

WLCG-LHCC meeting - CERN - 21 February 2017

CMS

Liz Sexton-Kennedy (Fermilab)

D. Bonacorsi (Univ. of Bologna / INFN)

(on behalf of CMS Sw/Comp)

Outline

Since last LHCC meeting:

- Nov 30th: deadline for WLCG pledges insertion → dialogue with FAs started
- CMS continued to work on model evolution and optimisations
 - ❖ explore ways to limit the increases in resource requests, while protecting the CMS Physics program
- request for a LHCC 'review' document
 - ❖ CMS delivered the doc on Feb 6th, one round of Q&A followed

In this talk (asked to be *short*):

- focus on overview of major activities and their scale over 2016 and changes w.r.t last LHCC meeting
- discussion on the 'review' doc later in the meeting
 - ❖ only one summary table in these slides

Update on Fall/Winter 2016-17 processing

Full 2016 data rereco **[A]**

- ~5B evts completed in ~5 weeks

Major MC re-DigiReco campaign **[B]**

- >10B evts completed in ~3 months

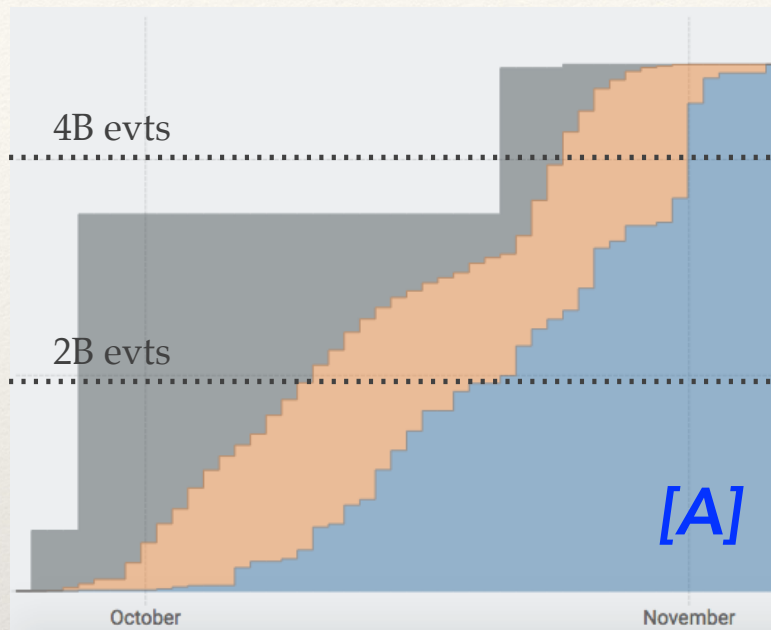
Re-MiniAOD for Moriond'17 **[C]**

- completed in ~1 week - it was crucial to still have AODs on disk

Now and next:

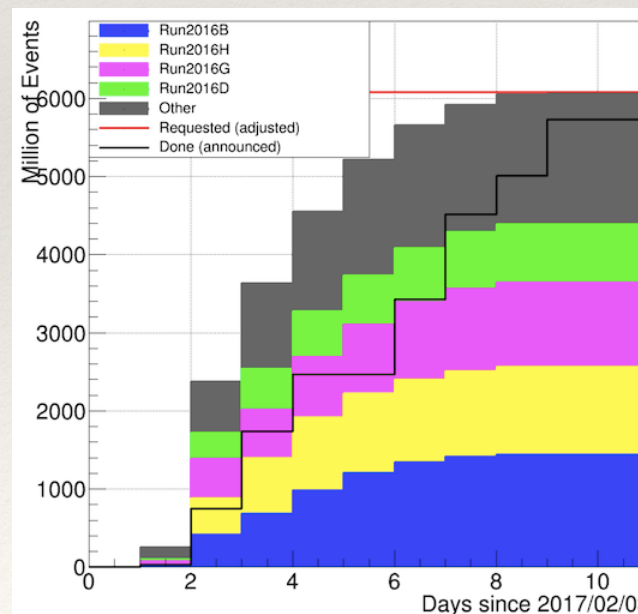
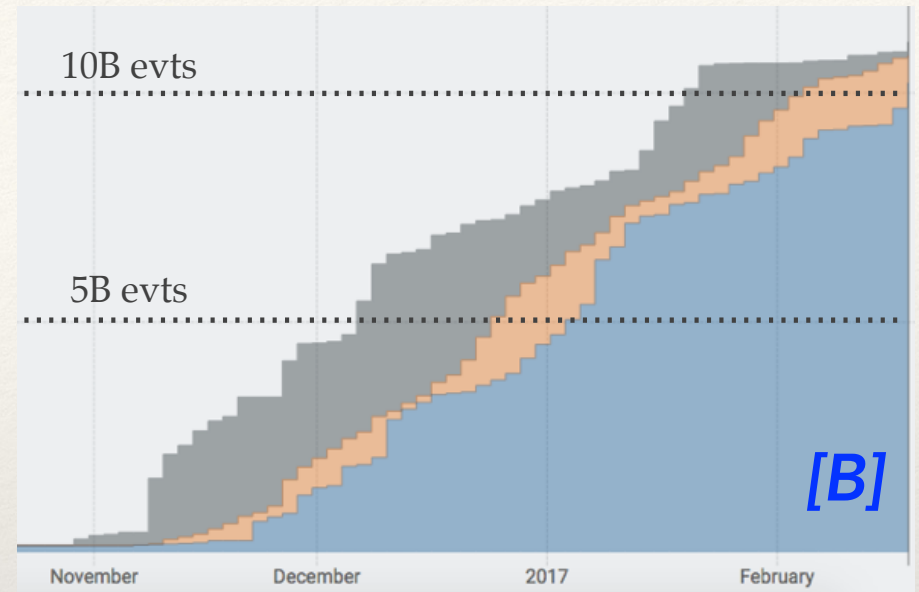
- Phase-I and Phase-II preparations flowing in
- legacy re-reco planned to start around mid/end of March

