconstructed within a day



In order to handle services losses without data loss, the system is built with substantial storage buffers

- ➔ In nominal data taking CMS could buffer data for several days at the experiment and roughly 2 days in IT after reconstruction
- ➔ If the accelerator starts out at a low duty cycle (20%), there will be pressure to take as many events as possible
 - Overdrive the system

Initial Running Conditions

In these early scenarios CMS opens the trigger and can take data at up to 2kHz (Nominal is 300Hz)

- ➔ The rate into the storage manager exceeds the rate that can be transferred to IT by a factor of two

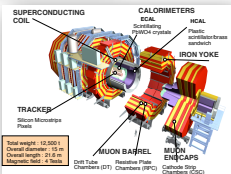


- ➔ The rate into the Export Buffer exceeds the ability of the Tier-1s to drain the data

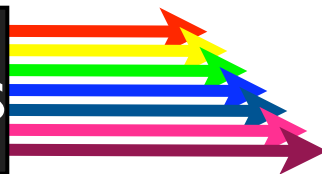
Provided CMS only takes data for a small percentage of the few day period the time in between runs allows recovery

- ➔ Allows CMS to collect additional commissioning data if it's interesting
- ➔ Requires a series of additional tests to ensure the system can be over driven for periods of time and brought back into stable operations

Equitably Distributing Data



Triggered Events



Primary Datasets

In CMS the Triggered events are divided into ~20 Primary Dataset Streams

- ➔ It should be possible to do an analysis with one stream
- ➔ Good physics reasons for dividing data and reasonable technical reasons for keeping the stream together for later reprocessing
 - Prioritize reprocessing of streams. Reprocess based on a new calibration that impacts a particular stream.
 - Goes to a family of tapes at Tier-I

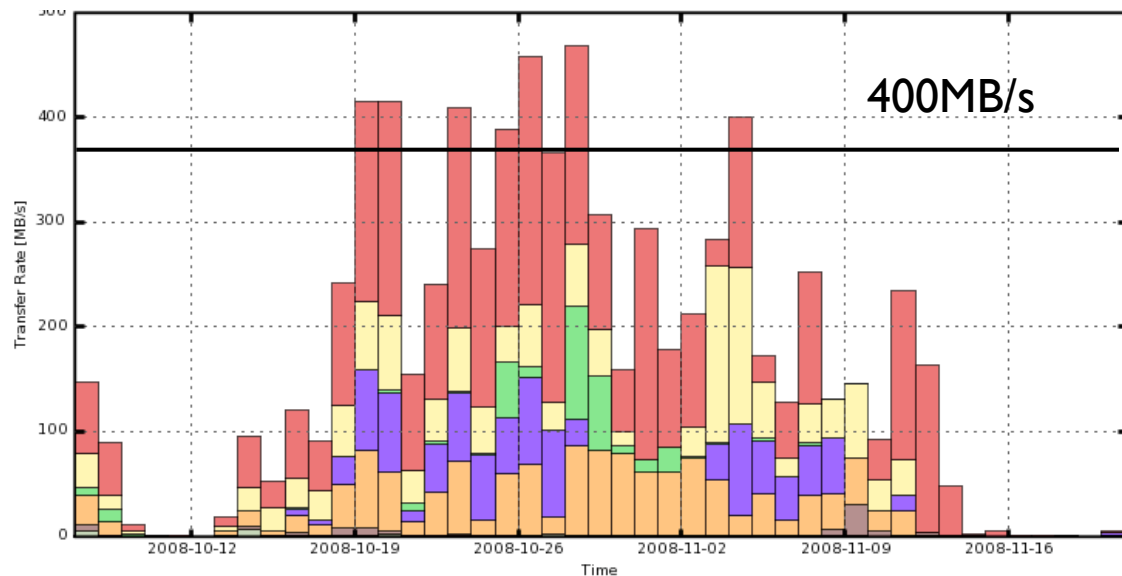
Challenging aspect is how to meet the technical needs of the experiment with the local desires of the Tier-I sites

- ➔ Spirited debate on data placement even for simulation
- ➔ Hosting a primary dataset at a Tier-I gives a site the responsibility for reprocessing

