

CMS Computing Operations

OSG All Hands Meeting
19. March 2012

Oliver Gutsche
for
Computing Operations



- ▶ CMS Computing Operations was re-organized Beginning of 2012
- ▶ Reflect the change from **Commissioning** to **Operations and Optimization** of the **Computing systems and infrastructure**
- ▶ Data Operations & Facilities Operations were merged → Computing Operations
- ▶ Lead by: Markus Klute (MIT), Oliver Gutsche (FNAL), Pepe Flix (PIC)

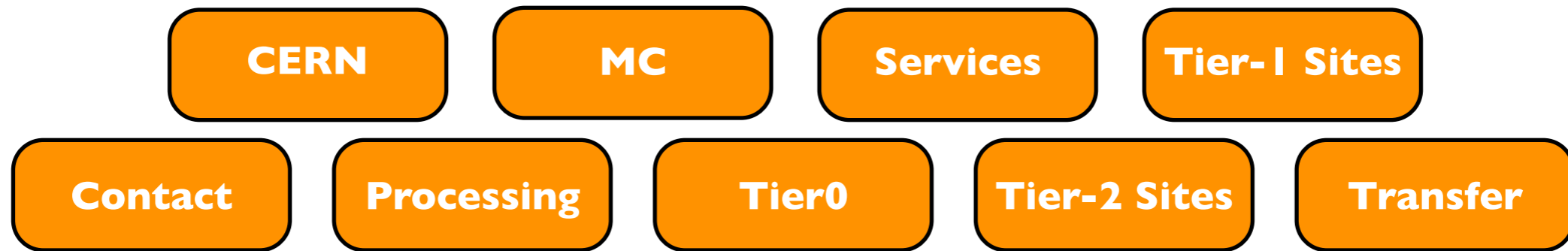


Leadership

Experience

L2

L3



CORE

Tier-0

Infrastructure

Workflow

Transfer

Monitoring

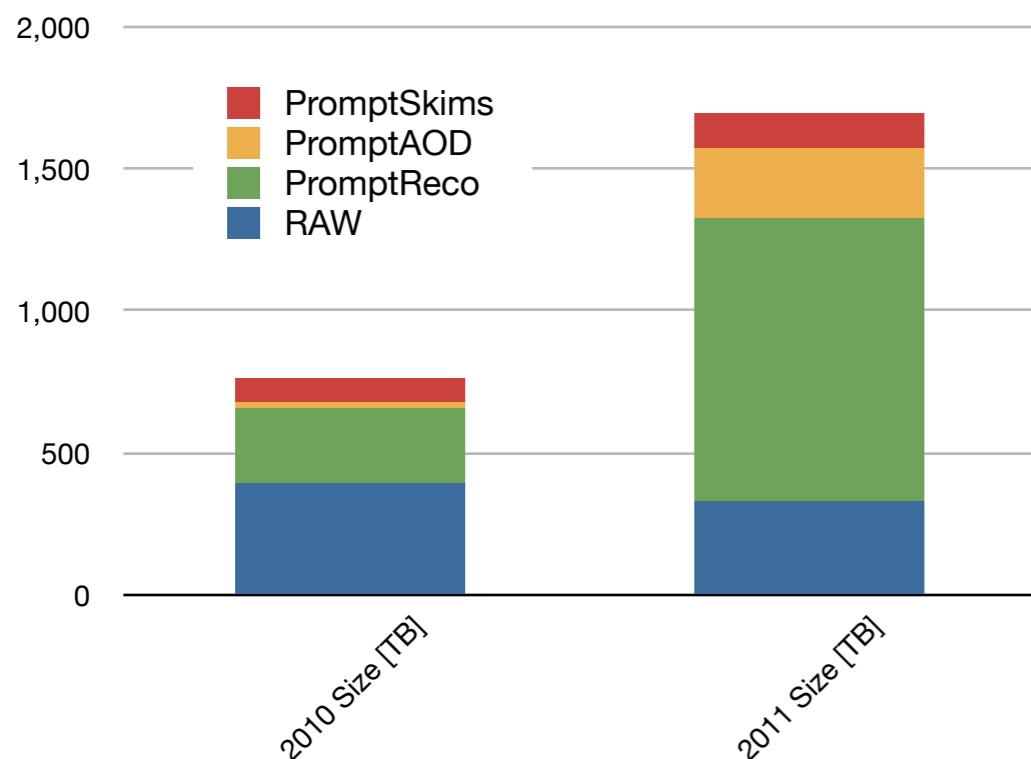


- ▶ Promising start of new project with a lot of gained synergy
 - ▶ Focus is on communication between the different parts of the Computing Project and to other projects of CMS
- ▶ Still places where we can improve processes and smoothen out operational procedures