Update on Fall/Winter 2016-17 processing



- Requested
- Delivered
- Valid (DBS)

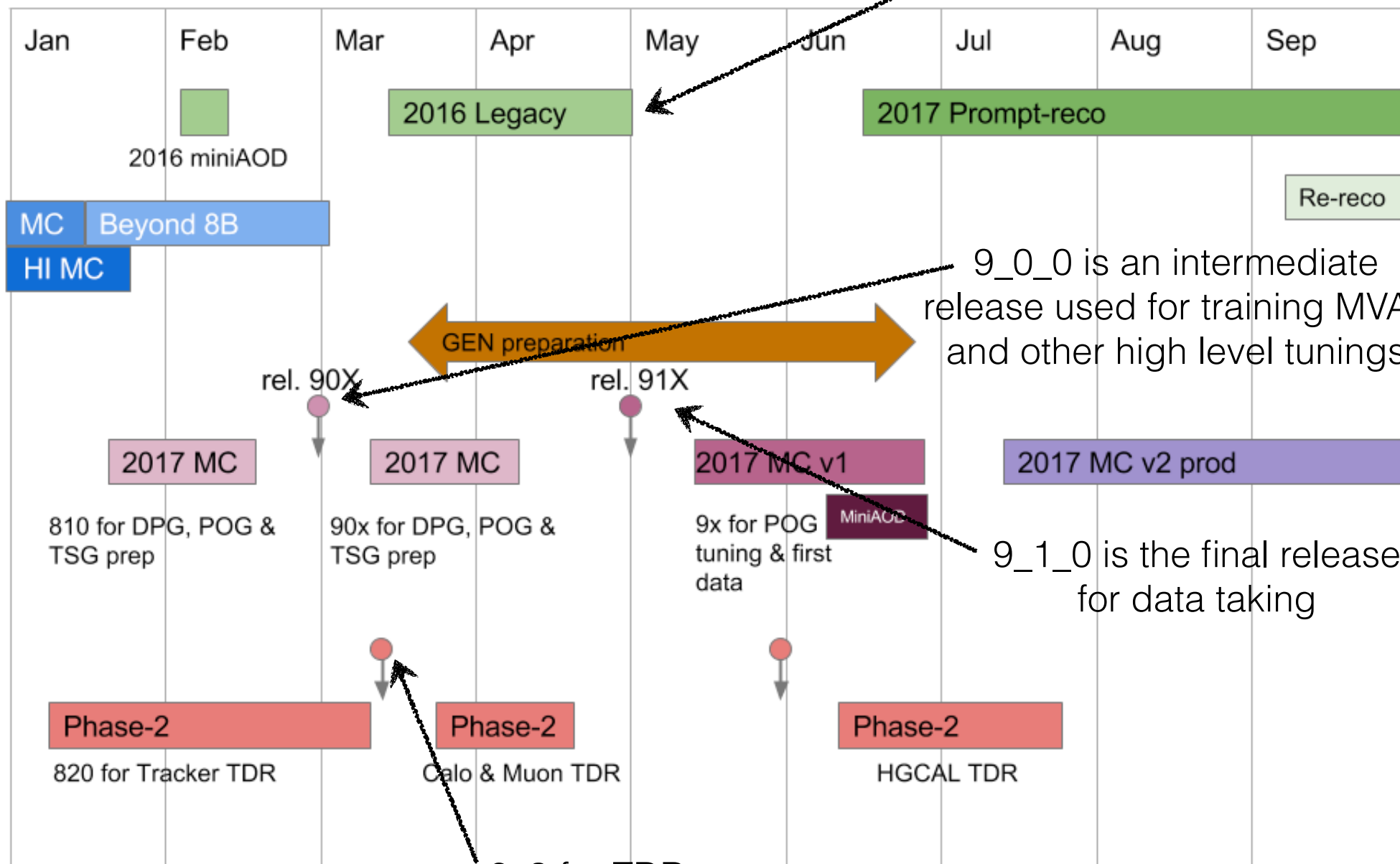
Remarkable slope of delivered events compared to requests injection rate





2017 production activities

Will use 8_0_X release with new conditions



9_0_0 is an intermediate release used for training MVAs and other high level tunings

