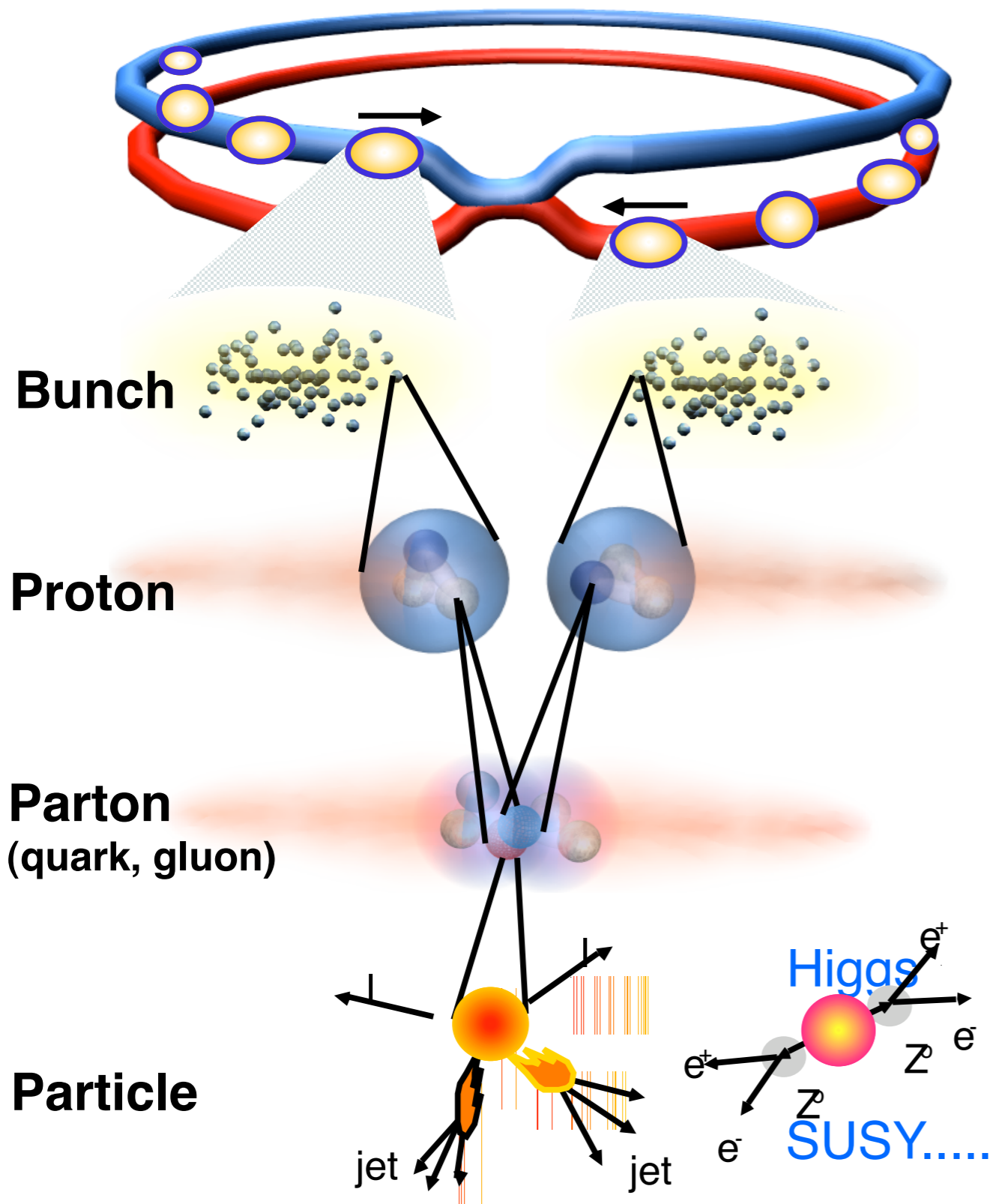


CMS Computing: Performance and Outlook

Ken Bloom
DPF 2011
August 11, 2011





Proton-Proton 2835 bunch/beam
Protons/bunch 10^{11}
Beam energy 7 TeV (7×10^{12} eV)
Luminosity 10^{34} cm⁻² s⁻¹

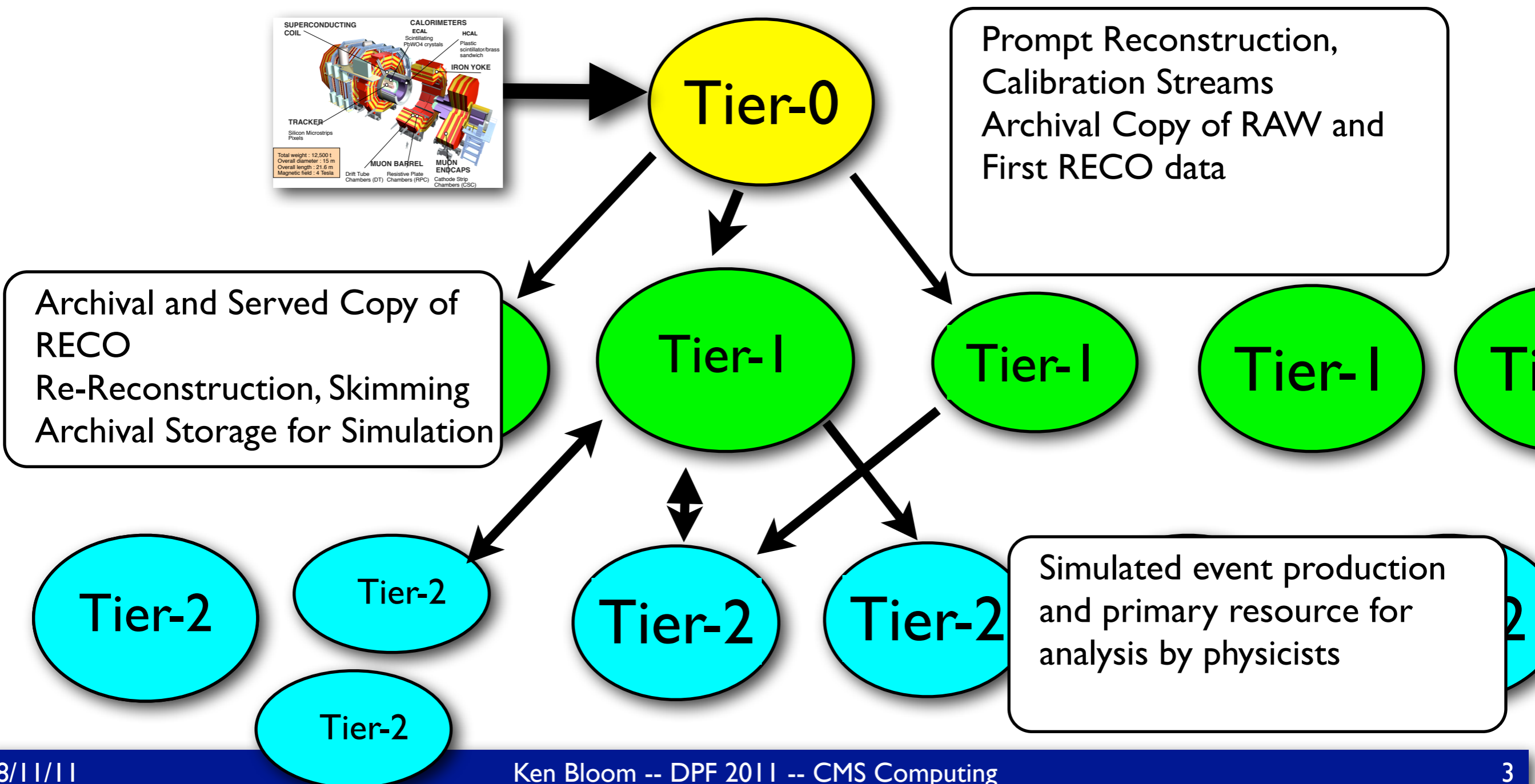
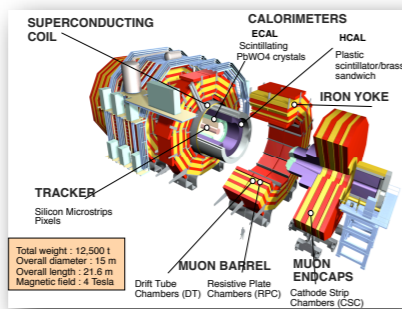
Crossing rate 40 MHz