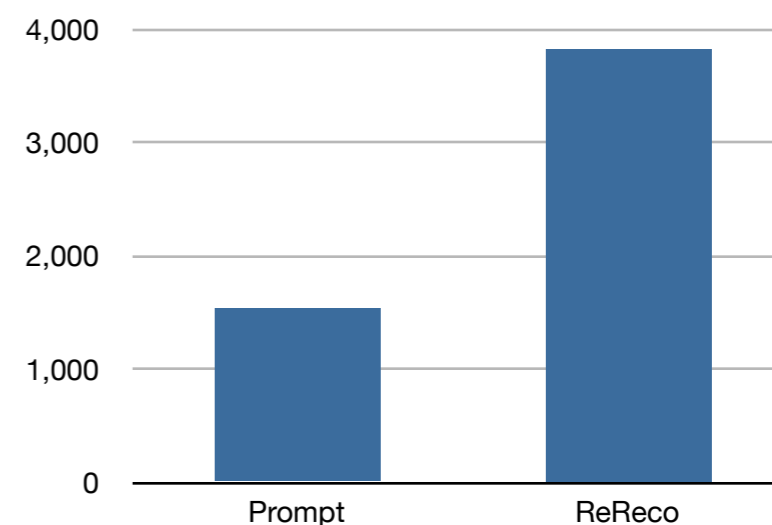
- ▶ Team is not complete, following positions are not filled:
 - ▶ L3_CERN (1): service and infrastructure coordination at CERN
 - ▶ L3_Services (2): CMS global services and infrastructure coordination
 - ▶ L3_Tier I Sites (1): liaison to Tier-I sites
 - ▶ L3_Transfers (1): Transfer monitoring and trouble shooting

- ▶ Please contact us if you are interested in helping us in these tasks

Data: Size



Million AOD Events in 2011

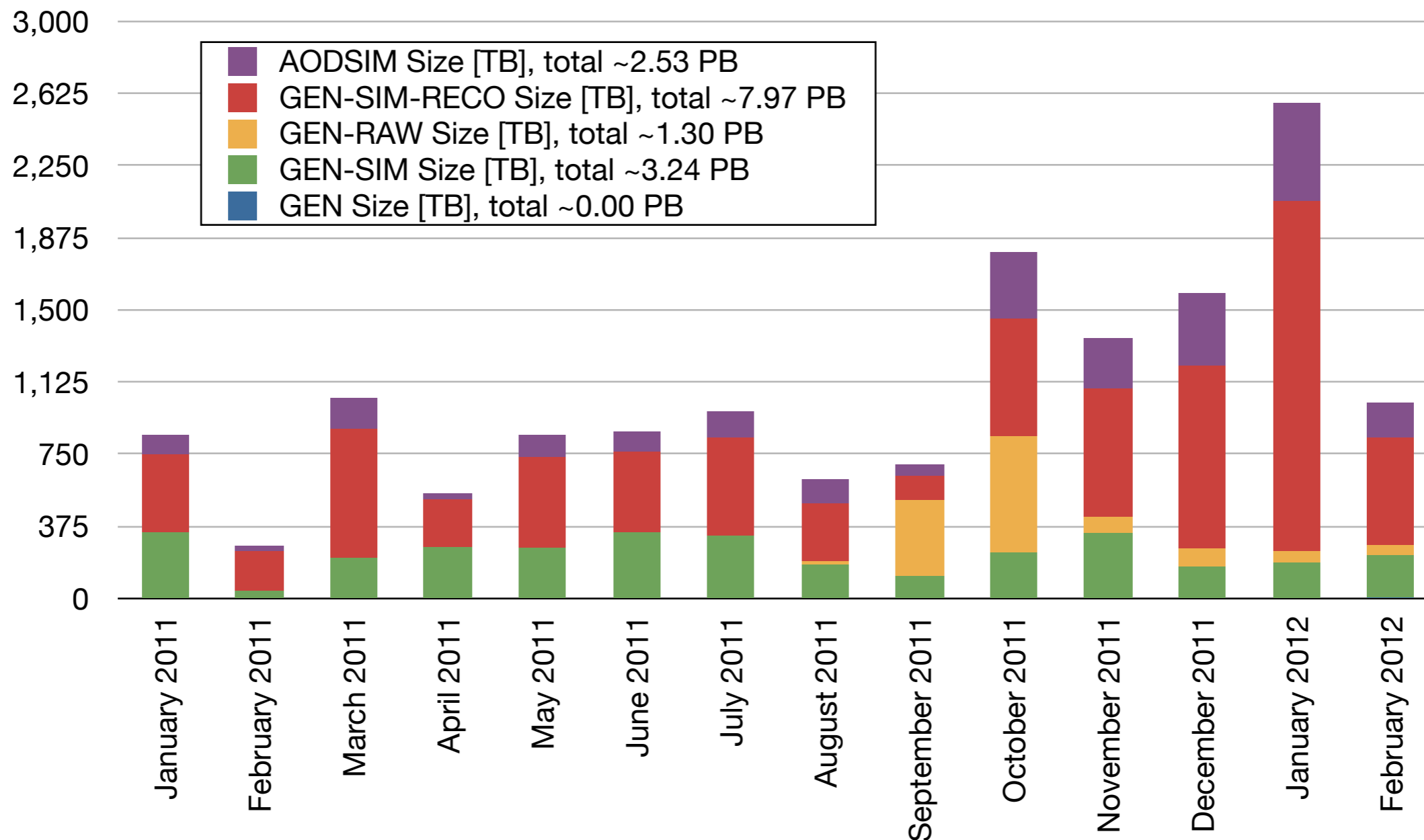


	RAW	PromptReco	PromptAOD	PromptSkims	Total
2010 Size [TB]	393	266	19	87	765
2011 Size [TB]	331	995	247	123	1,696
Total Size [TB]	724	1,261	266	211	2,462

	RAW	PromptReco	PromptAOD	PromptSkims
2010 [Million Events]	1,536	1,186	184	222
2011 [Million Events]	1,535	1,525	1,443	284
Total [Million Events]	3,071	2,711	1,627	506