CMS now has the concept of non-custodial data into the data management

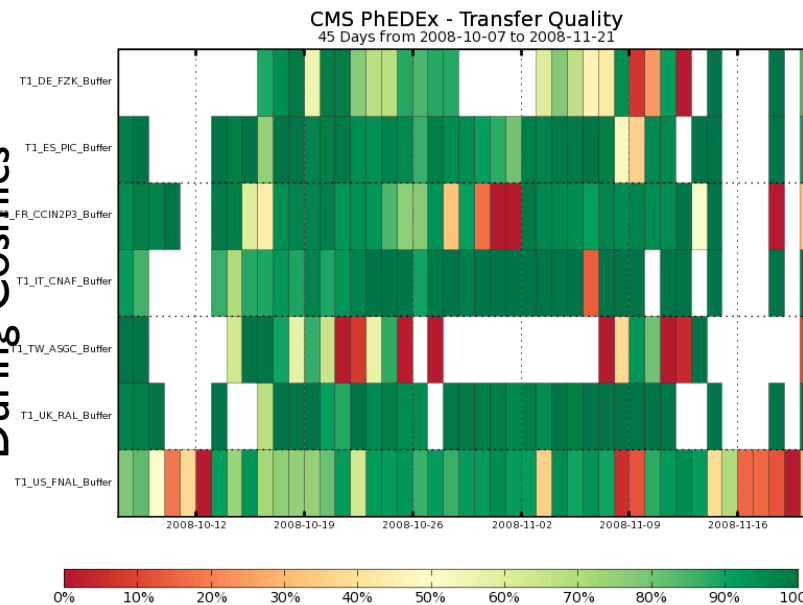
Choosing a Tier-I to Store the Data

Data Transfer from CERN to Tier-I's has become quite reliable



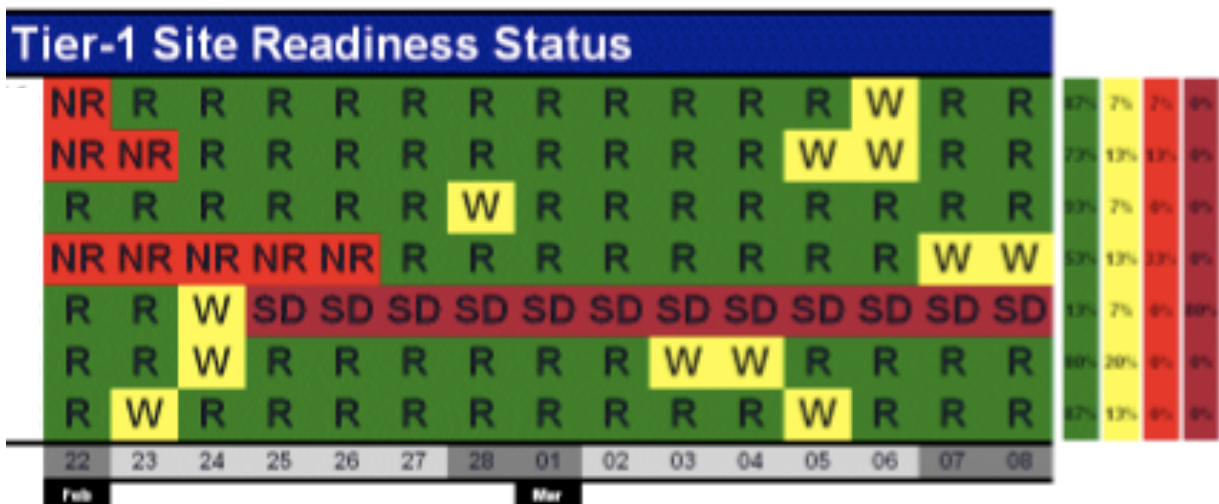
T1_US_FNAL_Buffer
T1_TW_ASGC_Buffer
T1_FR_CCIN2P3_Buffer
T1_IT_CNAF_Buffer
T1_UK_RAL_Buffer
T1_DE_FZK_Buffer
T1_ES_PIC_Buffer

Transfer Rate and Quality
During Cosmics



Trying to improve the ability to work on the site once the data is there

- ➡ Green is good
- ➡ Red is "Not Ready"