Collision rate $\sim 10^9$ Hz

New physics rate ~ 0.00001 Hz

Event Selection:
1 in 10,000,000,000,000

- ▶ Distributed computing model planned from the start
- ▶ Variety of motivating factors (infrastructure, funding, leverage)
- ▶ Challenges in making the distributed model work, but worth it



- ▶ All workflows ran at the designated facilities from Day 1!
- ▶ T0 handled many different datasets/workflows:
 - ▶ 100 different datasets, 13.9B events, 674 TiB
- ▶ Data re-processed many times at T1
 - ▶ 19 re-recos of data, 17.2B output events, 2.4 PB
 - ▶ 4 MC re-reco passes, 8.3B output events, 2.9 PB
- ▶ MC production at T2 and T1:
 - ▶ 3.6B events, 3.9 PB, maximum > 500M events/month
- ▶ Transfers throughout the system:
 - ▶ Kept up with data taking, analysis datasets to T2 within a day
 - ▶ Peak rates > 600 MB/s T0→T1, 1200 MB/s T1→T2 like those of original computing model
 - ▶ More T2→T2 transfers than originally envisioned

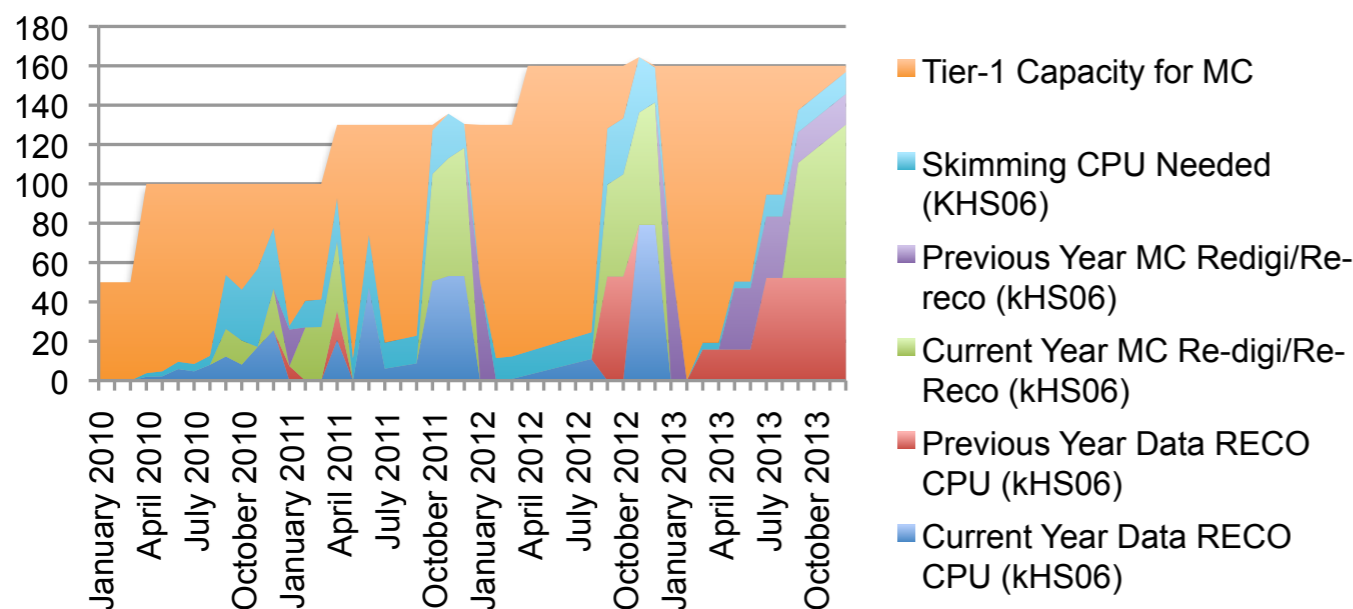
- ▶ Successful migration of analysis to T2 sites:
 - ▶ Did not know for sure that grid could handle hundreds of users
 - ▶ 450 unique analysis users/week, 150K analysis jobs per day
 - ▶ Both grew throughout the year
- ▶ 75 physics papers on 2010 data submitted/accepted published with more in pipeline; computing was never a bottleneck
- ▶ All of this during rapidly changing experimental conditions!

- ▶ All wonderful, but LHC only delivered 45 pb⁻¹.
 - ▶ Smaller than what the system was designed for
 - ▶ A good opportunity to shake down the system under relatively small load

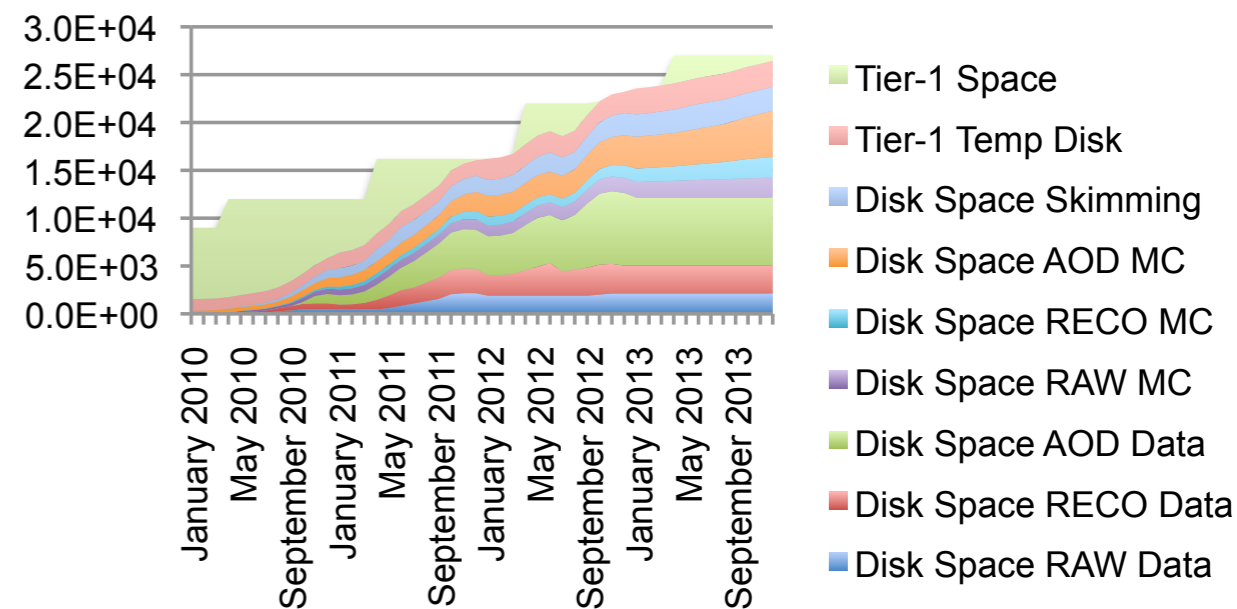
- ▶ LHC has reached 2011 target luminosity in June
 - ▶ 16 interactions/event anticipated before September technical stop
→ more pileup
 - ▶ Could be even more after that....
- ▶ Event sizes expected to double from 2010 values
 - ▶ 0.8 MB/event for RECO format, 0.2 MB/event for AOD
 - ▶ Processing time expected to quadruple to 96 HS06 sec/event
- ▶ Trigger rate nominally 300 Hz, but a challenge to keep it there
- ▶ This information, plus experience with real LHC operations in 2010, was the basis for a very thorough modeling effort
- ▶ Result: CMS computing expected to be resource-limited in 2011 and 2012, even after squeezing a lot of efficiency out of operations

- ▶ T I's extremely busy when re-processing, less so otherwise
- ▶ Make T I's primary site for MC production when not re-processing
- ▶ Keep fewer copies of data at T I
- ▶ Originally envisioned seven copies of AOD at T I, now just two
- ▶ Only one copy of RECO kept on tape
- ▶ Encourage physicists to move from RECO to AOD
- ▶ Need regular deletion campaigns to stay within resource envelope

