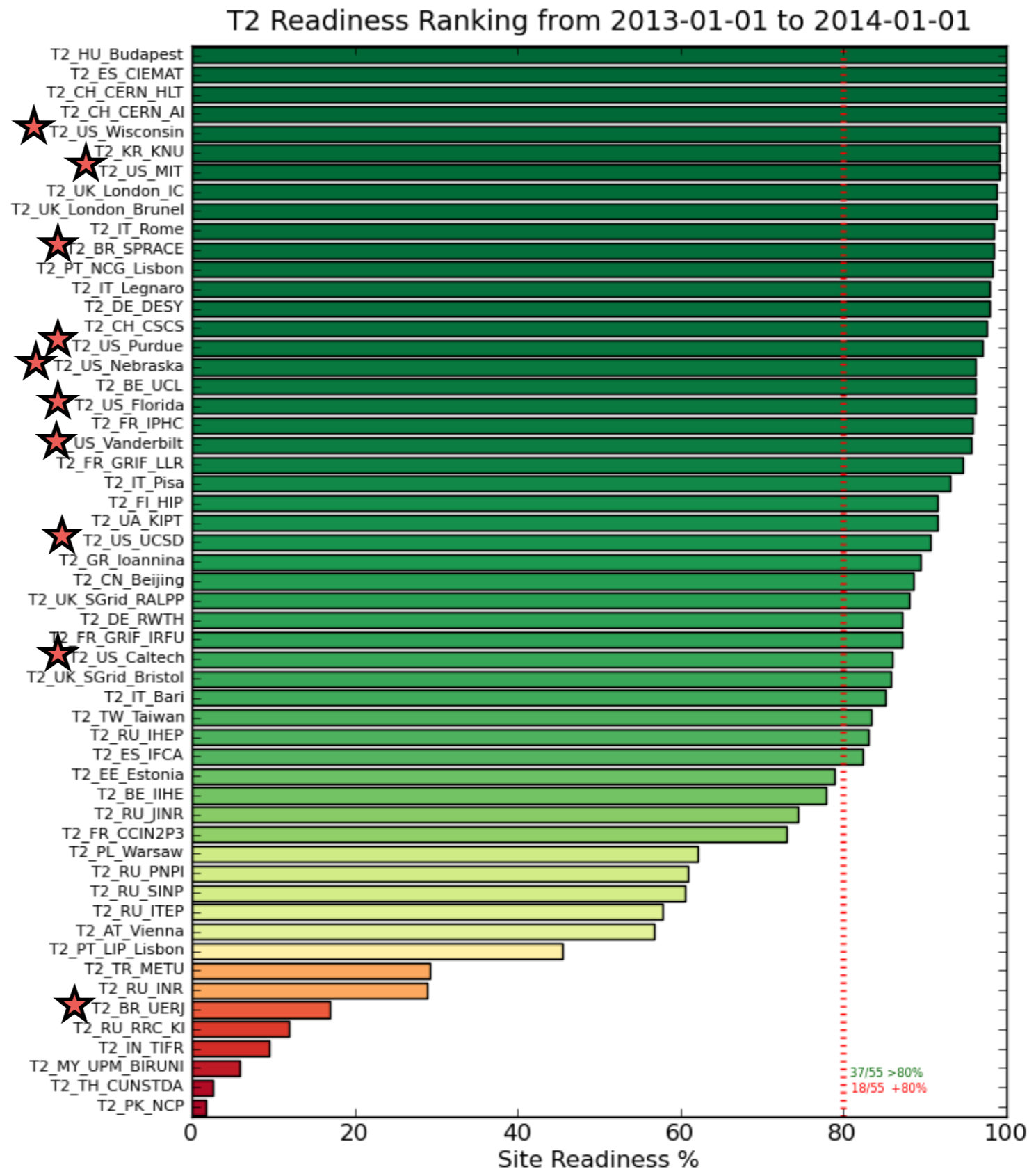


US CMS Tier-2 Status and Future

Ken Bloom
April 7, 2014



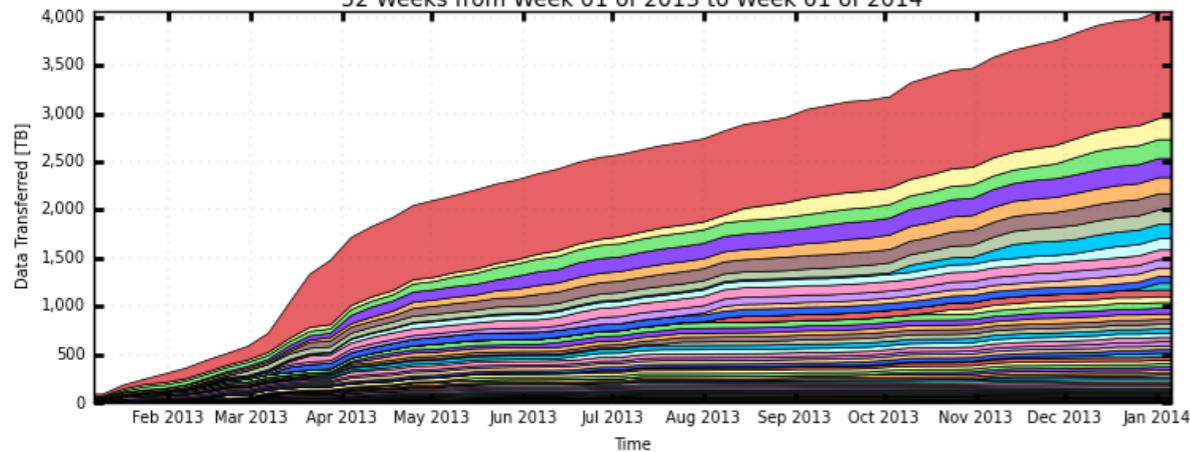
- ▶ The site readiness metric has its quirks, but is still our best measure of the usability of each site on any given day, on the basis of ability to execute basic functions, including data transfers
- ▶ Makes allowances for short-term problems, scheduled downtimes, weekends
- ▶ CMS goal for this is >80%, all but one OSG site above



Outbound: 4.1 PB

CMS PhEDEx - Cumulative Transfer Volume

52 Weeks from Week 01 of 2013 to Week 01 of 2014

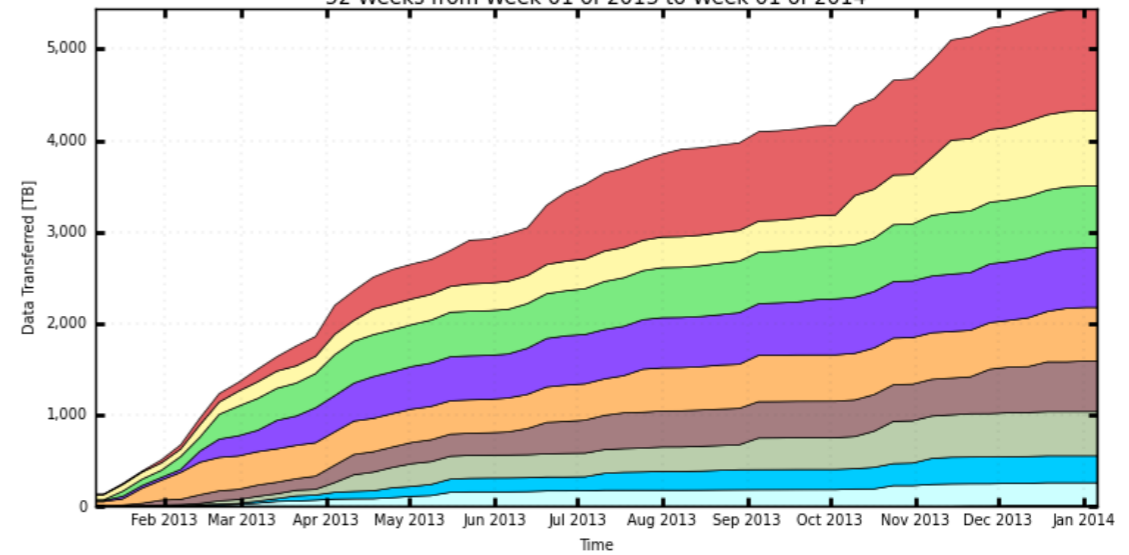


- | | | | | |
|-------------------|----------------------|-------------------|-----------------|------------------|
| T1_US_FNAL_Buffer | T1_FR_CCIN2P3_Buffer | T1_IT_CNAF_Buffer | T2_DE_DESY | T1_DE_KIT_Buffer |