9_1_0 is the final release for data taking

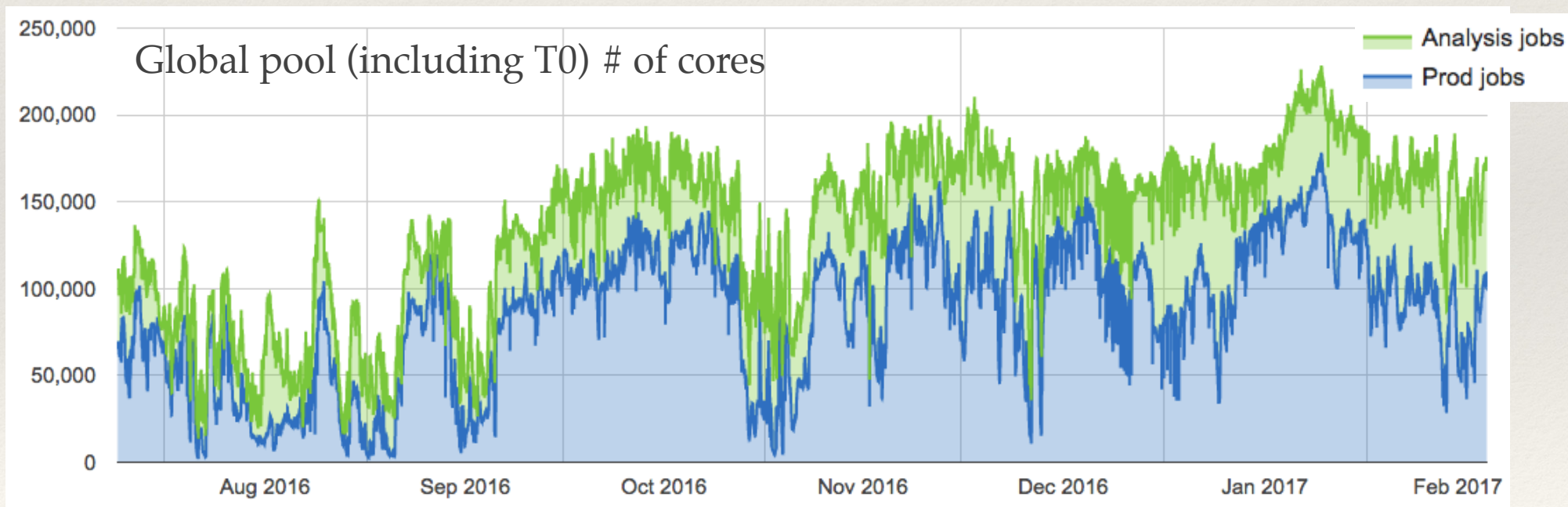
8_3 for TDR

Global Pool scaling up

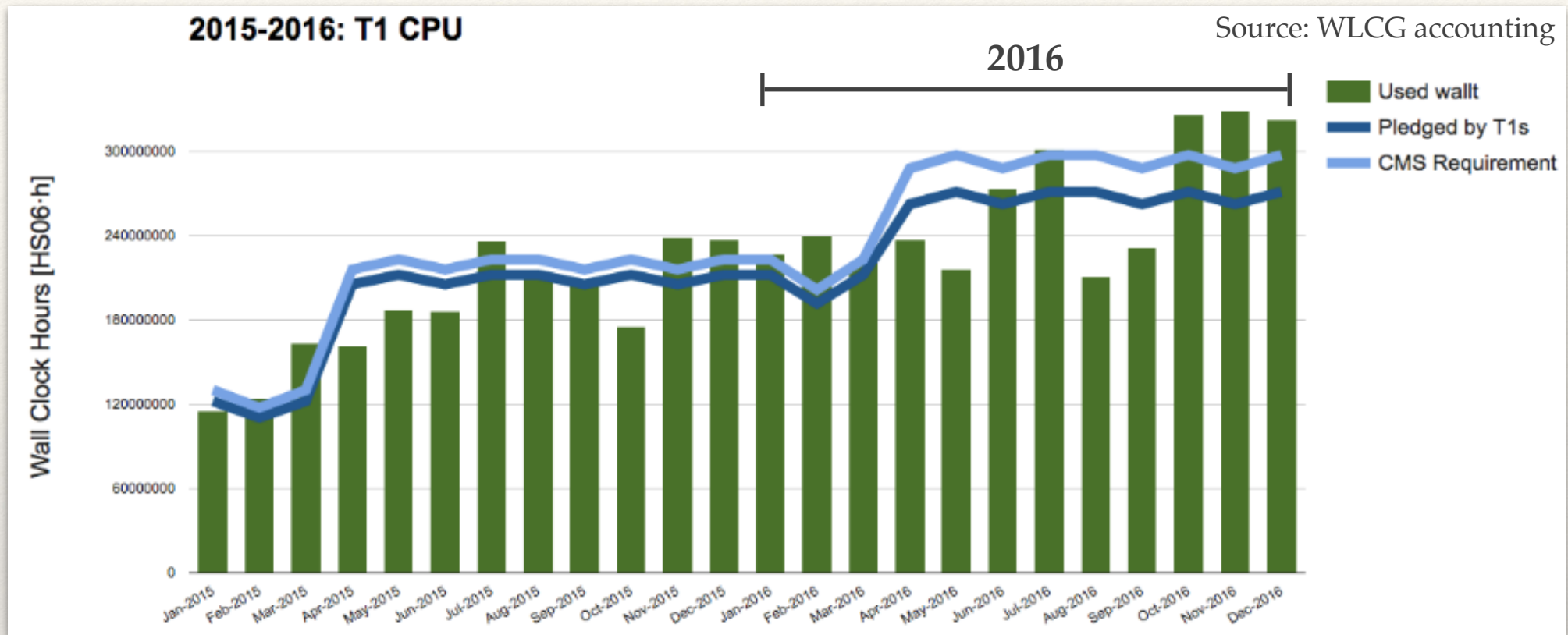
HTCondor Global Pool for resource provisioning via the glideInWMS infrastructure

- continues to be an EXCELLENT central control point for job priorities over resources
- average running cores in 2016 was ~130k, i.e. +50% as compared to 2015

Reached **>200k cores** scale on Grid/Cloud resources for the first time in early 2017