Tier-1 Processing Resources

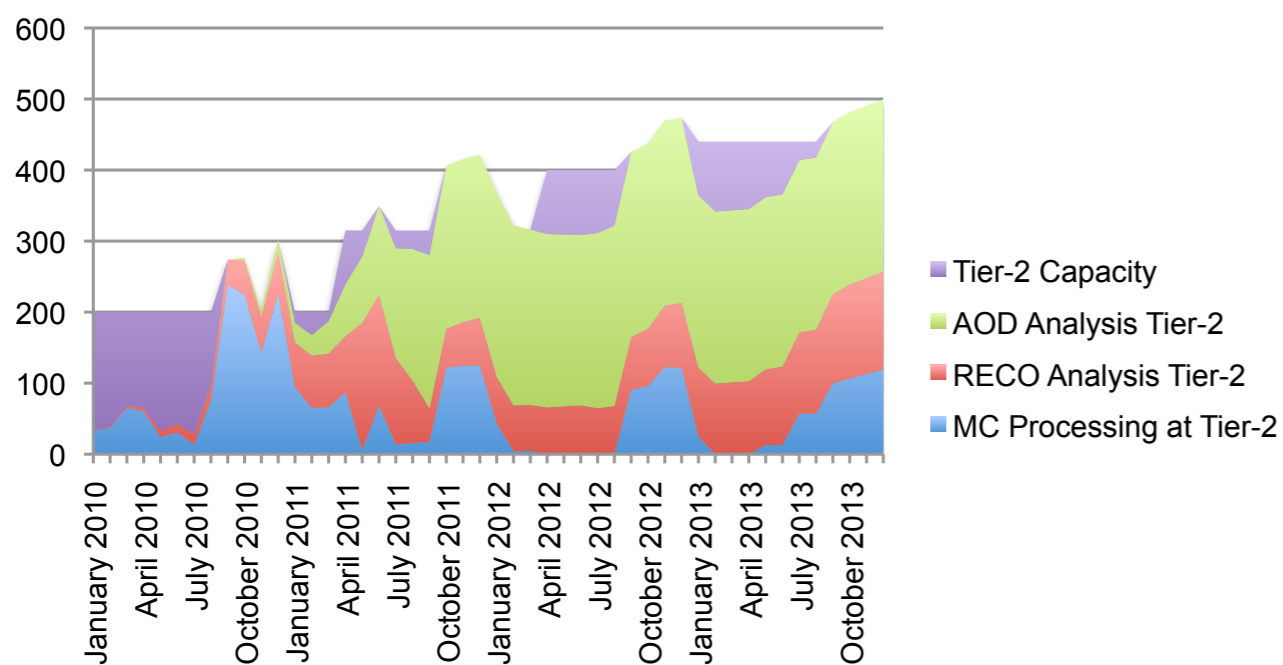


Tier-1 Disk Storage

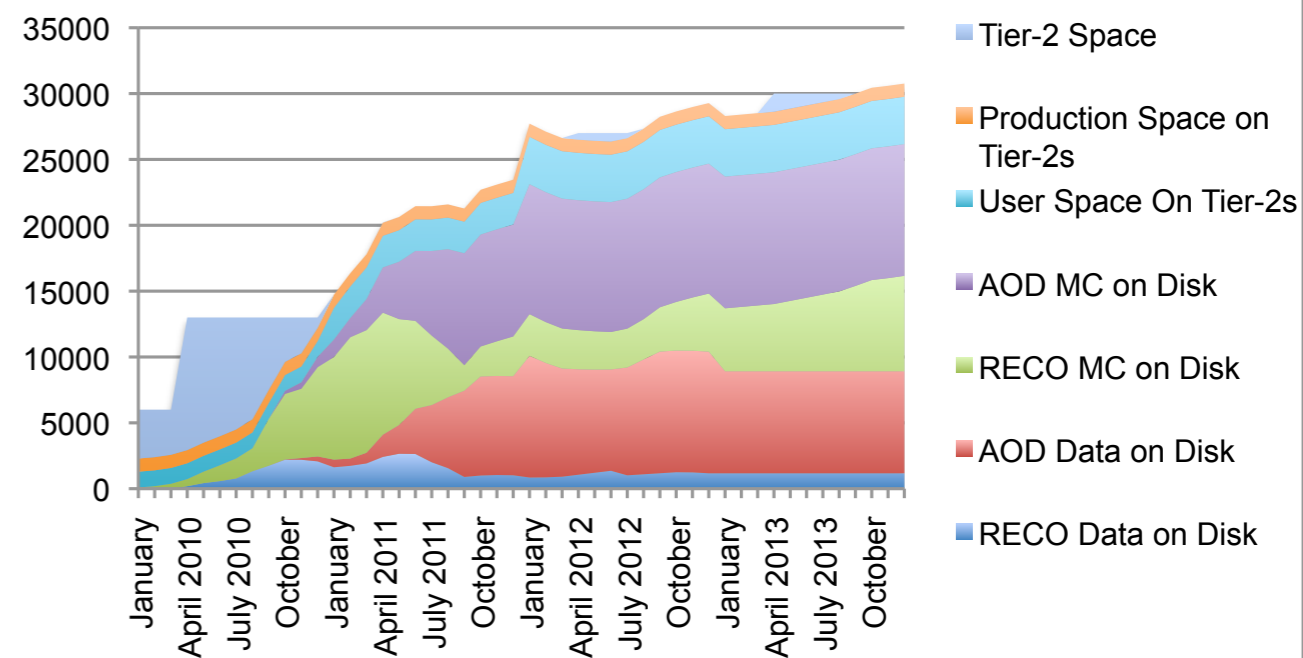


- ▶ Move MC production to T1 when possible, but when T1 is reprocessing, T2 CPU is heavily used
- ▶ 90% of user analysis needs to move from RECO to AOD
- ▶ Currently keep four copies of each analysis dataset across all 50 T2's, but need to be prepared to cut back
- ▶ Under any conditions, T2 resources are heavily committed