- ▶ Prompt data: 1.7 PB
- ▶ 2011 re-reconstruction passes
 - ▶ 27 partial
 - ▶ 1 complete
 - ▶ In total, re-reconstructed more than twice the events we recorded (once during the year, one time at the end of the year)

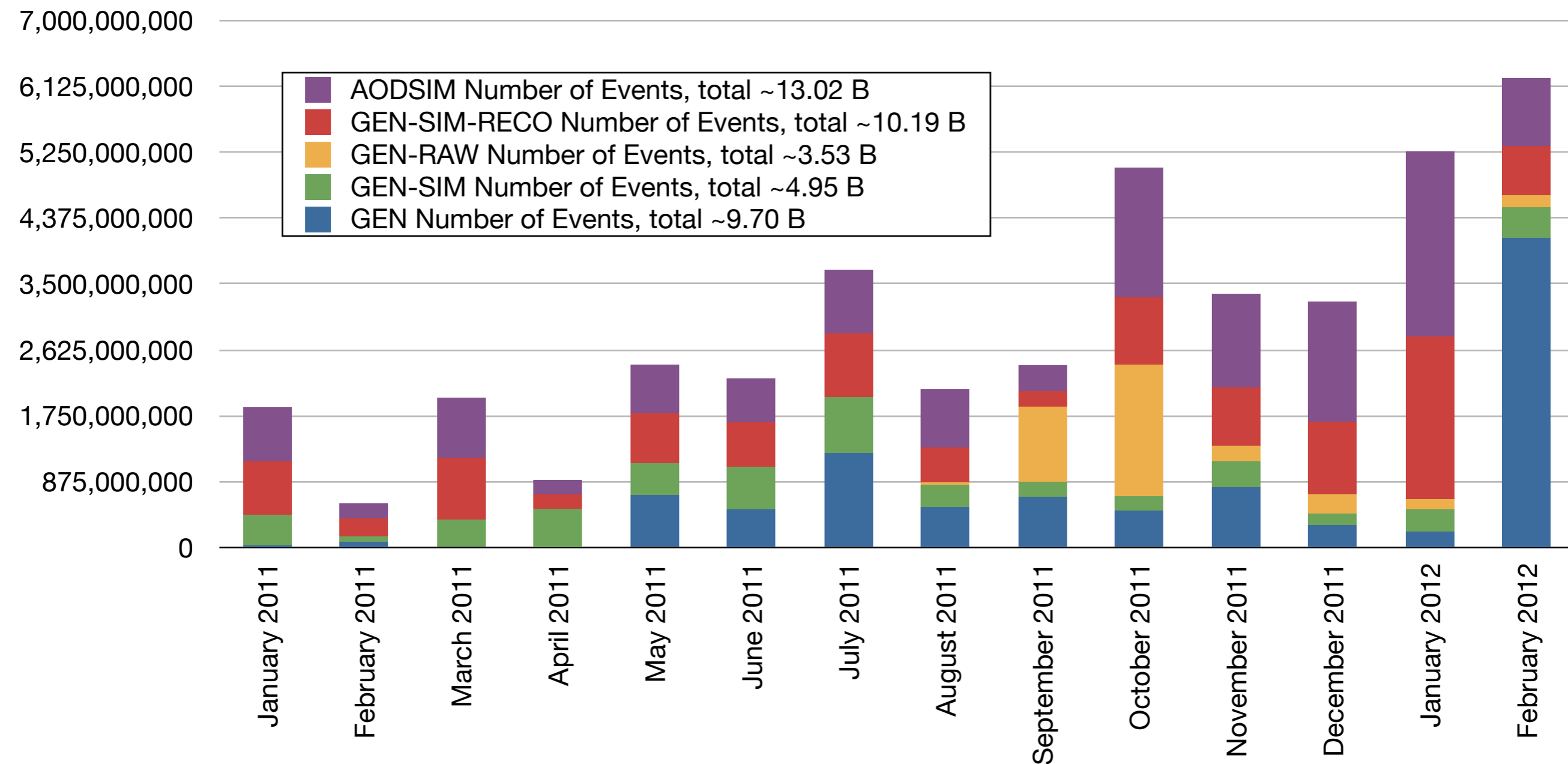
MC in 2011/2012: Size in TB per Month



▶ Total in 2011: 11.5 PB

▶ Total in 2012: 3.6 PB

MC in 2011/2012: Number of Events per Month



▶ 2011 GEN-SIM: 4.3 Billion → 2011 AODSIM: 9.7 Billion

▶ Re-reconstructed and partly re-digitized every GEN-SIM event more than twice

▶ <https://cms-popularity.cern.ch/>

▶ Statistics about dataset access through CRAB since June 2011

▶ Analysis since June 2011:

▶ **Summer11_R1**: CMSSW_4_2_X & PU_S4

▶ **Fall11_R2**: CMSSW_4_2_X, PU_S6

▶ **Fall11_R4**: CMSSW_4_4_X, PU_S6

▶ Datasets without any access:

	Summer11_R1		Fall11_R2		Fall11_R4	
	GEN-SIM-RECO	AODSIM	GEN-SIM-RECO	AODSIM	GEN-SIM-RECO	AODSIM
Valid Datasets in DBS	2160	2182		1651	2612	2636
Datasets accessed once	240	1946		799	3	118
Datasets not accessed at all	1920	236		852	2609	2518

▶ Observations:

Details: <https://hypernews.cern.ch/HyperNews/CMS/get/comp-ops/76.html>

▶ AOD switch worked very well! according to the RECO numbers

▶ Sizable number of AOD datasets not accessed at all, also in Fall11_R2

▶ Caveat:

▶ Only looked at number of datasets, not at number of events (could be that many small datasets have not been accessed)

▶ **But: Popularity Service very powerful, also for site admins!**

Data taking rate:

- ▶ 300 Hz Prompt reconstructed at Tier-0

- ▶ ~300 Hz additional reconstructed at Tier-I sites

- ▶ 300 Hz parked and reconstructed in 2013

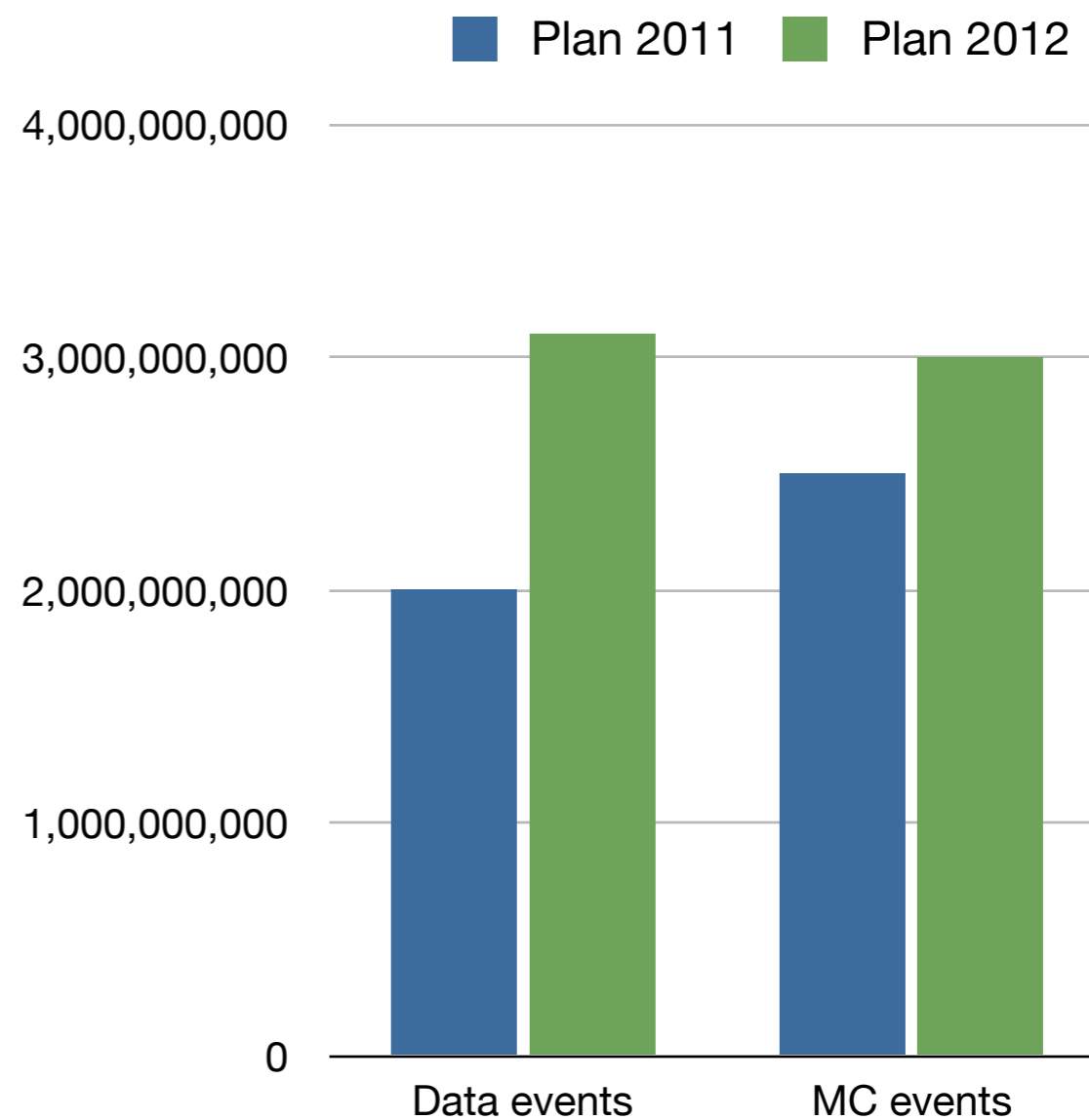
- ▶ **Important:** these are the planning numbers, the current estimates of used bandwidth are somewhat lower.

- ▶ 2012 will see more than 1 1/2 times the data events of 2011

- ▶ Impact for analysis

- ▶ 2012 will see a physics MC sample of 3 Billion events compared to 2.5 Billion in 2011

	Plan 2011	Plan 2012
Data taking [s]	5,200,000	5,200,000
Data taking rate [Hz]	300	600
Overlap	1.25	1.25
Data events	2,000,000,000	3,100,000,000
MC events	2,500,000,000	3,000,000,000



Data re-reconstruction

- ▶ Assume that small re-reconstruction passes like in 2011 will re-reconstruct the whole 2012 dataset
- ▶ No plan for End-Of-Year re-reconstruction pass

MC re-digitization/re-reconstruction

- ▶ 2011 saw 2 complete MC re-digitization/re-reconstruction campaigns
- ▶ 2012 will only see one complete campaign

PileUp will increase from 16 to 30 PU events

- ▶ Large impact on data reco and MC re-digitization/re-reconstruction times
- ▶ Small impact on AOD analysis times