| T2_CH_CERN | T1_ES_PIC_Buffer | T2_US_Nebraska | T3_US_Colorado | T2_BE_IHE |
| T2_US_Wisconsin | T1_IT_CNAF_Disk | T0_CH_CERN_Export | T1_US_FNAL_Disk | T1_UK_RAL_Disk |
| T1_DE_KIT_Disk | T2_IT_Bari | T2_US_Florida | T2_IT_Pisa | T2_EE_Estonia |
| T2_US_Vanderbilt | T2_UK_London_IC | T2_US_MIT | T2_US_UCSD | T1_UK_RAL_Buffer |
| T2_PT_NCG_Lisbon | T2_UK_SGrid_RALPP | T3_US_TAMU | T2_AT_Vienna | T2_RU_JINR |
| T2_BR_SPRACE | T3_US_NU | T2_DE_RWTH | T2_ES_CIEMAT | T3_TW_NCU |
| T2_US_Caltech | T2_IT_Rome | T2_RU_SINP | T2_FR_CCIN2P3 | T2_KR_KNU |
| T2_IT_Legnaro | T2_IN_TIFR | T2_FR_IPHC | T3_US_Minnesota | T2_CH_CSCS |
| T2_FR_GRIF_LLJ | T3_US_FNALLPC | T1_TW_ASGC_Buffer | T3_US_Brown | ... plus 14 more |
- Total: 4,063 TB, Average Rate: 0.00 TB/s