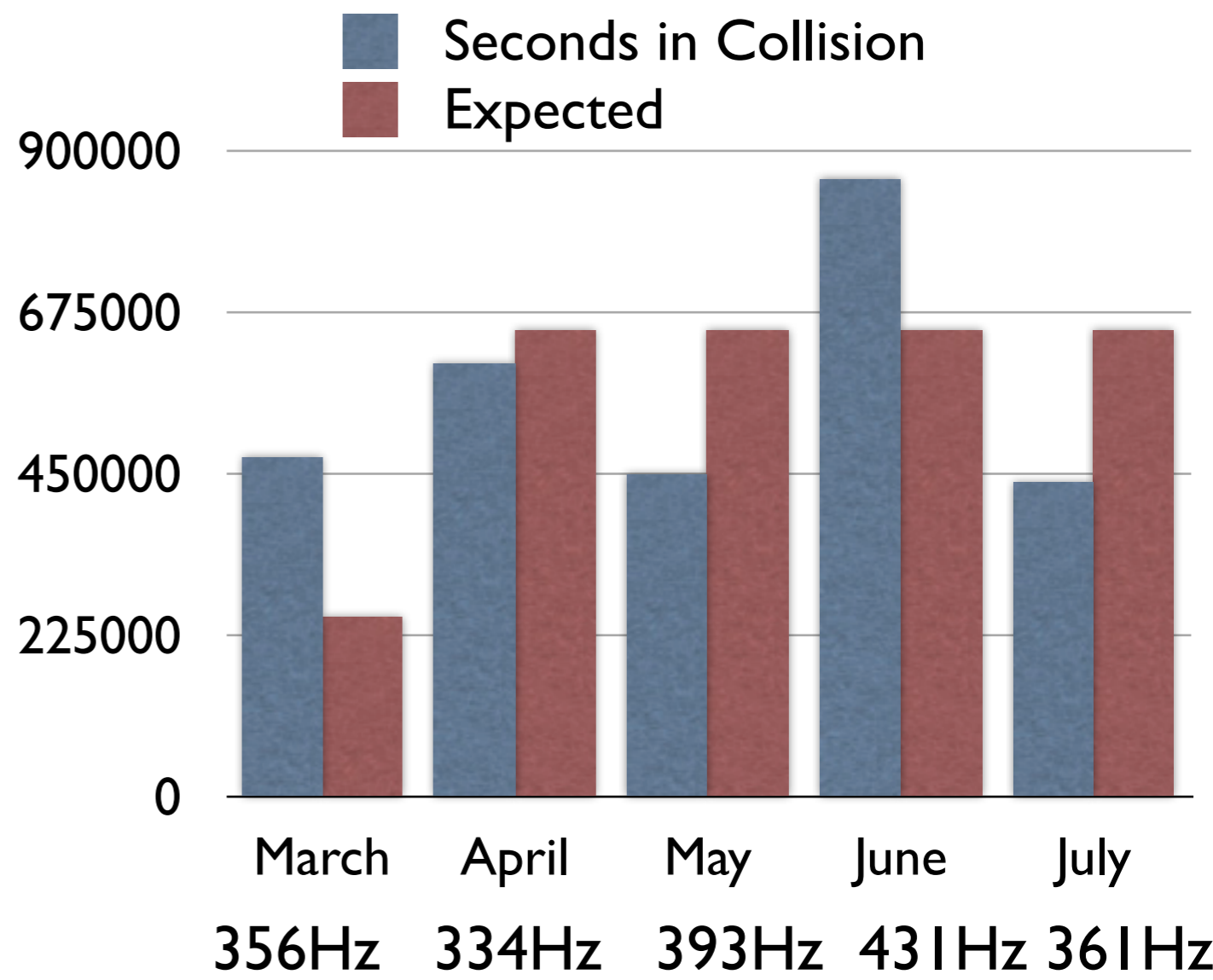
Tier-2 CPU Needs



Tier-2 Disk Storage



- ▶ How well does real CMS life match up with the plan?
- ▶ LHC duty cycle lower than anticipated, but CMS trigger rate above 300 Hz.
- ▶ Trigger rate includes overlap in primary datasets, planned to be 25%
- ▶ Recorded 1.1B events, compared to 1.3B in the planning
- ▶ Small contingency gained
- ▶ Re-reconstruction of full 2011 data should be ~1 PB

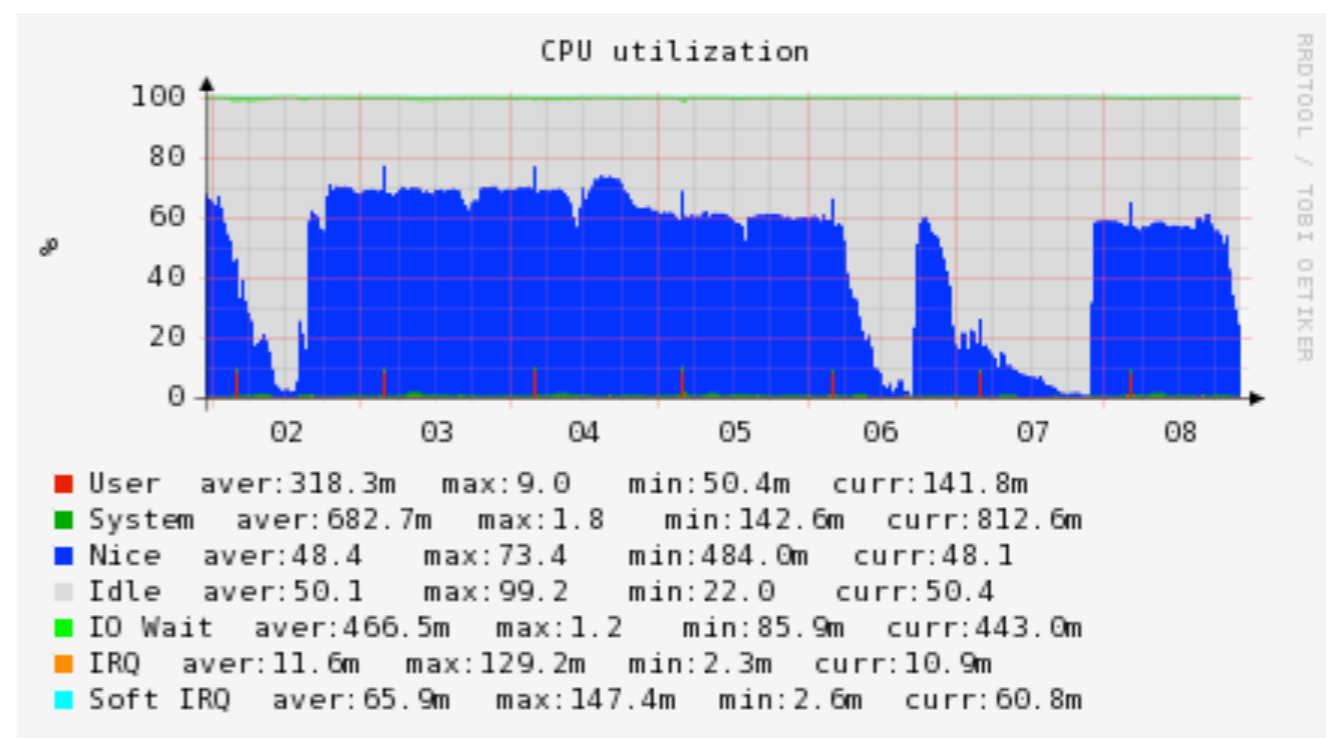
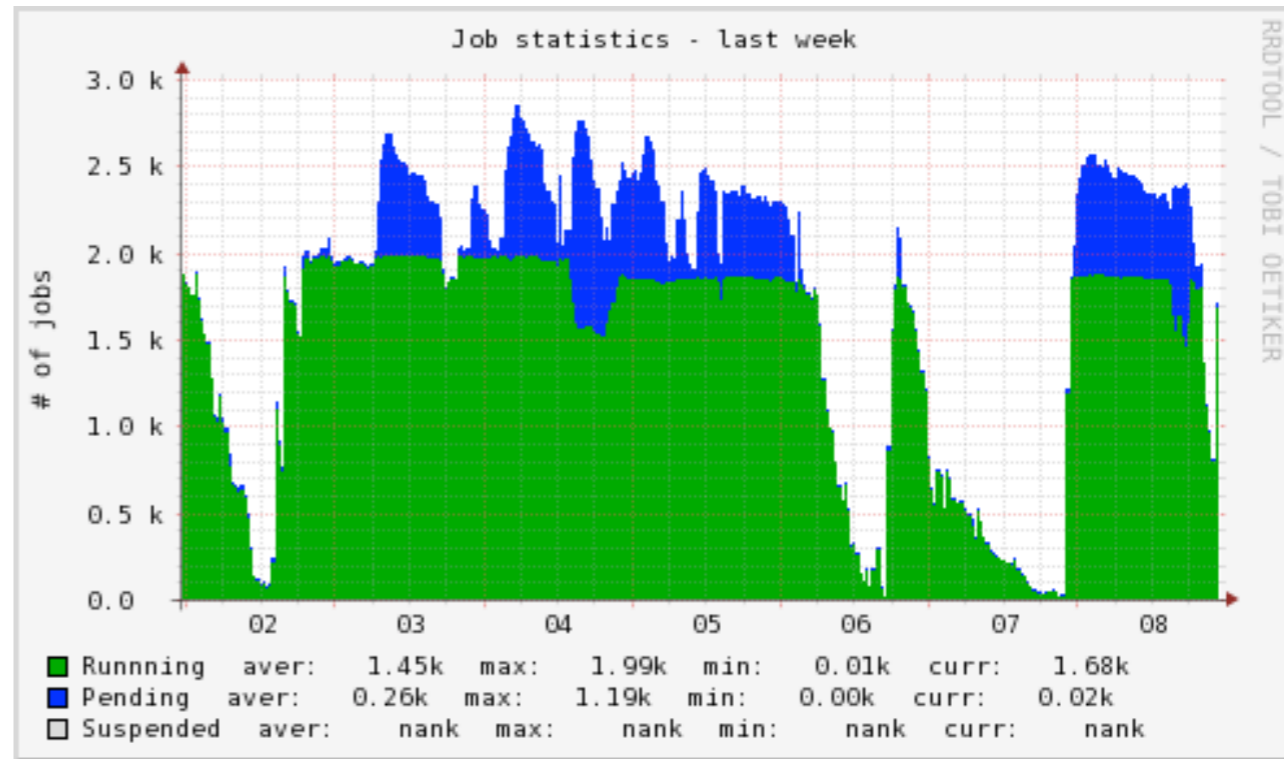