	Plan 2011	Plan 2012
PileUp	16	30
RECO Time Data (HS06s)	92	280
Re-digi/Re-RECO Time MC (HS06s)	164	400
RECO Analysis Time	35	50
AOD Analysis Time	11	12

▶ Christmas 2011/2012

- ▶ Produced GEN events using higher order Generators: 4.5 Billion events

▶ February 2012:

- ▶ Started MC Simulation in CMSSW_5_0_X

▶ Mid March 2012:

- ▶ Started MC Digitization/Reconstruction in CMSSW_5_1_X

▶ April 2012:

- ▶ Switch Tier-0 to CMSSW_5_2_X
- ▶ Re-start MC Digitization/Reconstruction in CMSSW_5_2_X
- ▶ Goal: 1 Billion MC events for ICHEP conference analysis (taking place Beg. Of July 2012)

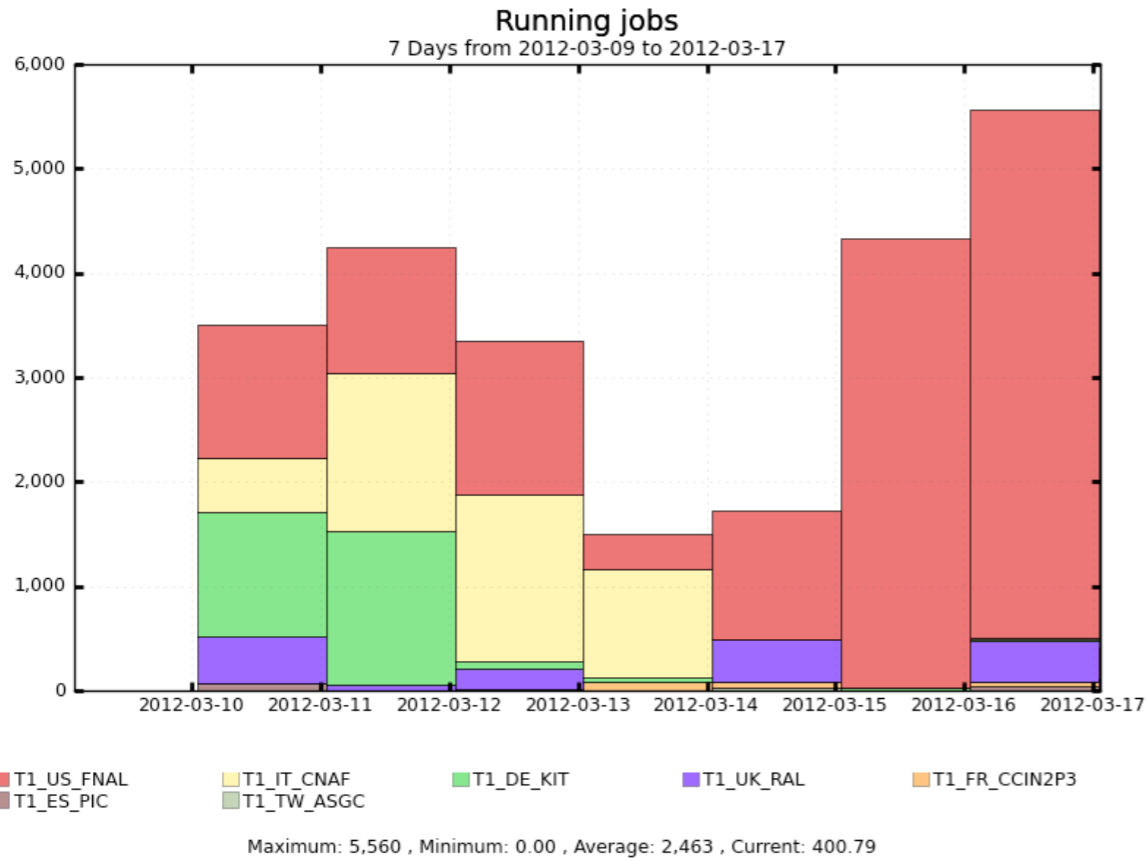
▶ May 2012:

- ▶ Possible checkpoint to decide if sufficiently large progress in physics performance of software was made, if yes:
 - ▶ Validate CMSSW_5_3_0 and deploy at Tier-0 during technical stop around ICHEP
 - ▶ After ICHEP, start 2012 data re-reconstruction
 - ▶ Decide if MC has to be re-digitized and/or re-reconstructed
 - ▶ Decide if 2011 data and MC has to be re-digitized and/or re-reconstructed