Interacting with Mass Storage

Tier-I Processing Workflows

Workflow	Processing Per Event	Nominal Tier-I	Processing Rate	Rate from Mass Storage	Number of Files per day
Reconstruction	25kSI2k*s	~1000 Cores of 2kSI2k Each	70Hz	~100MB/s	3k files (9TB)
Skimming	0.25kSI2k*s		7000Hz*	10GB/s**	300k files (900TB)

*Skimming was never expected to use the whole farm

For reprocessing the amount of data and rate from mass storage is manageable

➡ RAW data is primarily on tape and needs to be staged for processing

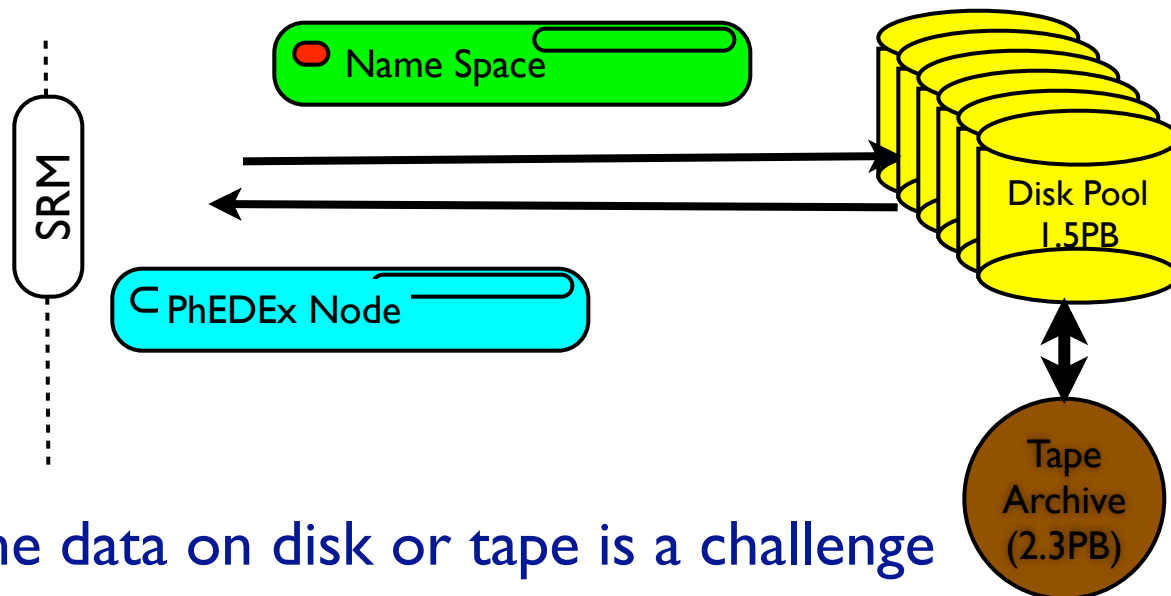
Skimming reads a portion of data to make a selection, so the IO from mass storage is substantially lower provided the files are primarily on disk