Average Trigger Rate

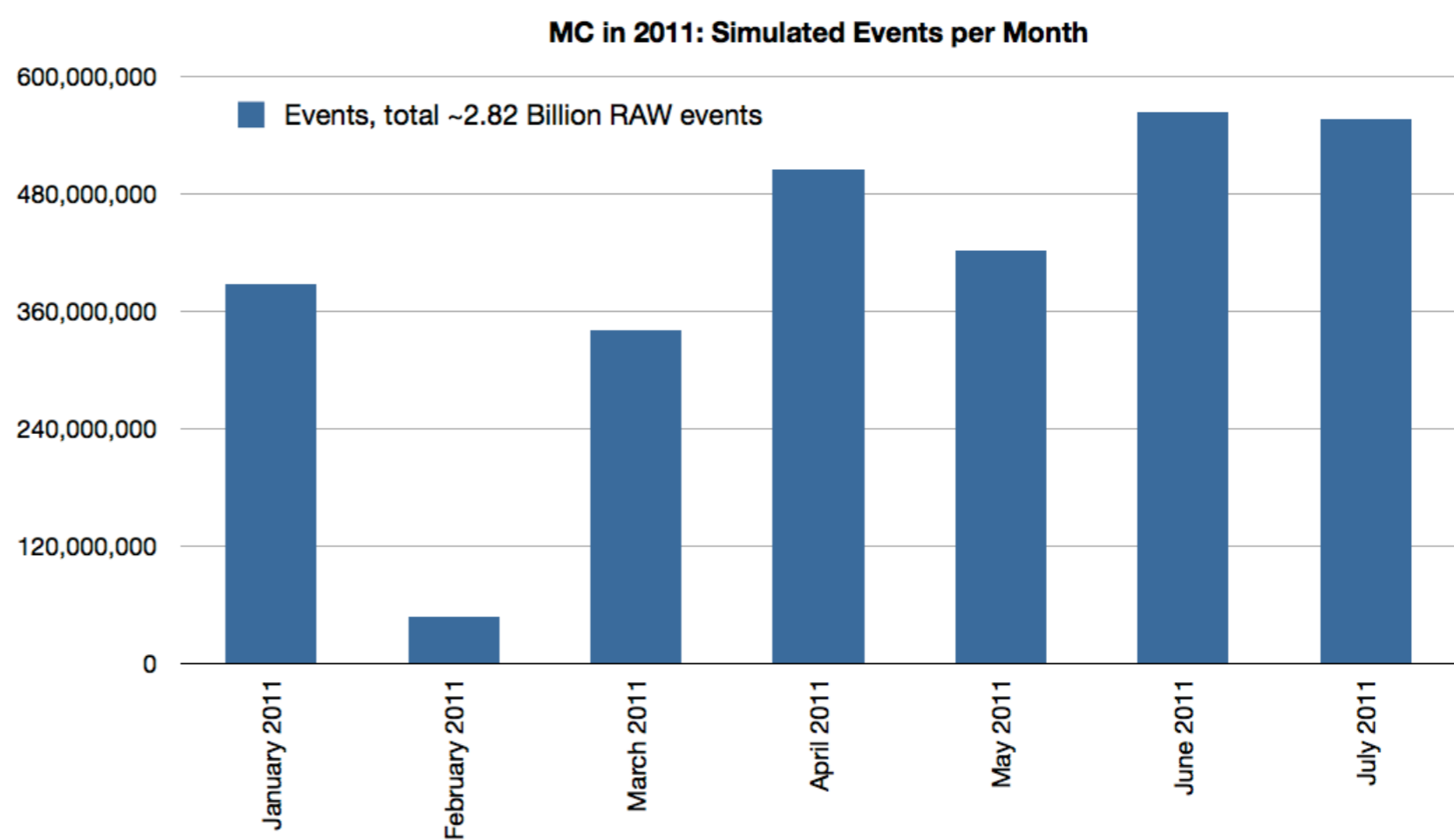
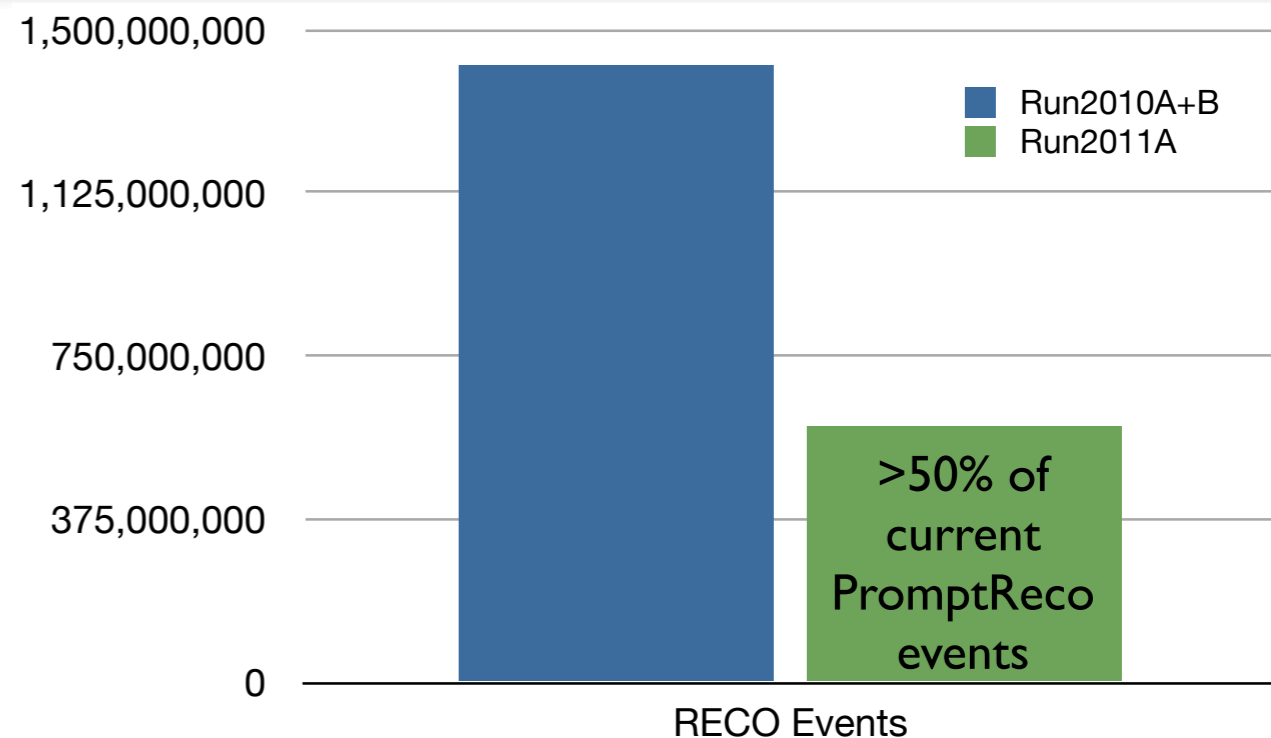
- ▶ In general pileup has been lower than anticipated due to larger number of smaller proton bunches
- ▶ Processing time about as expected for min-bias, 20% more than planned for other datasets
- ▶ Event sizes are generally smaller, have been roughly constant over time so far
- ▶ But everything expected to get bigger/longer as beam currents increase later this year

| Tier | Size | Expectation |
|-----------|-------|-------------|
| Data RAW | 200kB | 390KB |
| Data RECO | 500kB | 530kB |
| Data AOD | 100kB | 200KB |
| MC Reco | 970kB | 600kB |
| MC AOD | 250kB | 265kB |

- ▶ 40% LHC livetime in early June led to saturation of T0
- ▶ But could not fully use CPU:
 - ▶ Switch to 64-bit and new ROOT gives large memory footprint
 - ▶ Working to reduce exe size, take advantage of whole-node scheduling for shared read-only memory across multiple reconstruction jobs
- ▶ However, keeping up well enough with incoming data