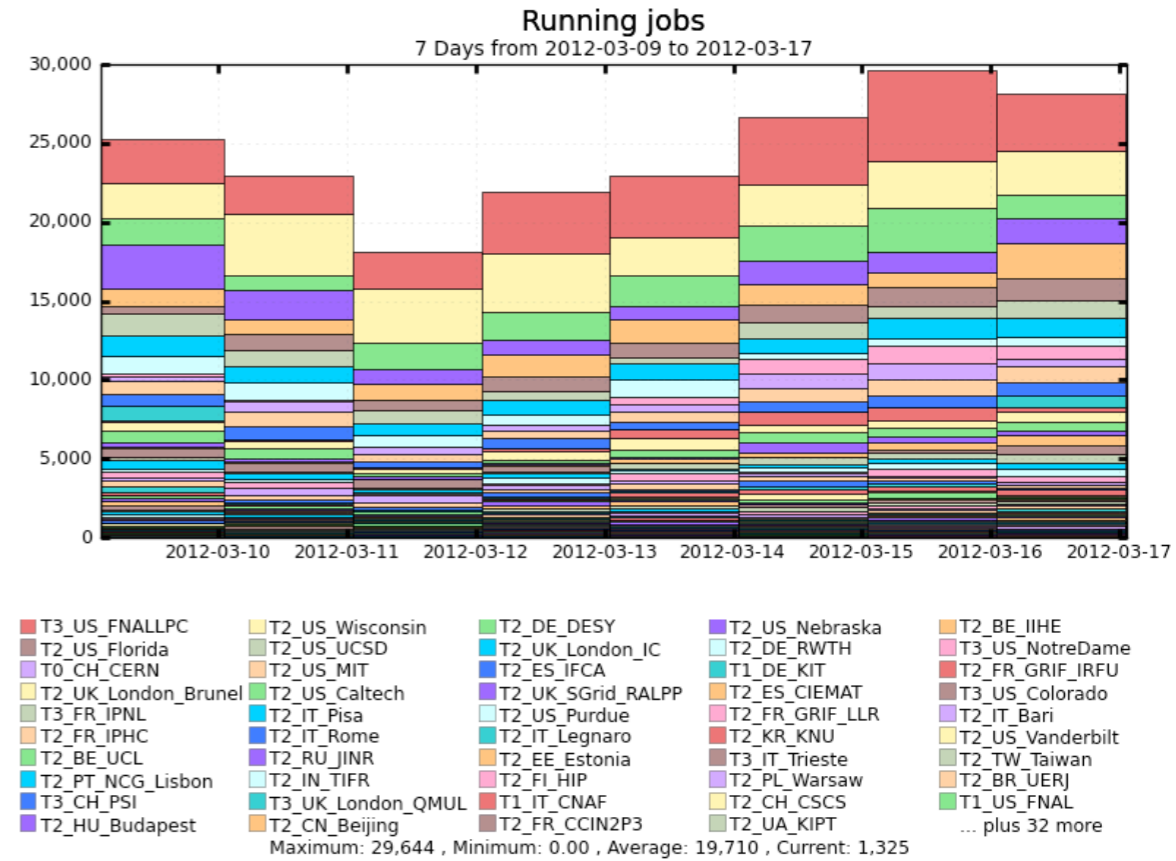
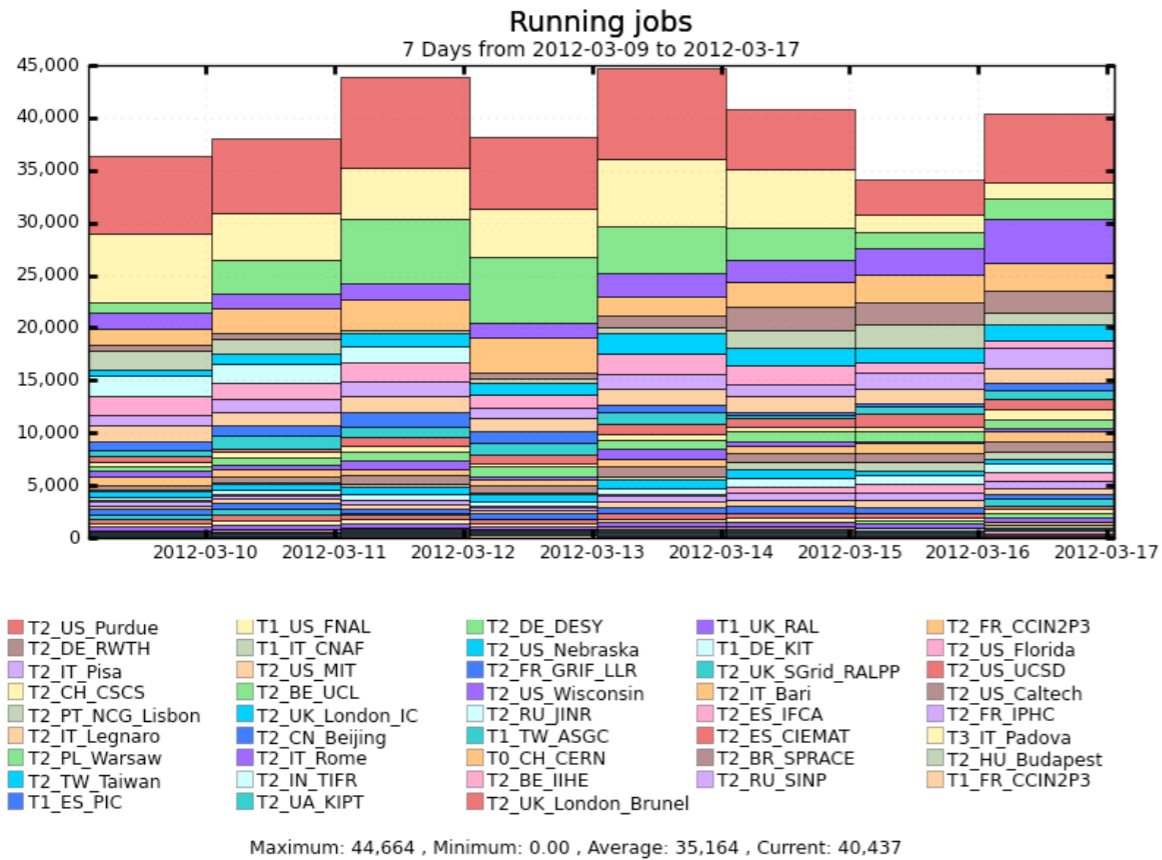
▶ Tier-2 sites can expect to see significant fluctuation of samples going in the site for analysis!