Inbound: 5.4 PB

CMS PhEDEx - Cumulative Transfer Volume

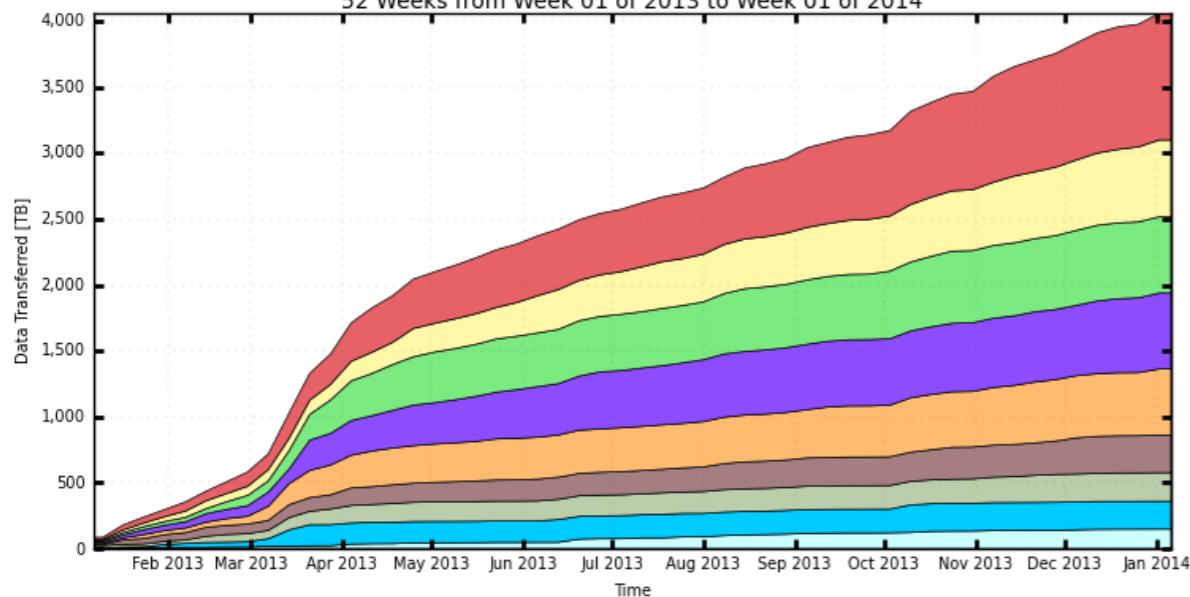
52 Weeks from Week 01 of 2013 to Week 01 of 2014



- | | | | | |
|-----------------|----------------|---------------|--------------|------------------|
| T2_US_Wisconsin | T2_US_Nebraska | T2_US_UCSD | T2_US_MIT | T2_US_Vanderbilt |
| T2_US_Florida | T2_US_Purdue | T2_US_Caltech | T2_BR_SPRACE | T2_BR_UERJ |
- Total: 5,441 TB, Average Rate: 0.00 TB/s

CMS PhEDEx - Cumulative Transfer Volume

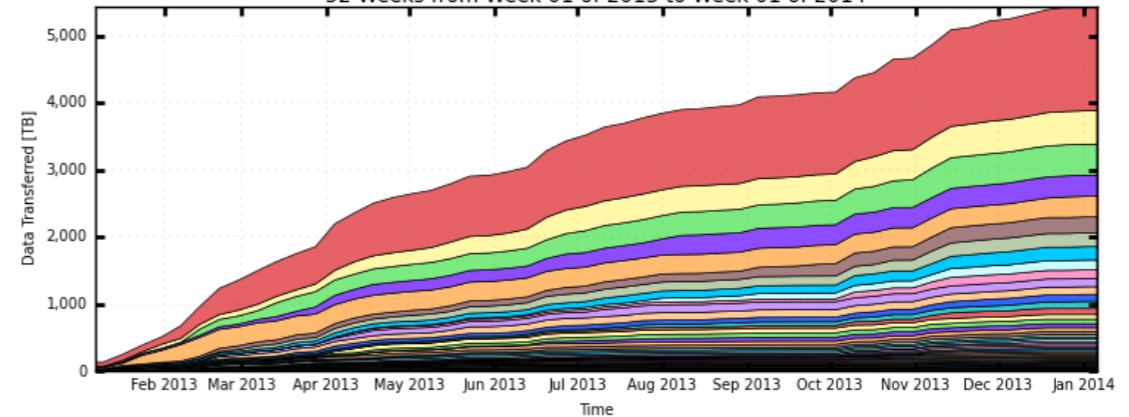
52 Weeks from Week 01 of 2013 to Week 01 of 2014