➡ Reading the whole file would not be possible -> 10GB/s

Interface to Mass Storage

SRM very successfully implemented in WLCG

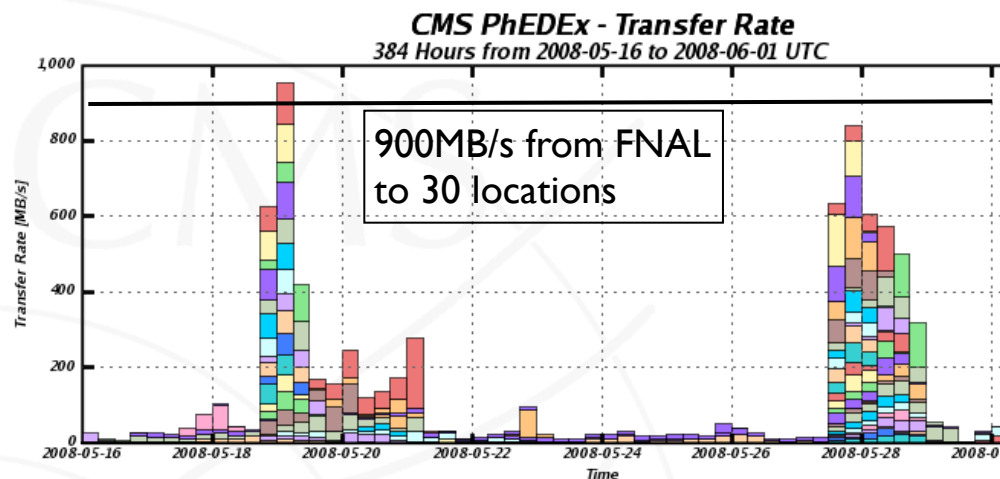
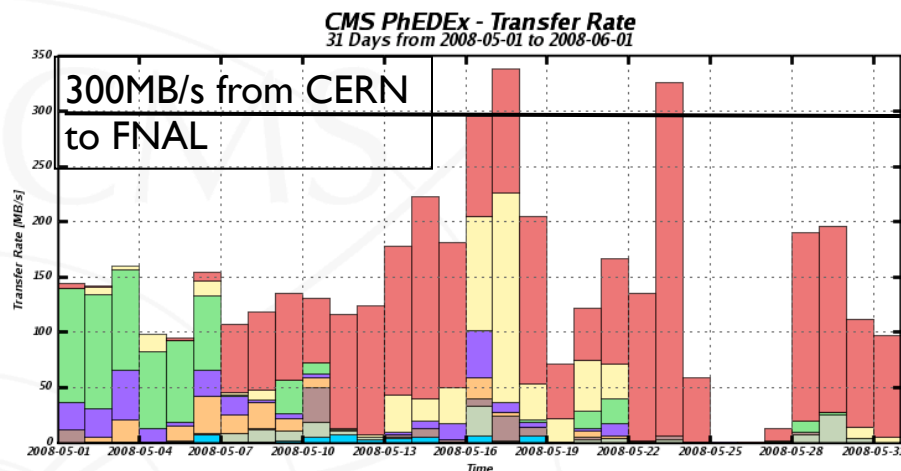
- ➔ Consistent protocol
- ➔ Load balancing to physical hardware
- ➔ Good transfer rate



How to efficiently manage the data on disk or tape is a challenge

- ➔ SRM provides a consistent interface, but it may not scale at all locations if its used to monitor what data is successfully staged and trigger staging
 - ➔ Staging requests currently are sent to administrators
- ➔ CMS has a VO box for PhEDEx at all sites to handle data management
 - ➔ Trigger transfers, verify data consistency, publish data blocks
- ➔ Anticipating challenges for staging data and already seeing issues with large scale processing of data that may have been migrated to tape
 - ➔ Work ongoing to improve processing systems

Data Serving

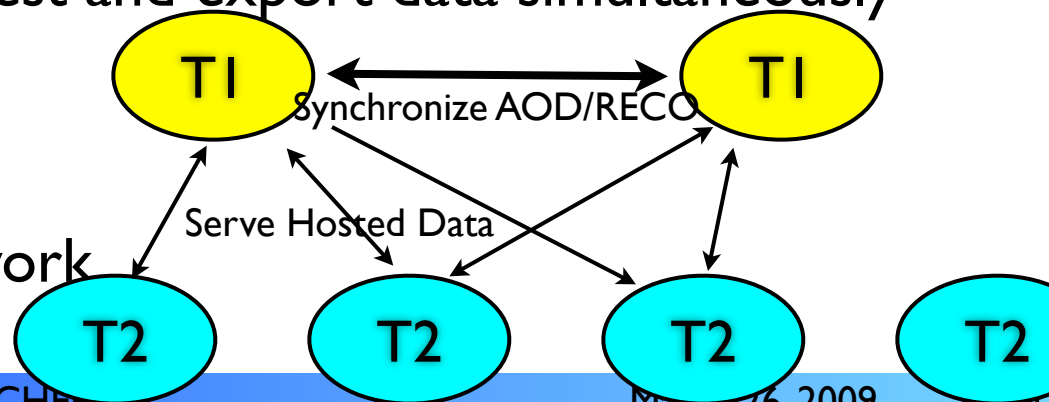


In the CMS model the Tier-1s serve the analyzed copy of the data

- ➡ While data is written once, it can be read many times
- ➡ The data serving requirements of the T1s can exceed that of CERN
 - More locations and a higher rate for bursts
 - Like CERN the T1s need to ingest and export data simultaneously