Tier-2 sites

T1 processing FNAL glideln WMS factory
Full scale: up to 20k running jobs



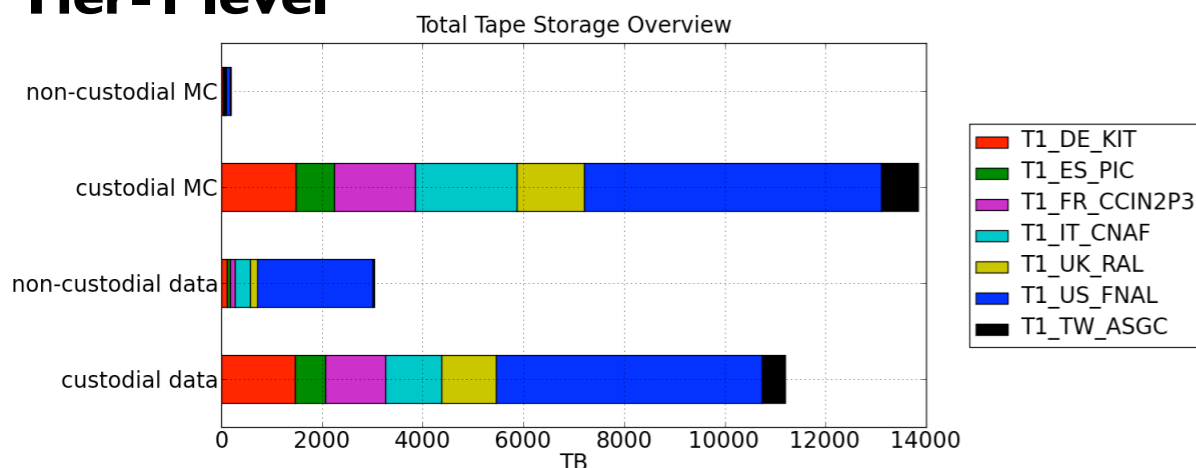
MC production CERN & UCSD glideln WMS factories
Reached: 45k running jobs on all tier levels



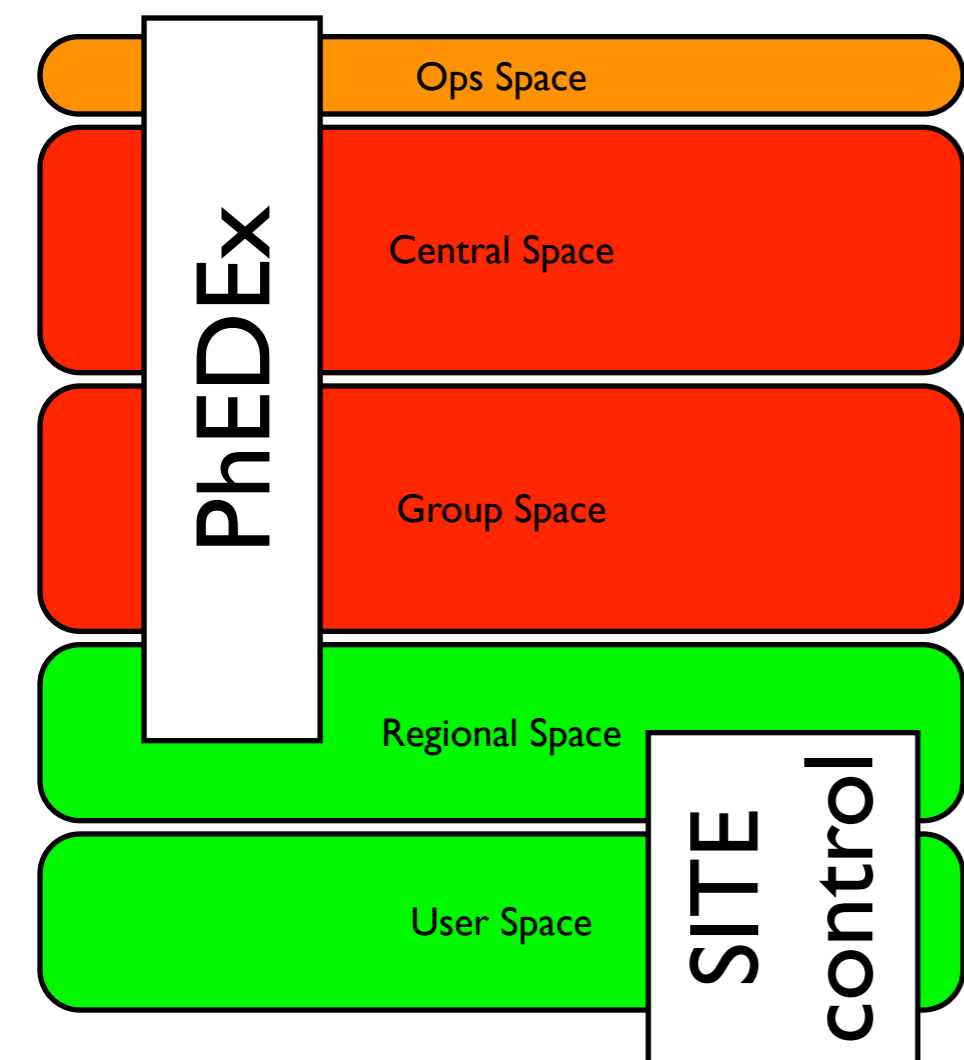
Analysis
gLite WMS & CERN & UCSD glideln WMS
factories
Routinely 30k++ running jobs