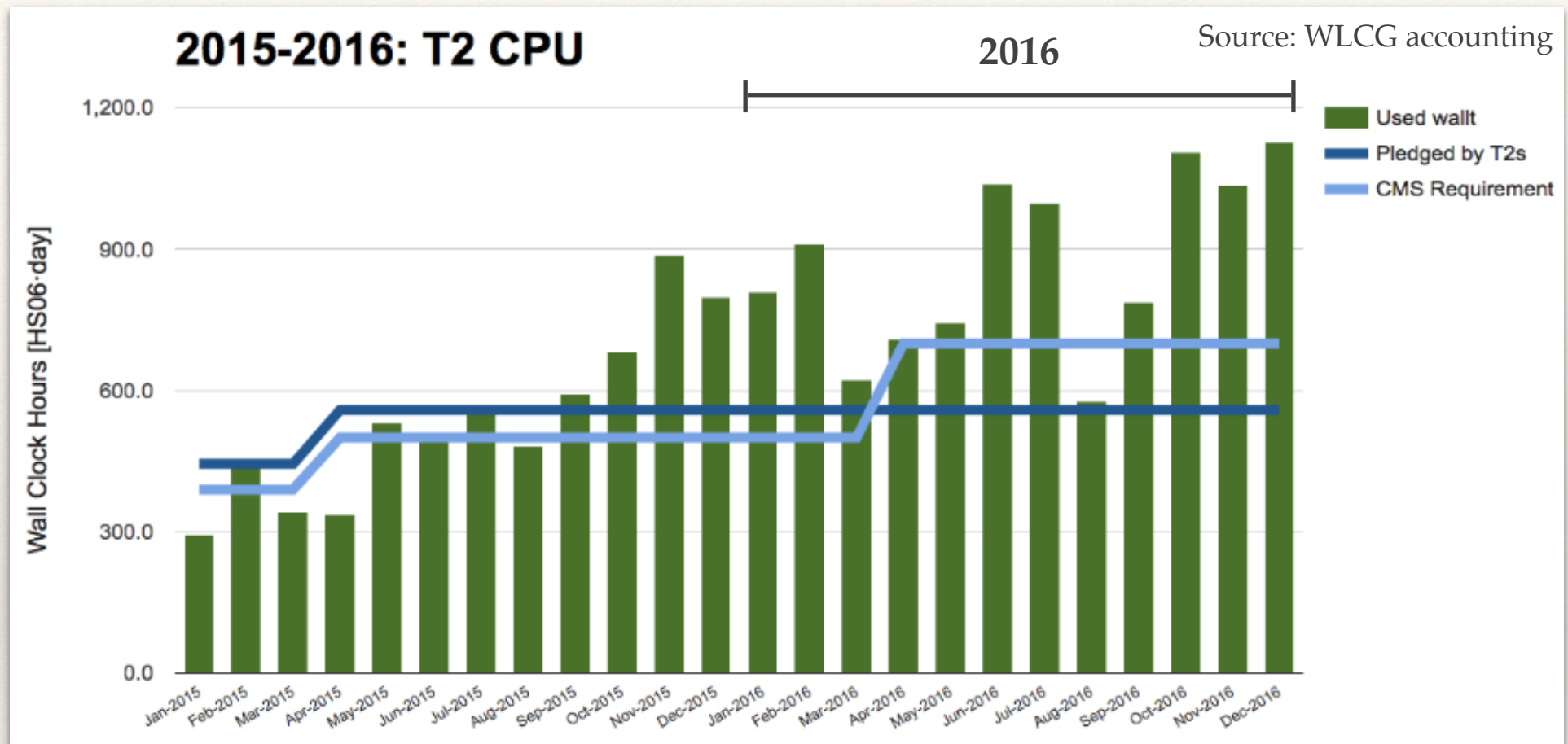


Resource utilisation: T1 CPU



In 2016, on average we used **103%** of the T1 CPU pledges (**96% as compared to the CMS requirement**)

Resource utilisation: T2 CPU



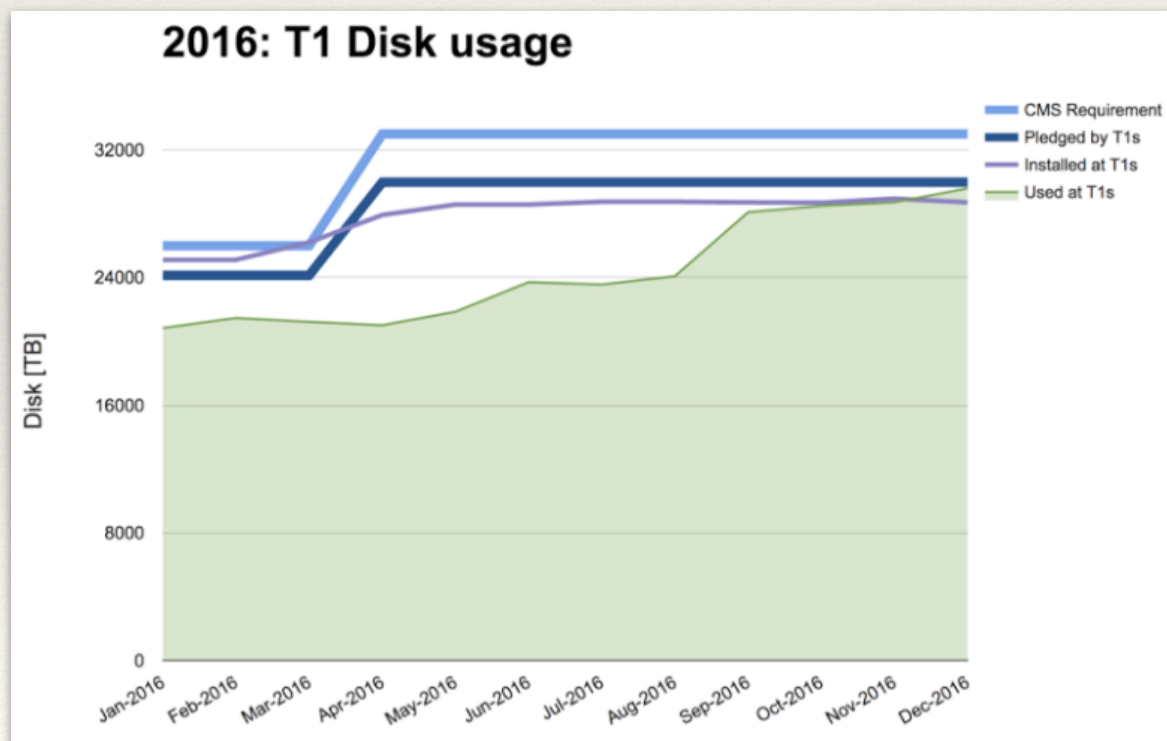
In 2016, on average we used **129%** of the T2 CPU pledges (**134%** as compared to the CMS requirement)