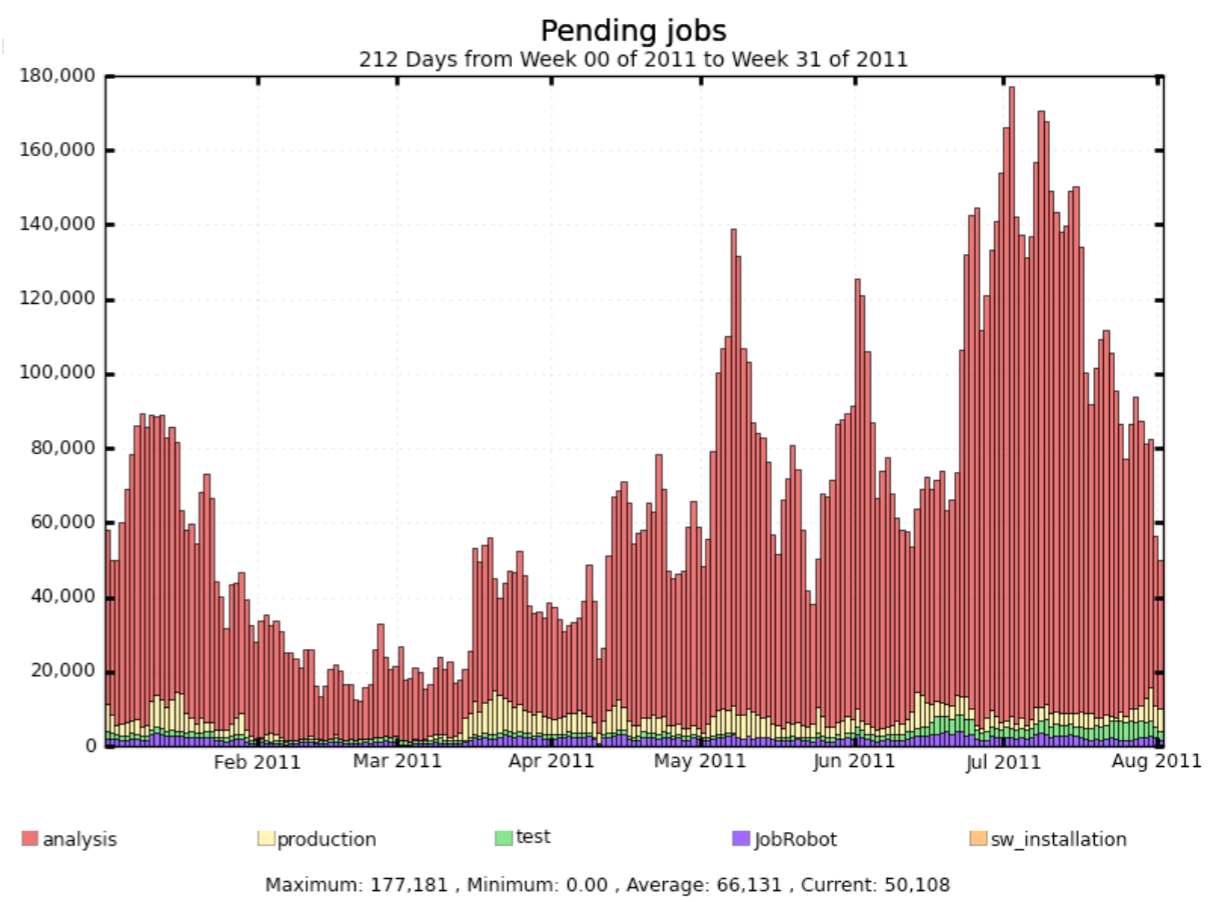
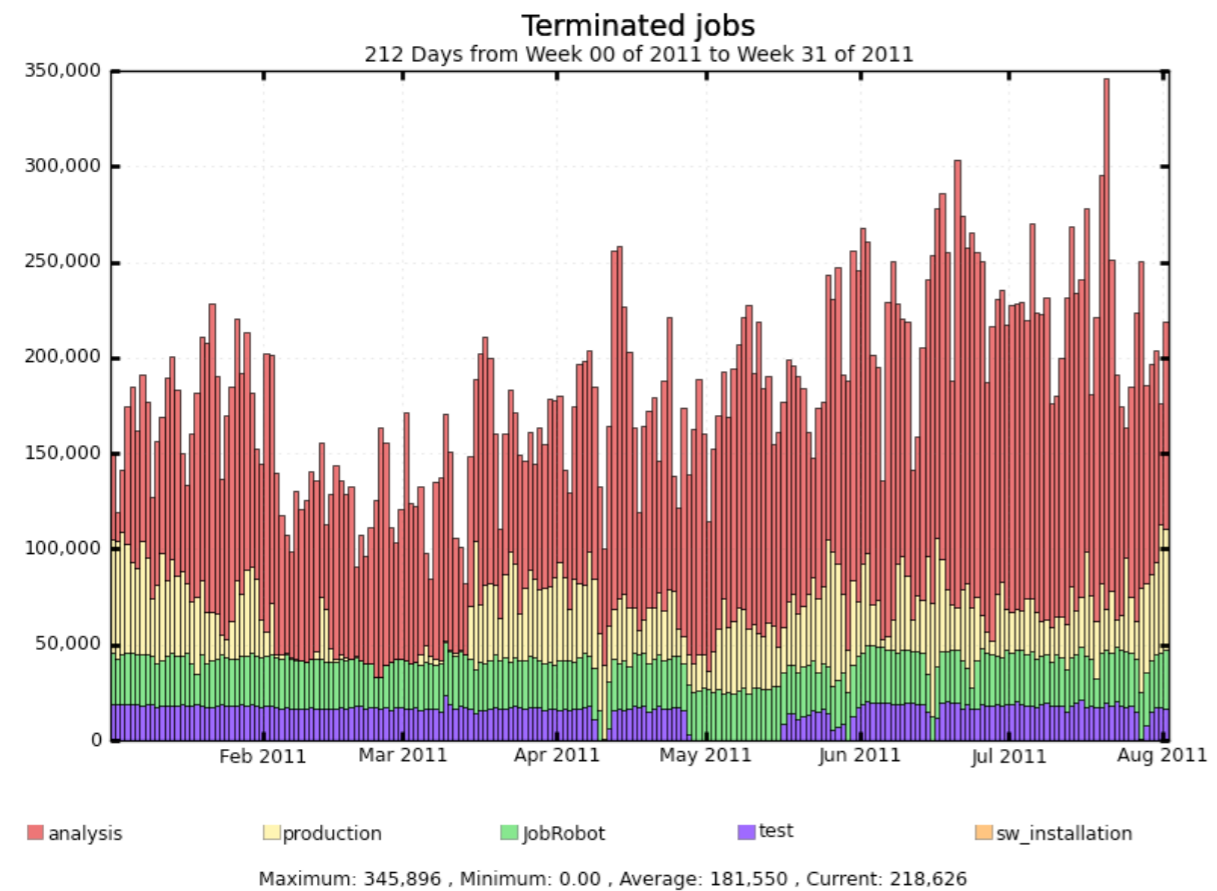
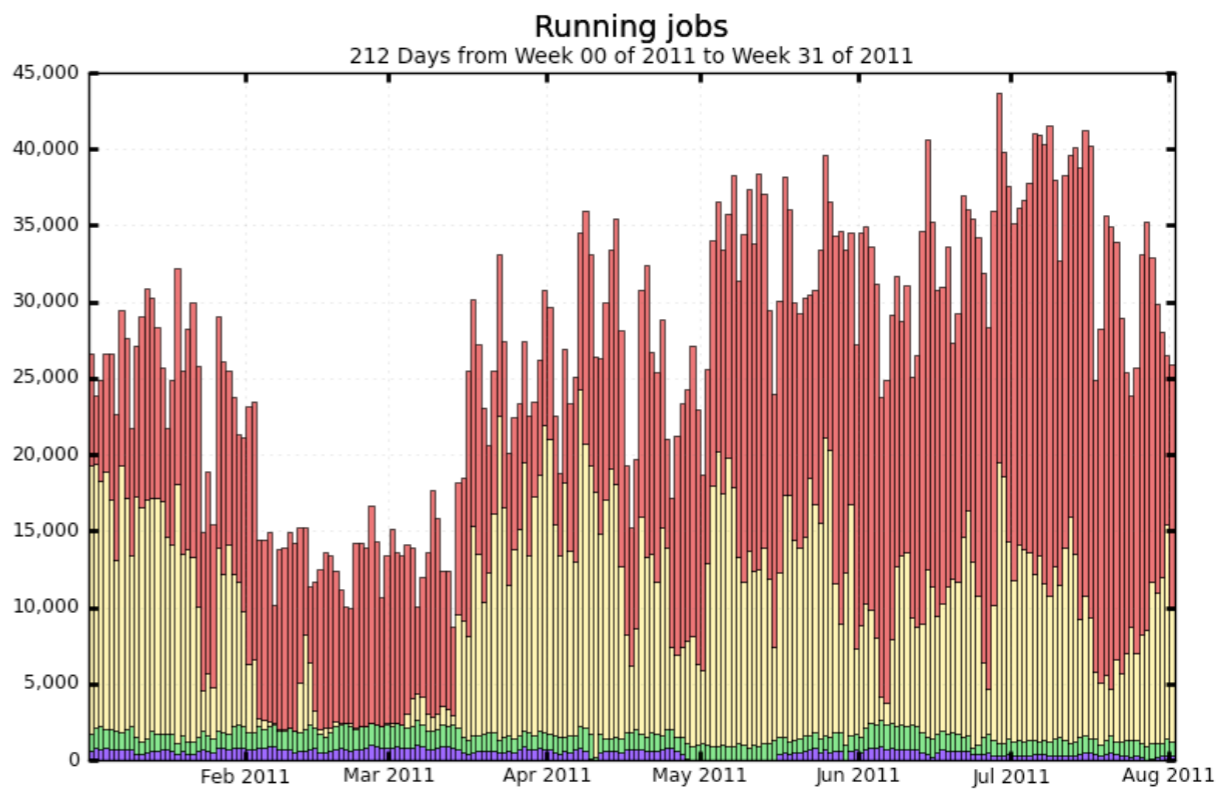