- ▶ glideln WMS in CMS
 - ▶ Used in Central Operations since Spring 2011
 - ▶ Since 2012 also MC production via glideln
 - ▶ Using for Analysis since the CRAB servers at UCSD came online in 2011
- ▶ gLite WMS
 - ▶ Used only for analysis
- ▶ 2012/2013:
 - ▶ Atlas will be moving to glideln WMS
 - ▶ WLCG TEG said that beyond 2012 very little gLite WMS submission
- ▶ Currently using 2 main VOMS roles important for Tier-2 sites
 - ▶ MC production
 - ▶ 50% of the resources of a Tier-1 site
 - ▶ Production role
 - ▶ Analysis
 - ▶ 50% of the resources of a Tier-1 site
 - ▶ Priority role is used to prioritize analysis jobs before jobs with no role
- ▶ Important for Tier-2 sites
 - ▶ Correct prioritization of these roles and following the respective resource percentages!

Tier-1 level



Tier-2 level

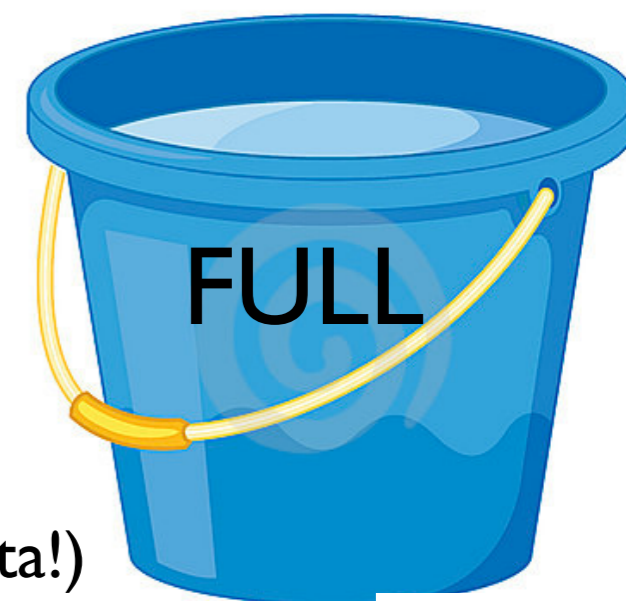


- ▶ Data Manager responsibilities very important
 - ▶ Ops, Central and Group space needs to be guaranteed
 - ▶ Regional and User Space have to be checked regularly and overfilling of the site has to be avoided
 - ▶ Data Manager is responsible that central PhEDEx spaces can be used at agreed quotas at all times

- ▶ CMS is working on tools to allow central accounting of used space outside PhEDEx control (Regional and User space)
 - ▶ Based on regular complete storage dumps at sites

- ▶ CMS is also establishing consistency checks for missing files and orphans
 - ▶ Tier-1 sites are checked every month
 - ▶ Pilot program with test sites on Tier-2 level running
 - ▶ Planned to have monthly consistency checks on Tier-2 level as well

- ▶ Every CMS user has at least one /store/user area where he can store ntuples through GRID access
 - ▶ Handled through regional organization, ask your institute leader where your /store/user area is hosted
 - ▶ Every Tier-2 sites supports a number of users (about 40 per US T2)
 - ▶ Usually 1-3 TB of space is allocated per user
- ▶ CERN:
 - ▶ CERN still provides tape backed user directories for CMS
 - ▶ /castor/cern.ch/user
 - ▶ Last year, CMS users wrote 2.7 PB (compared to 1.7 PB prompt data!)
 - ▶ Access to these stored files is impacting data taking and transfers
 - ▶ Files need to be staged from tape, available disk pool very small, significant activity impacting the whole system
 - ▶ CMS will work in the next month to transition users to official /store/user storage and close /castor/cern.ch/user
 - ▶ Expect increased demand for /store/user areas at your sites, CERN users can also qualify for space on EOS at CERN (T2_CH_CERN)
 - ▶ FNAL provides /store/user space (2 TB per user) for every US collaborator



- ▶ Software deployment
 - ▶ Transition to model where external software installation is not necessary anymore → access to new software releases is provided through self-updating systems
 - ▶ 2 solutions for sites: CVMFS or CRON installation
- ▶ CRON:
 - ▶ Site runs CRON job that installs new software releases automatically
- ▶ CVMFS (CERN virtual file system based on SQUID caches)
 - ▶ Site mounts CVMFS on all workernodes
 - ▶ Software is installed centrally at CERN and distributed through CVMFS to all sites
 - ▶ Status: Preparation of central CERN installation about to be migrated to final production hardware.
 - ▶ First Tier-2 sites in UK tests CVMFS served from the current installation base at CERN, then we expand testing to more sites and in the end migrate all sites to CVMFS or CRON

- ▶ Computing Operations reorganized beginning of 2012
- ▶ 2011 saw a lot of activity, Tier-2 sites backbone of CMS analysis and a vital part of the whole analysis chain
 - ▶ Thanks for exceptional performance in 2011!
- ▶ 2012 will be a busy year with even more data as 2011 and a lot of analysis activity
- ▶ Data manager role is very important to guarantee space availability and avoid overfilling of site storage
- ▶ /store/user areas become more and more important, especially when CERN Castor user area is closed
- ▶ 2012 will be very exciting, so expect a lot of activity of users!