Disk usage and pledge deficit

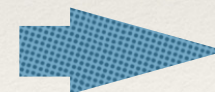
Pledge deficit in 2016 at the level of $\sim -7\%$

- high data taking rates in 2016 forced us to keep this constantly on the radar
- significant efforts to monitor and manage disk space at T1s. In the Fall, 4 Tier-1 sites were at risk to stop operation for lack of available disk space, **disk caches clean-up** triggered

By the end of 2016, CMS was using **99%** (!) of the disk pledges at T1s

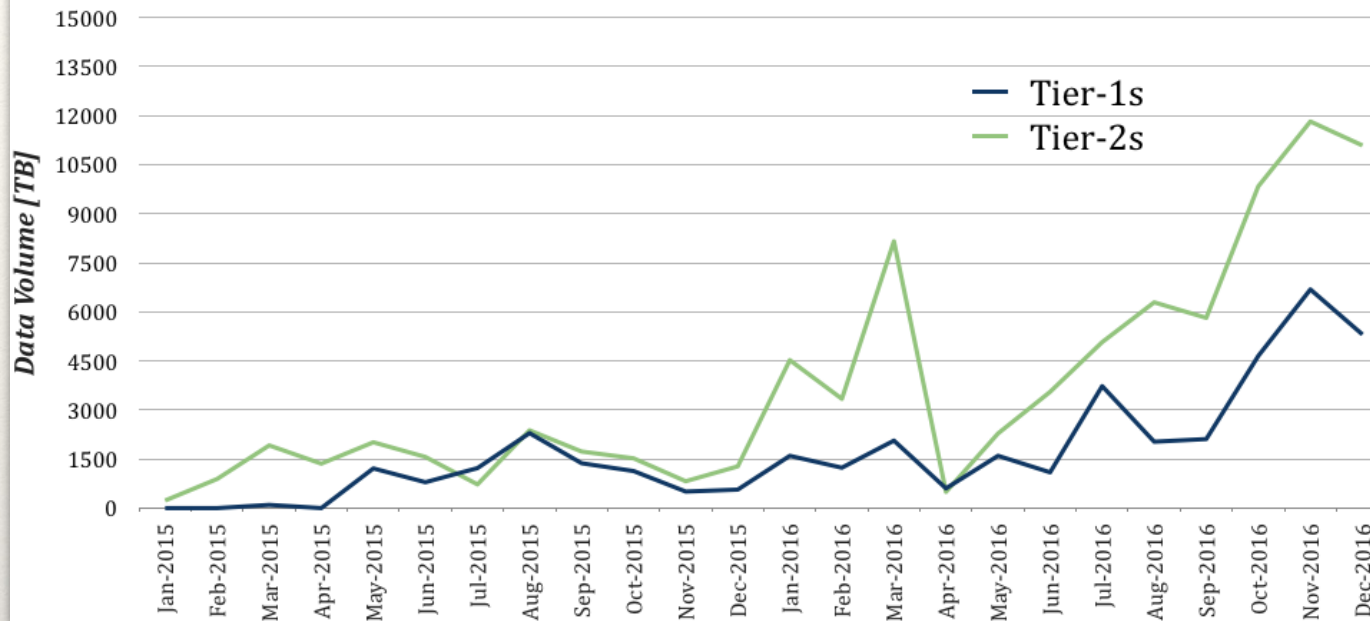


Frequent samples turnover
via dynamic DM



Dynamic disk space management

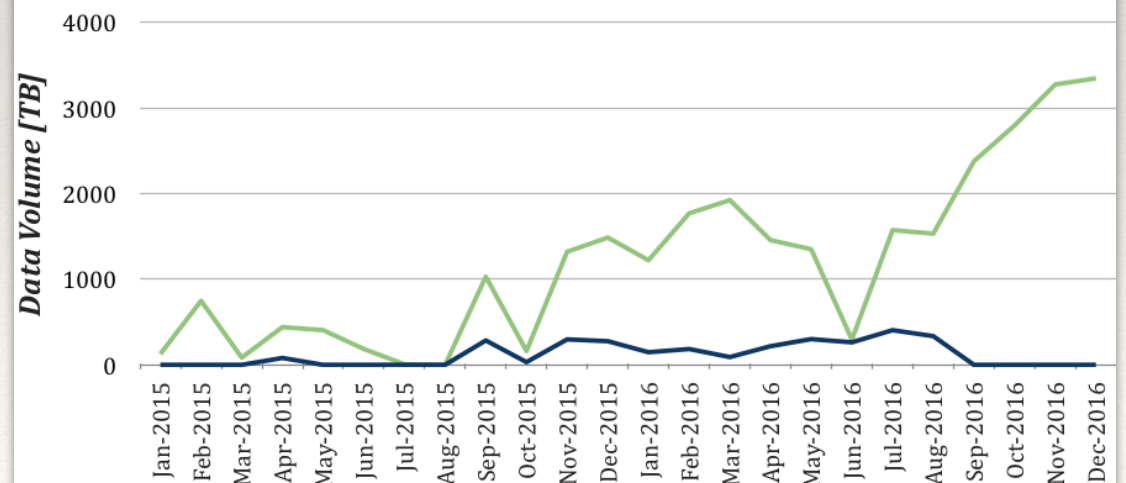
Data Deletions by DDM (2015-2016)