- ▶ Did full re-reco pass of 2010 data in April and all available 2011 data in May
- ▶ Consistent with planning
- ▶ Might not do full re-reco again until end of 2011 LHC run
- ▶ 2.8 billion MC events produced in 2011
- ▶ Latest simulation includes out-of-time pileup
- ▶ Had planned on 0.22B/month, in fact capable of much more

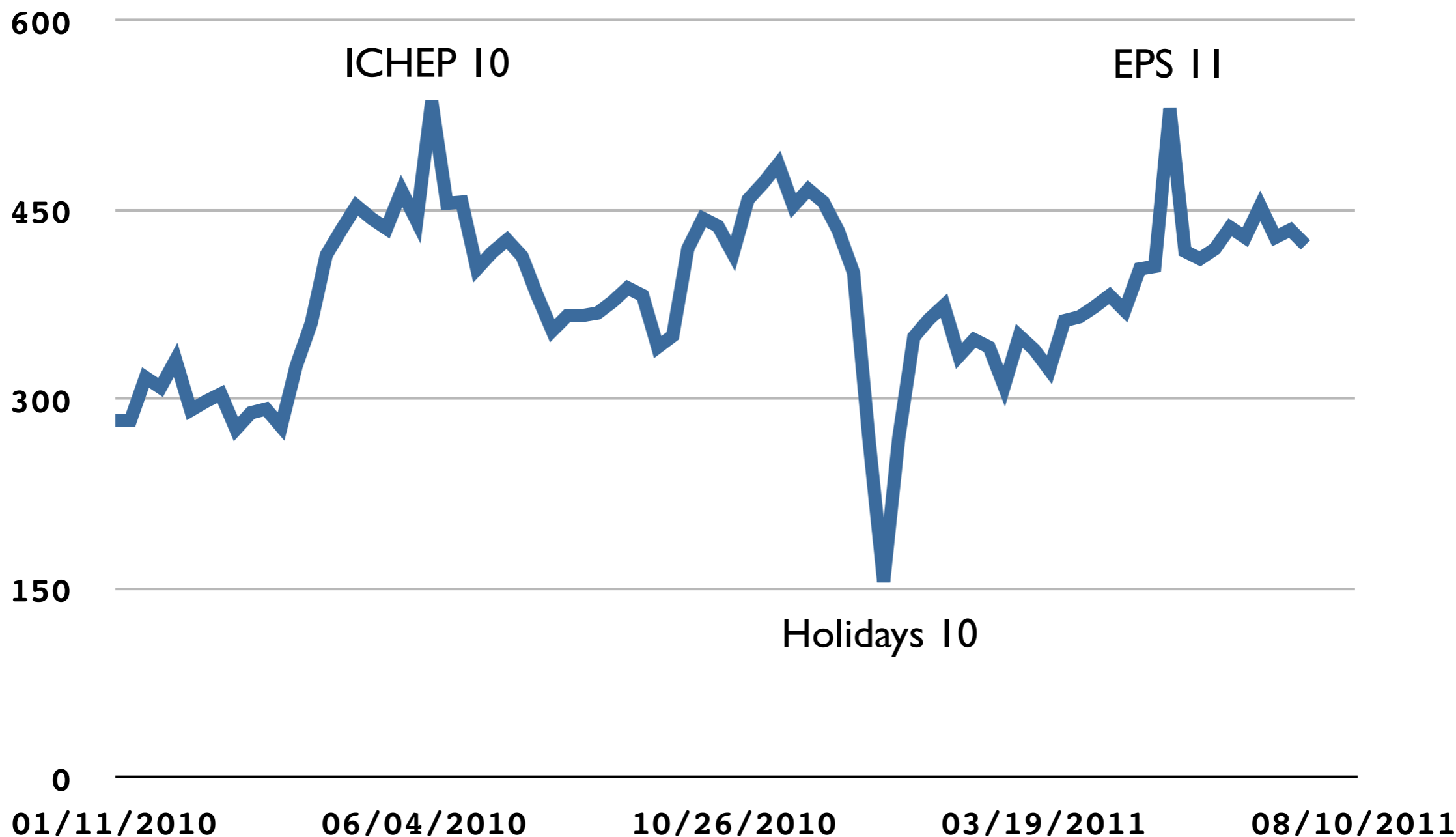


- ▶ Workflow management system for data re-processing was designed for MC production
 - ▶ OK if you lose some MC events, can always make more
 - ▶ Not OK to be losing data events!
- ▶ New WMAgent system is much more robust
 - ▶ State machine rather than messaging system
 - ▶ 100% accountability for all events processed
 - ▶ Current version of software uses more memory than before, jobs are running longer, more failed jobs
 - ▶ But WMAgent can re-do failed jobs straightforwardly
 - ▶ Has also allowed for more efficient MC production
- ▶ Whole-node scheduling at TI's will also bring operational efficiencies, aim for 50% of resources used this way by end of year

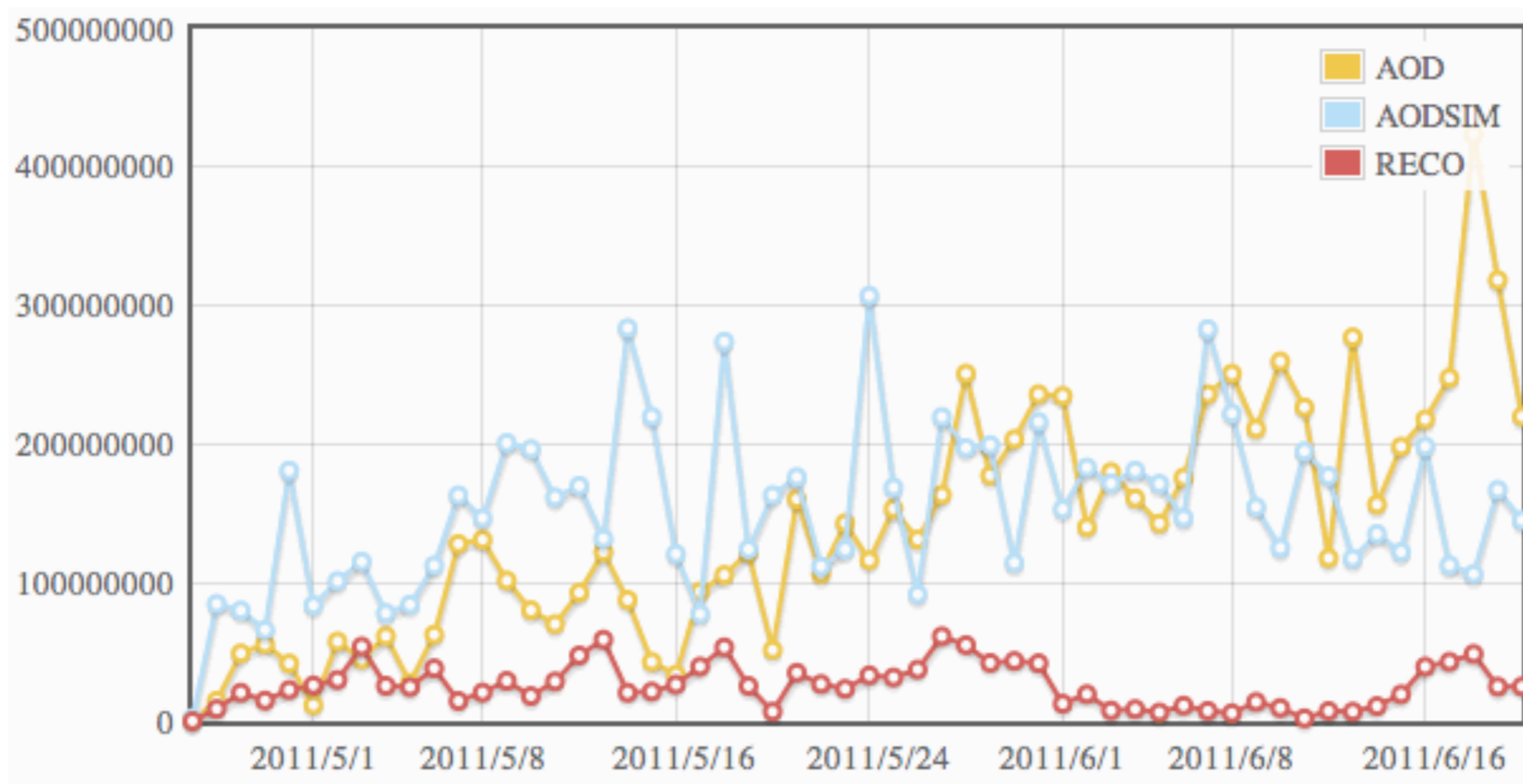
- ▶ 30K cores for analysis, continually
- ▶ More MC is moving to T1
- ▶ ~250K analysis jobs/day
- ▶ More than original computing model
- ▶ Still, many jobs pending....