Full mesh of transfers improves data access

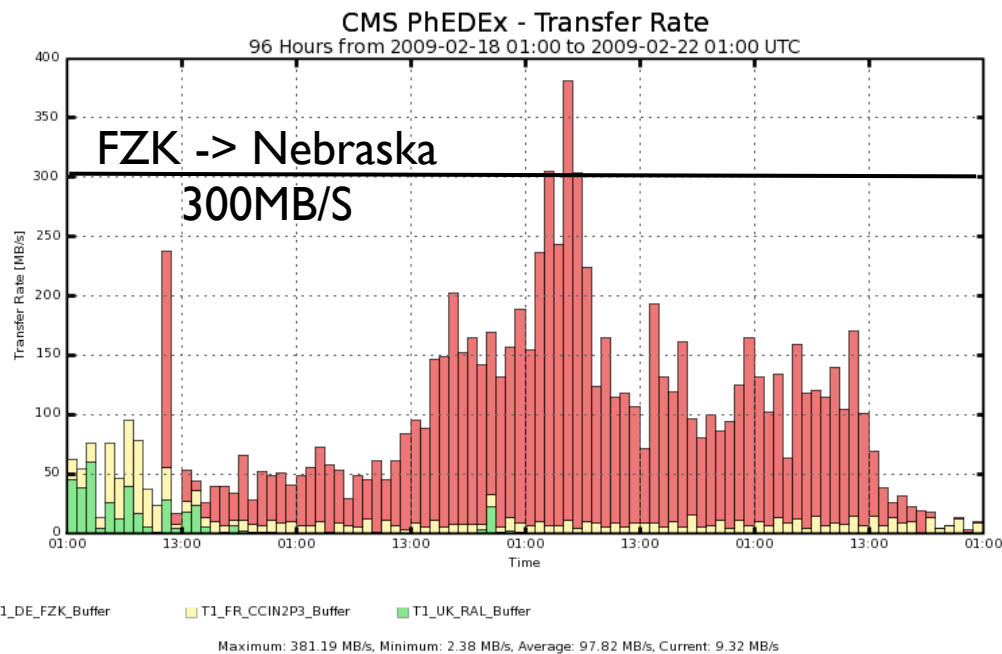
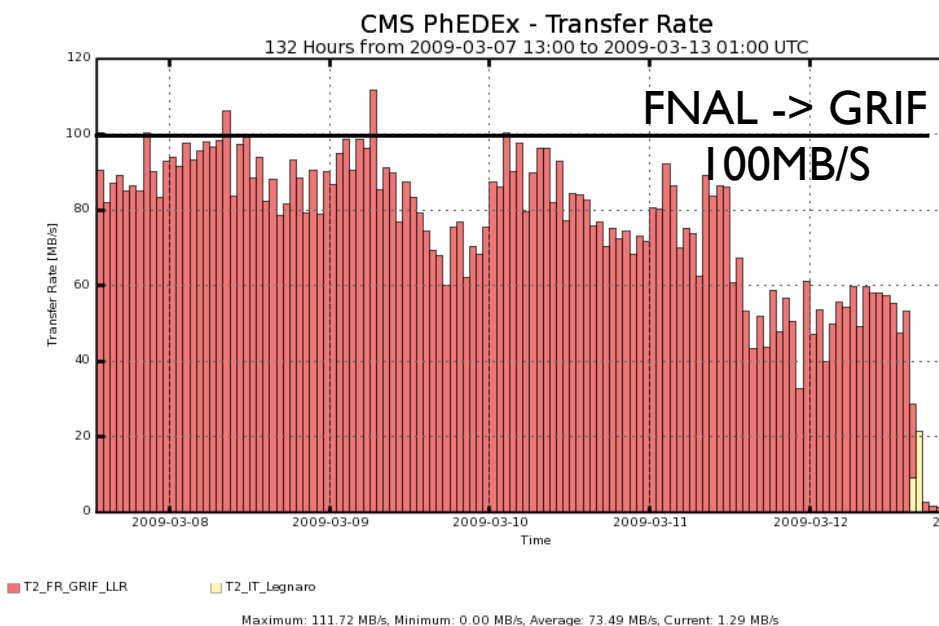
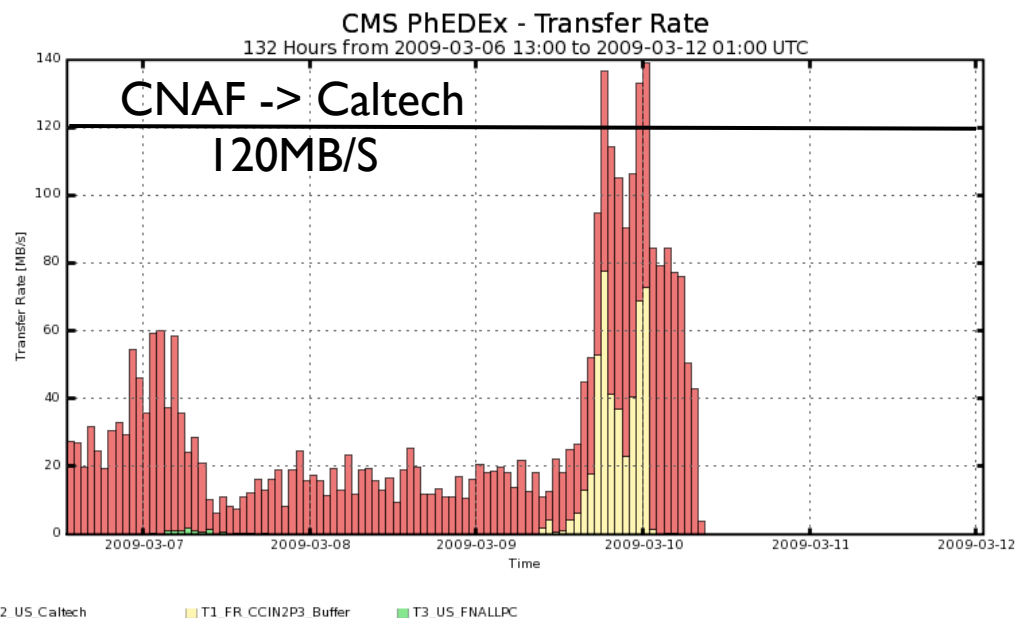
- ➡ Also increases commissioning work