More aggressive deletions from disk starting in Summer '16

Re-population based on access needs (on T2s only, basically)

Data Subscriptions by DDM (2015-2016)



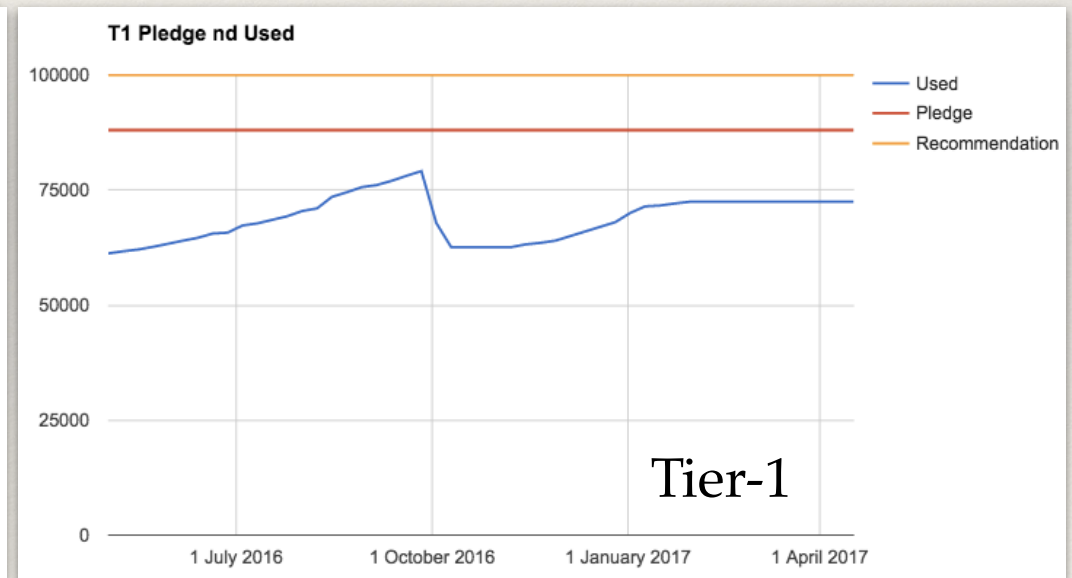
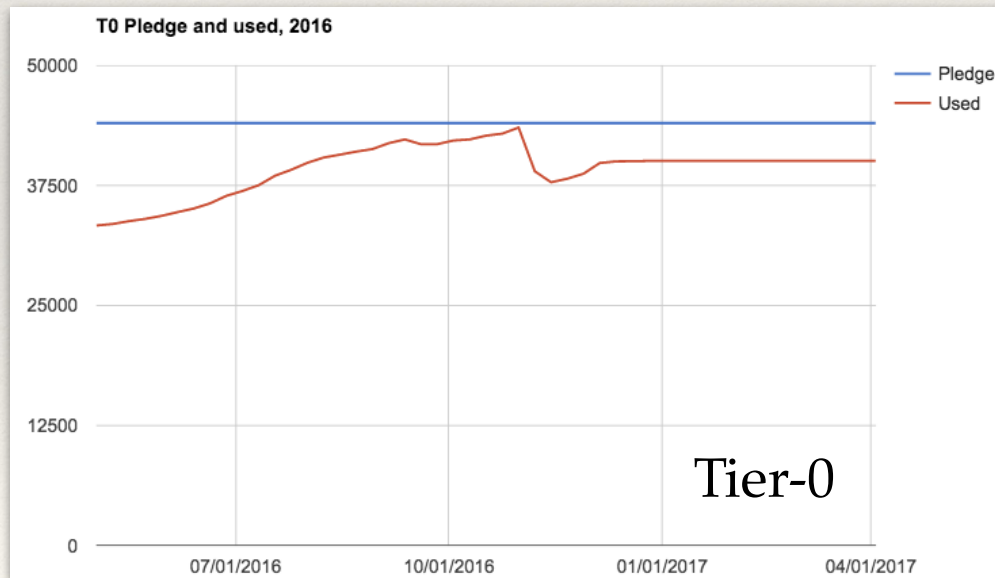
Tapes

Pledge deficit in 2016 at T1s at the level of $\sim -12\%$

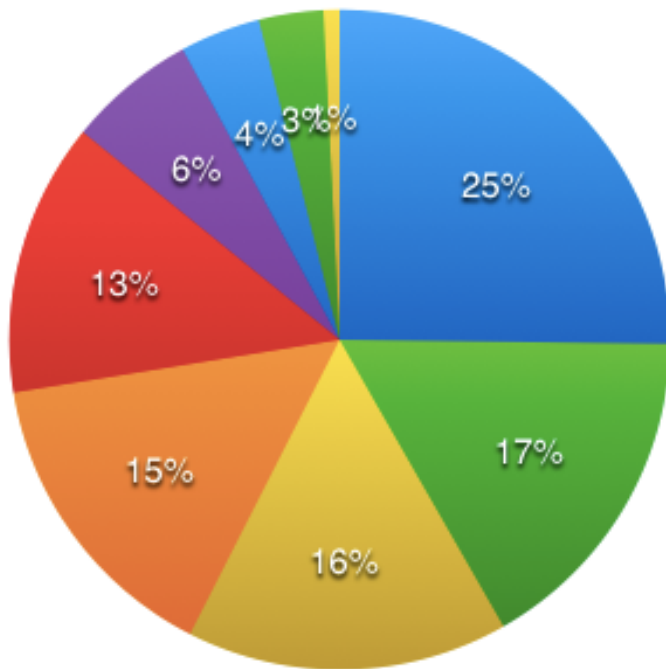
LHC performances in 2016 required quick reaction:

Massive **tape** deletion campaign

- prepared in advance of the trouble (Summer): very careful review of all data on tape
- executed (Fall) \rightarrow e.g. **~ 30 PB** deleted across all Tiers



Data transfers and access in 2016



Collected data on transfers are becoming statistically interesting for modelling studies

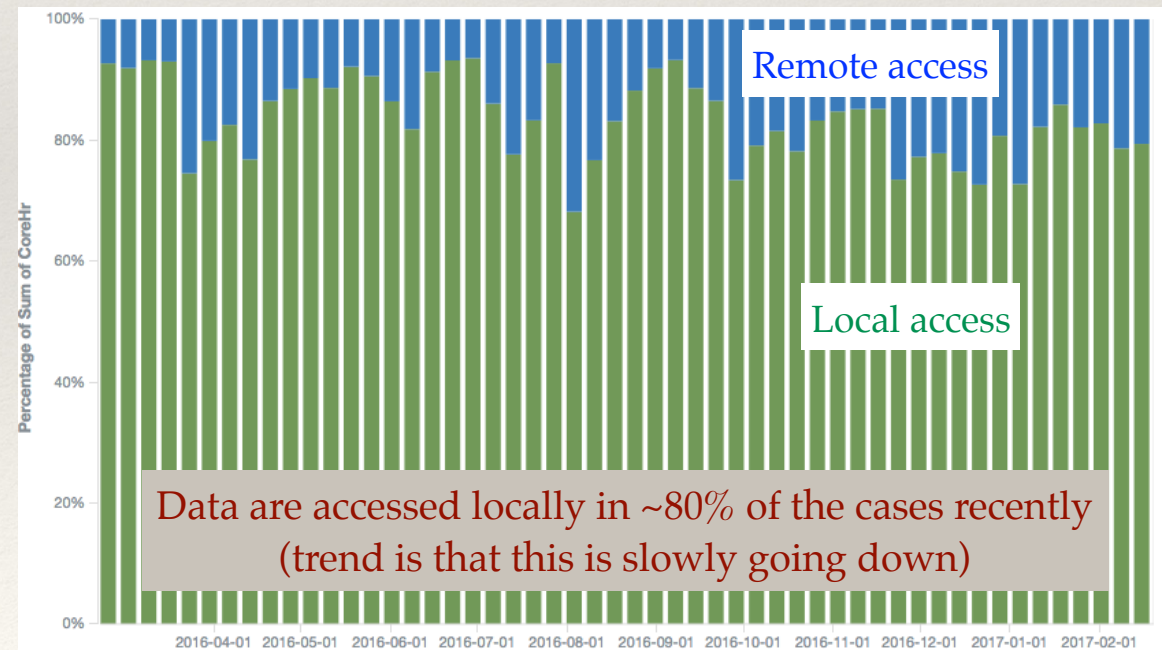
More than 3500 links between PhEDEx nodes on the overall Tiers topology

On average:

- globally, close to 3 PB/week. Weekly average >1GB/s on most busy routes

Peaks (weekly):

- up to almost 2.5 PB in one week in T0-T1 alone. Weekly peaks at ~4 GB/s in T0-T1, ~3GB/s in {T1,T2}-T2