► User community steadily growing, with a significant fraction of the entire collaboration (800 unique users/month) making use of grid resources for analysis



- ▶ Now have improved ability to track dataset usage
- ▶ Users are in fact making necessary transition from RECO to AOD

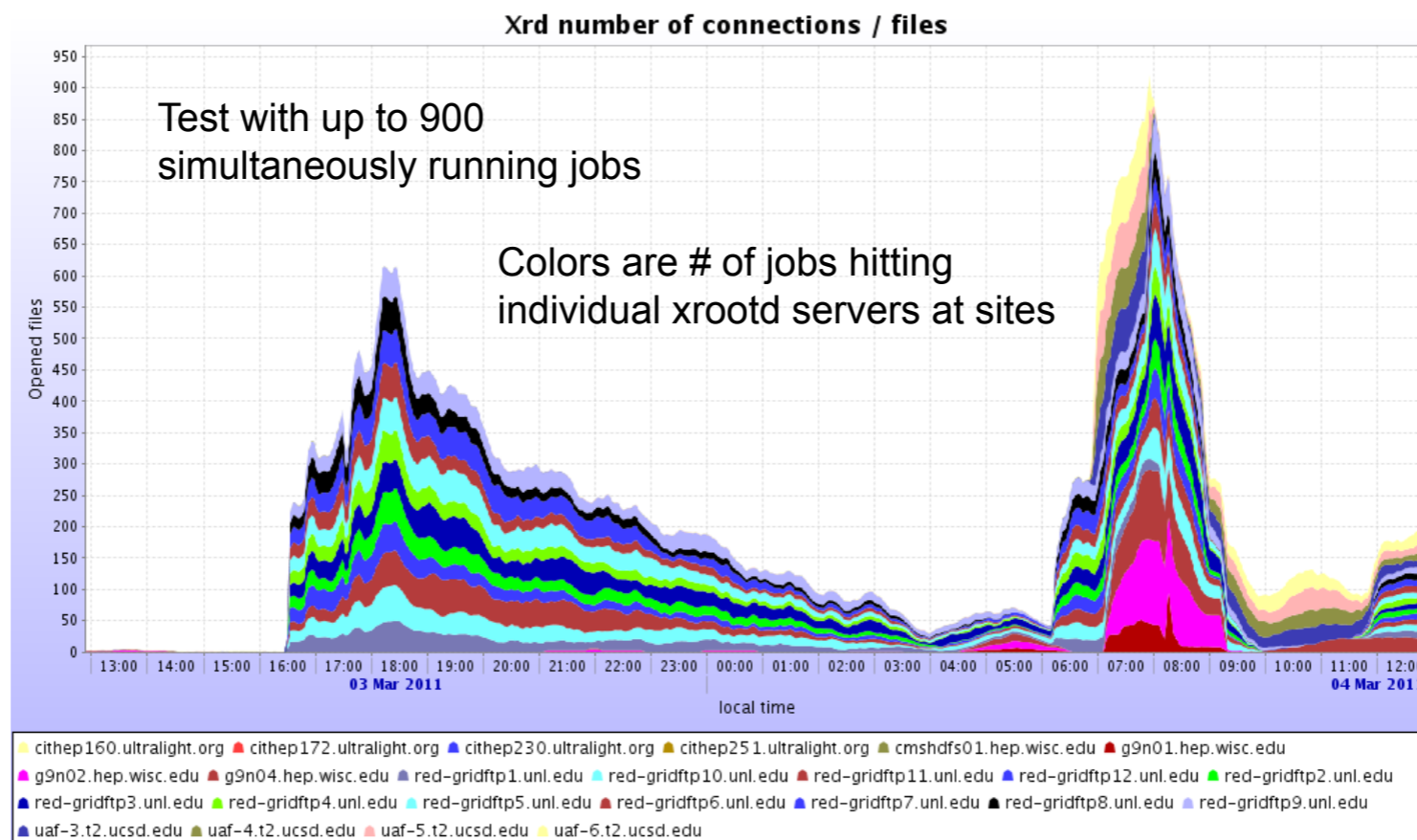


- ▶ Some other experiment previously had these tracking tools
- ▶ Will now help us manage data distribution and more

- ▶ The CMS Remote Analysis Builder (CRAB) that analysts use to submit grid jobs will have a significant revision
 - ▶ Install WMAgent underneath to take advantage of its features
 - ▶ User interface will also change, requiring some user re-education
- ▶ Greater use of pilot jobs (glide-ins) for analysis
 - ▶ Could allow for a prioritization of user jobs across entire distributed system, not just at individual sites
 - ▶ Also potential for balancing usage across sites

- ▶ Key limitation of computing model: CPU and storage co-located
 - ▶ Must place the data where the processing resources are
 - ▶ Difficult to optimize, need to guess analyst preferences
- ▶ But WAN is more reliable than anticipated in the MONARC days
- ▶ And CMS has optimized reading data files over the network
- ▶ Forget co-location and think big -- what if you could analyze data in one place with a CPU that's some other place?
 - ▶ Data placement hardly matters anymore
 - ▶ Users insulated from storage problems at sites: if a file is corrupt at one site, failover to network and access elsewhere transparently
 - ▶ Enable users who don't have large storage systems: small clusters can still have access to any data in the world, "diskless T3"
 - ▶ Access experiment data using cloud resources?

- ▶ Prototype systems for this have already been deployed
- ▶ Key element: redirectors that allow jobs to find data at remote sites
- ▶ Fallback to WAN access already enabled at US T2 sites
- ▶ Need to test/operate at scale, develop monitoring/accounting/throttling systems



- ▶ In related work, exploring how to migrate jobs between sites to optimize use of processing resources

- ▶ The actual use of CMS computing resources is largely in line with the model that was created based on 2010 experience
 - ▶ Some parameters higher or lower, within about 20%, but variances have tended to compensate each other
- ▶ The model predicts that CMS will be limited by its computing resources during this year. Some of this is being seen:
 - ▶ Some analyses slowed by wait for MC samples
 - ▶ Significant demand for processing resources at T2
- ▶ If CERN runs LHC at very high luminosity, could get worse
 - ▶ But mini-Chamonix workshop says it will be gradual....
- ▶ Physicists will need to adapt to this new environment
- ▶ Good news: the resource limitations reflect the fact that LHC datasets are growing rapidly and provide the opportunity for new physics discovery

- ▶ 2010 was an extremely good year for CMS computing
 - ▶ Computing was a strategic asset for producing physics results
 - ▶ Perspective: scales that were bleeding edge a few years ago are now every-day operations
- ▶ Strong performance has continued in 2011, but CMS has now entered an era of resource constraints
- ▶ But continuing technology developments are giving some operational breathing room
- ▶ Some of these developments have the potential to change the paradigm of computing at the LHC, and of data-intensive, high-throughput computing in general.