Full Mesh Examples

Link Commissioning in CMS has been a long effort intensive process,

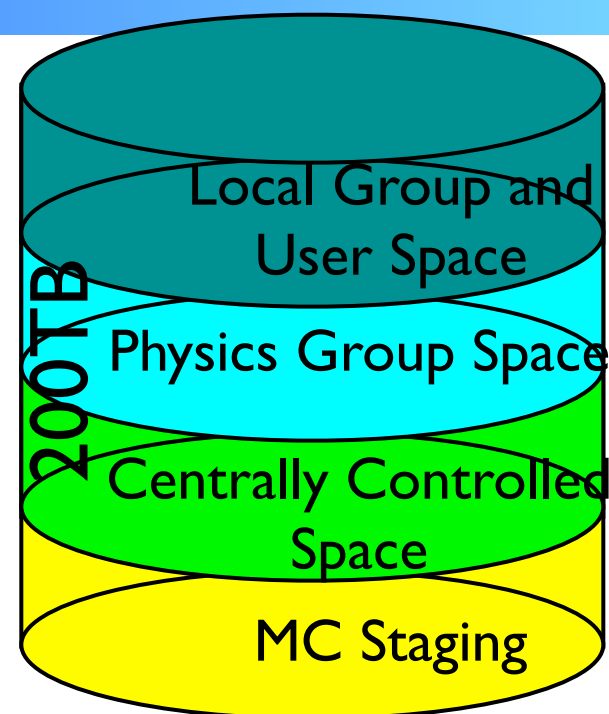
- ➔ Good performance achieved in both directions across the Atlantic
- ➔ Work ongoing



Organizing Data

A nominal Tier-2 has 200TB of disk

- ➔ Total storage at Tier-2s is enormous
 - Making sure the appropriate version of the data is being hosted at a location with resources for the community that needs to access is challenging
 - Huge collaboration with many sites
- ➔ Beginning in summer 2008 CMS began assigning blocks of storage to physics groups associated with sites
 - Increases the number of people participating in data management and puts the control closer to those working
 - Also a challenging political process
 - Assessment of how well this is working is still ongoing