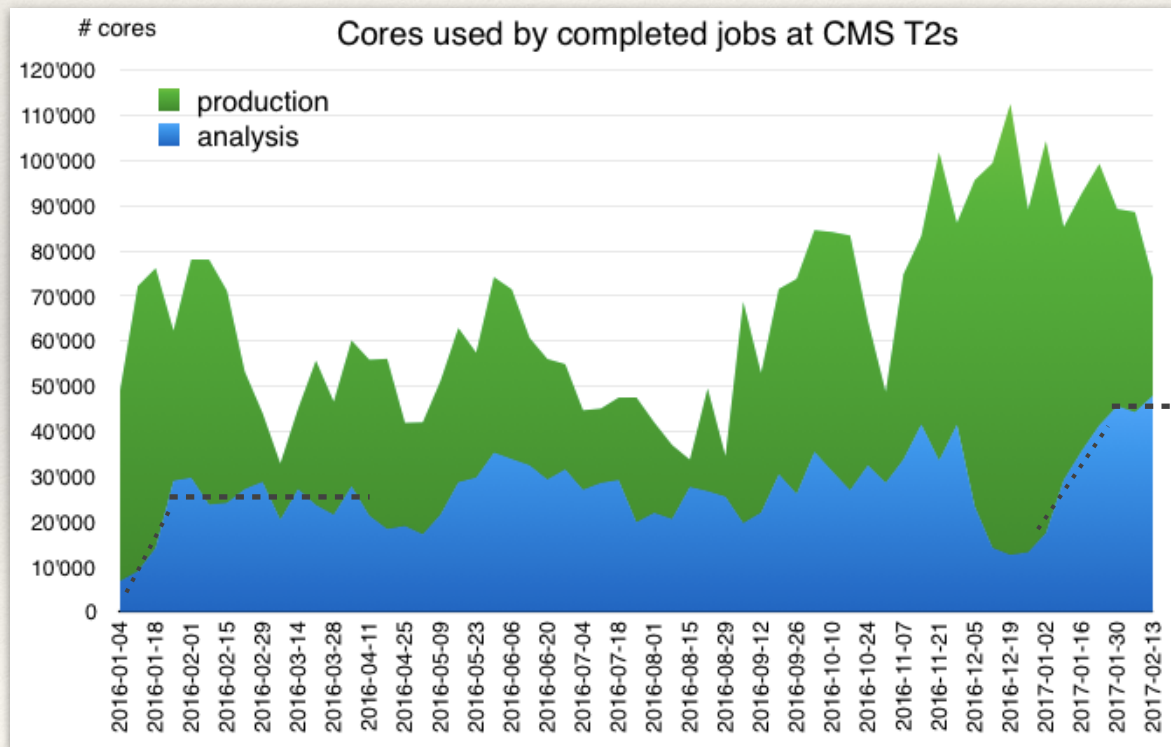


Distributed analysis

2016 vs 2015: increase in distributed analysis load

- in the number of distinct individuals per week submitting analysis jobs
- in the number of slots used at all Tier levels (shown below for T2s only)

CRAB-3 routinely at ~60k cores simultaneously used (at all Tiers)



Note: we changed the prod-analysis share in the global pool from 50-50 to 75-25 at the T2s for a couple of weeks, but during the holidays - now back to 50-50

Summary

CMS in 2016 continued to exploit the pledged computing resources at very high utilisation levels

- Global pool scaling up - according to the development/ops plan

In 2016 we managed to deliver thanks to planning and quick counter-measures

- disk caches clean-up, massive tape deletions, Ops load, ..

Careful planning is a must for 2017-18

- more in the CRSG doc due soon

For the discussion: [see next](#)

NOTE: a talk at the CMS session just later today. Full details with request tables in the CRSG document.

We now have a CMS resource request model that should basically fit within financial boundaries for 2017-18 (first glance at 2019)

- not at zero cost: the execution of such model puts additional daily load on Operations manpower - additional work and extra vigilance in operations takes time away from developing new methods, since that has been key to dealing with scarce resources. The rigorous management to deal with the present does have implications for the ability to respond to future needs.
- we have no more safety margins left to be prepared for the unexpected
- we think the resource request model is fully inline with the requests coming from **LHCC** and **CRSG**
 - ❖ **LHCC**: input document delivered on Feb 6th, Q&A, discussion at the WLCG-LHCC session
 - ❖ **CRSG**: full doc due end of February, discussion in March, in preparation for April RRB

for the discussion

The LHCC doc in one table

Full doc sent by mail to LHCC referees on Feb 6th

Areas definitions (by LHCC)	Priority	Status	Few remarks
Optimization of workflows	very high	DONE	Premixing mode for high PU simulations at large scale. Deployment Nov 2016, used in production for MC DigiReco for Moriond'17. Local IO reduction + reduce 2x amount of CPU needed for DigiReco → <u>minimise the CMS requests for CPU in 2018</u> . Performed studies to exclude statistical biases. Pressure on remote access ops/monitoring (but a robust data federation is beneficial anyway).
Technology improvements	high	R&D	Potential to <u>reduce complexity by large factors</u> , but in which technology area and the timeline are largely unpredictable. GPU integration, HPC centres exploitation, opportunistic cloud extension of WLCG centres, overall orchestration of diverse resources, new FA mechanisms to offer these as pledges are all aspects to consider. Integration efforts (and manpower) not negligible.
Data / CPU / Tape management	very high	partially DONE	CMS computing model evolved towards higher flexibility in LS1. Main workflows can be submitted (almost) at any Tier level. Commissioned processing chains for better streamlining of global processing efforts (e.g. GEN to miniAOD can be run as a single step) → <u>save CPU and especially tape</u> . A limitation comes from being impractical for all GS, but very useful for a fraction of them.
Triggers thresholds tuning	medium	not DONE	<u>Negative impact of rate reductions on physics output is potentially large</u> . May be justified resource-wise only if reductions are sizeable. Rough estimates indicate that 1kHz → 800 Hz yields relatively modest savings, and put some physics programs at risk. Requires careful scrutiny and guidance by ECoM. In general, CMS would not suggest to pursue this path.
Amount of simulation	medium	not DONE	MC/data ratio tuned at 1.3 in the CMS resources model. Recently needed to do more than expected (both 2015 and 2016). Rough calculations for a 130%→100% reduction in 2018 yield savings of -8% (CPU), -3% (disk), -3% (tape). Extreme caution needed to avoid impact on physics by such reduction. New assessment for optimal tuning is in the ECoM mandate.
Parking / Delayed processing	medium	1. DONE 2. not DONE	Distinction between delayed processing and parking+scouting . CMS can do (and does) the latter, but would discourage to pursue the former. Gain of parking 200 Hz to be rereco'ed later would be quantitatively similar as the estimates for Triggers thresholds tuning above, with no gain on tape space as RAW will still be written. Devastating impact in the former case for B physics, for instance
Copies / Formats versatility / Analysis Frameworks	high	almost DONE	MiniAOD format introduction, ~8x gain in size, used by ~80% of the analyses today. Larger adoption is planned, but it will take time. → <u>reduce the disk needs</u> . Caveat: during the transition, miniAOD plus a fraction of AODs need to stay on disk to support all analyses. Dynamic use of storage space, load on Ops, mitigated by more automation. Remaining need for AOD must be small.

Backup

Tier-0 and HLT