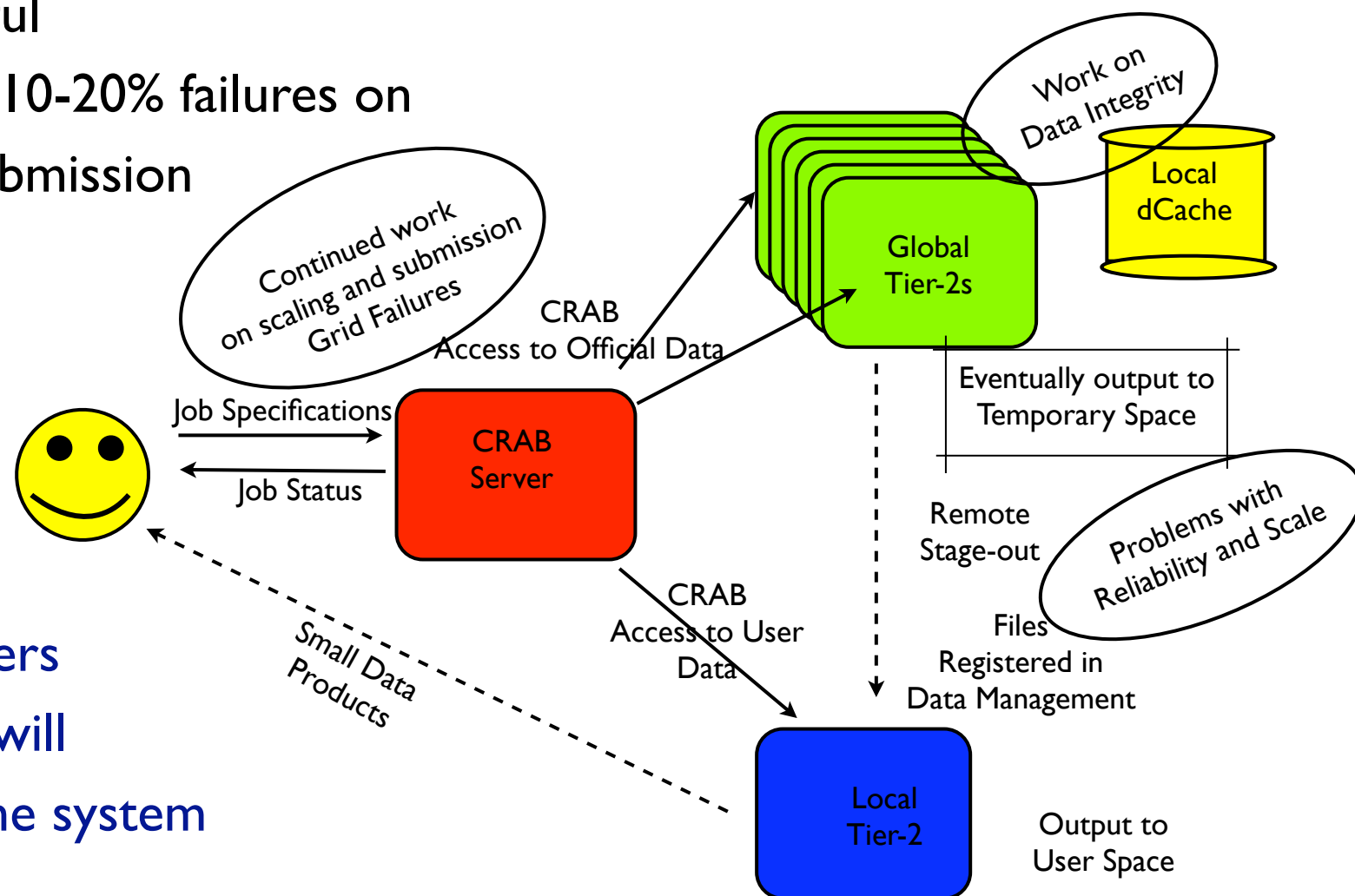
T2_US_UCSD Group Usage

Group	Subscribed	Resident
DataOps	4.24 TB	4.24 TB
ewk	5.71 TB	5.71 TB
higgs	5.84 TB	5.84 TB
susy	498.88 GB	498.88 GB
top	58.60 TB	55.50 TB
tracker	551.37 GB	490.06 GB
undefined	63.72 TB	63.69 TB
	139.12 TB	135.94 TB

Analysis

How the system will work with 2000 collaborators?

- ➔ CMS Remote Analysis Builder (CRAB) shields the user from the underlying complexity, but a many things have to succeed for analysis to be successful
- CMS sees 10-20% failures on analysis submission



Clear adding users and workflows will further stress the system

Analysis Scale and Data Access

CMS had analysis submissions from 700 individuals in 2008

➡ ~40% of the collaboration

Activity level is still much lower than expected from the Model design

➡ We expect the number of jobs to jump in the first year

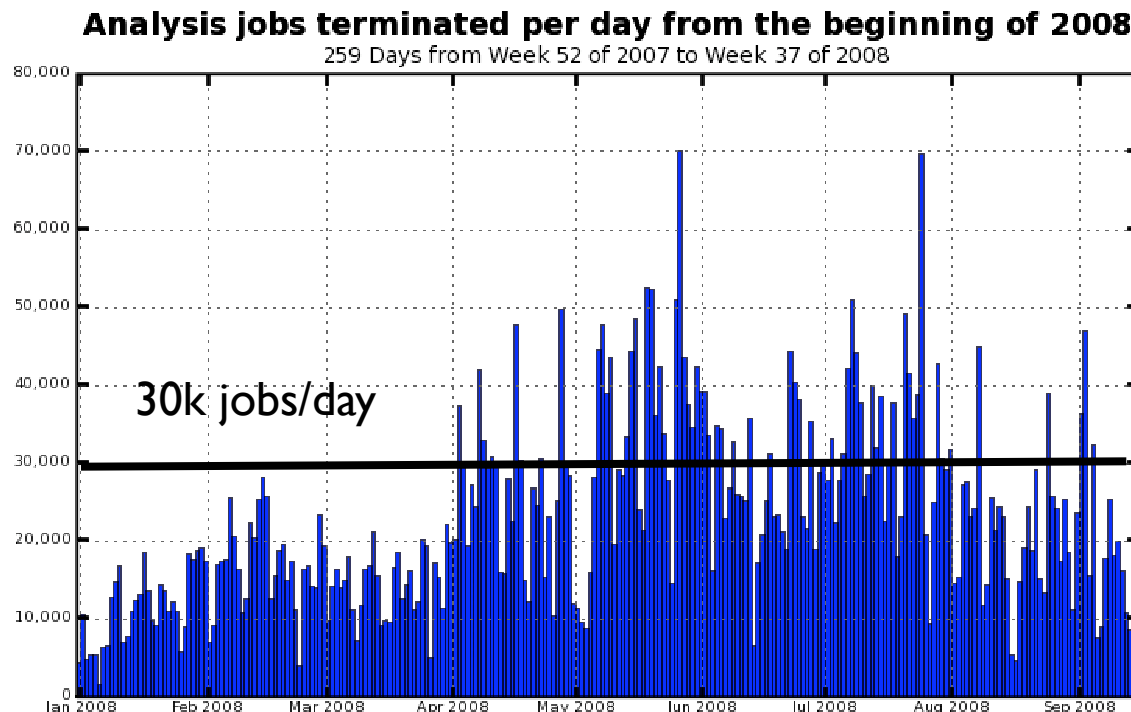
Eventually CMS will do analysis primarily on summarized (AOD) analysis objects

➡ In the first year we expect to be able to host two complete copies of the RECO data out to the Tier-2s

➡ In the presence of a long run this may not be possible

➡ Some user access to RAW data will be unavoidable

➡ Looking at VOMS roles to provide structured access to TIs



CMS has many challenges in the first year

- ➡ CMS will be commissioning a large distributed computing system while we commission the detector
 - We've worked on many computing challenges and activities to prepare
- ➡ The Run is longer than expected
 - This is a great development for physics, but a new challenge for operations
- ➡ How one distributes the data, access the events for processing, and distribute them to Tier-1s for analysis is well understood in theory
 - We will learn new lessons as we do this in practice

Will be an exciting year.