- | | | | | |
|---------------|-----------------|------------------|----------------|------------|
| T2_US_Purdue | T2_US_Wisconsin | T2_US_MIT | T2_US_Nebraska | T2_US_UCSD |
| T2_US_Florida | T2_US_Caltech | T2_US_Vanderbilt | T2_BR_SPRACE | T2_BR_UERJ |
- Total: 4,063 TB, Average Rate: 0.00 TB/s

CMS PhEDEx - Cumulative Transfer Volume

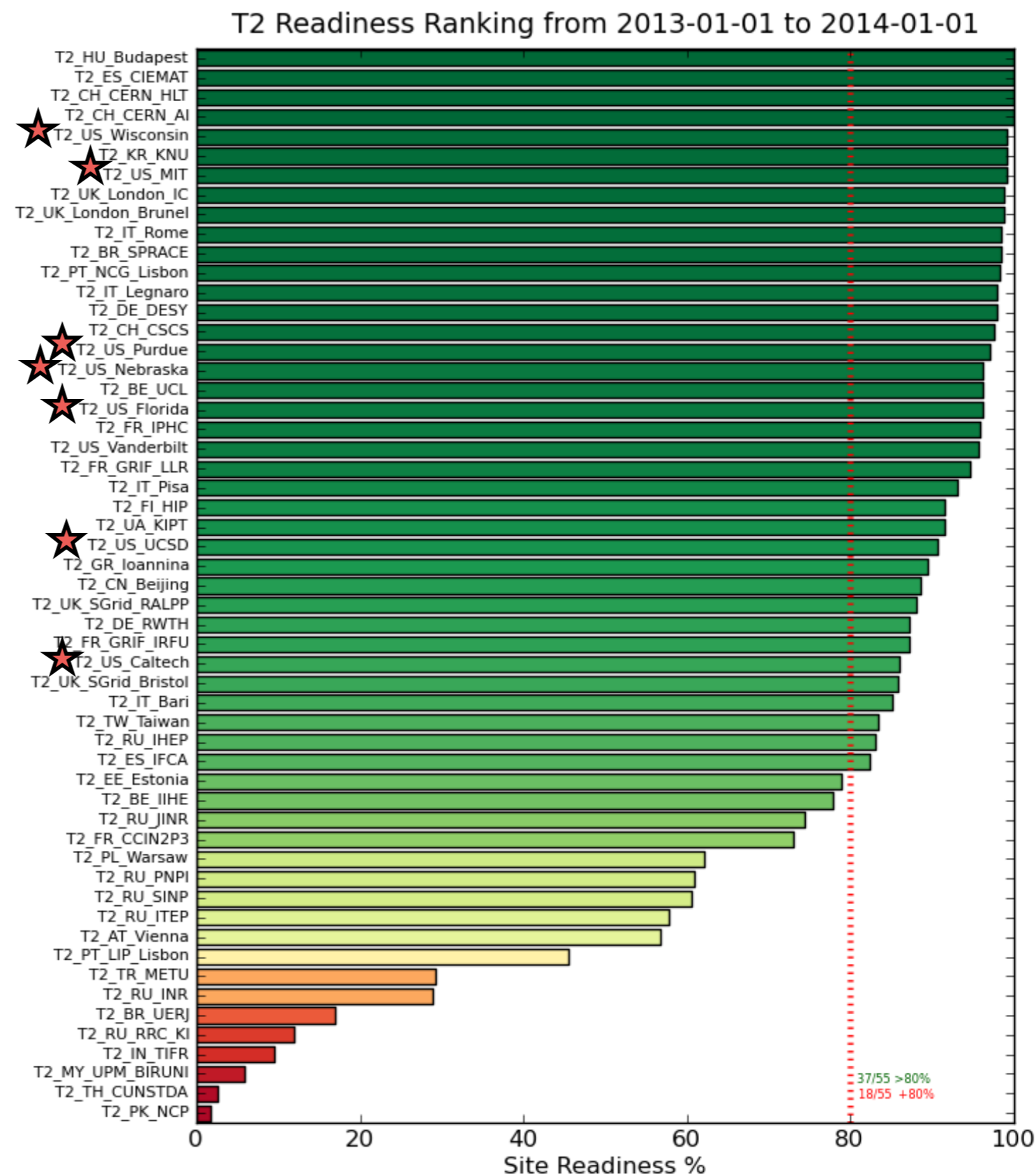
52 Weeks from Week 01 of 2013 to Week 01 of 2014

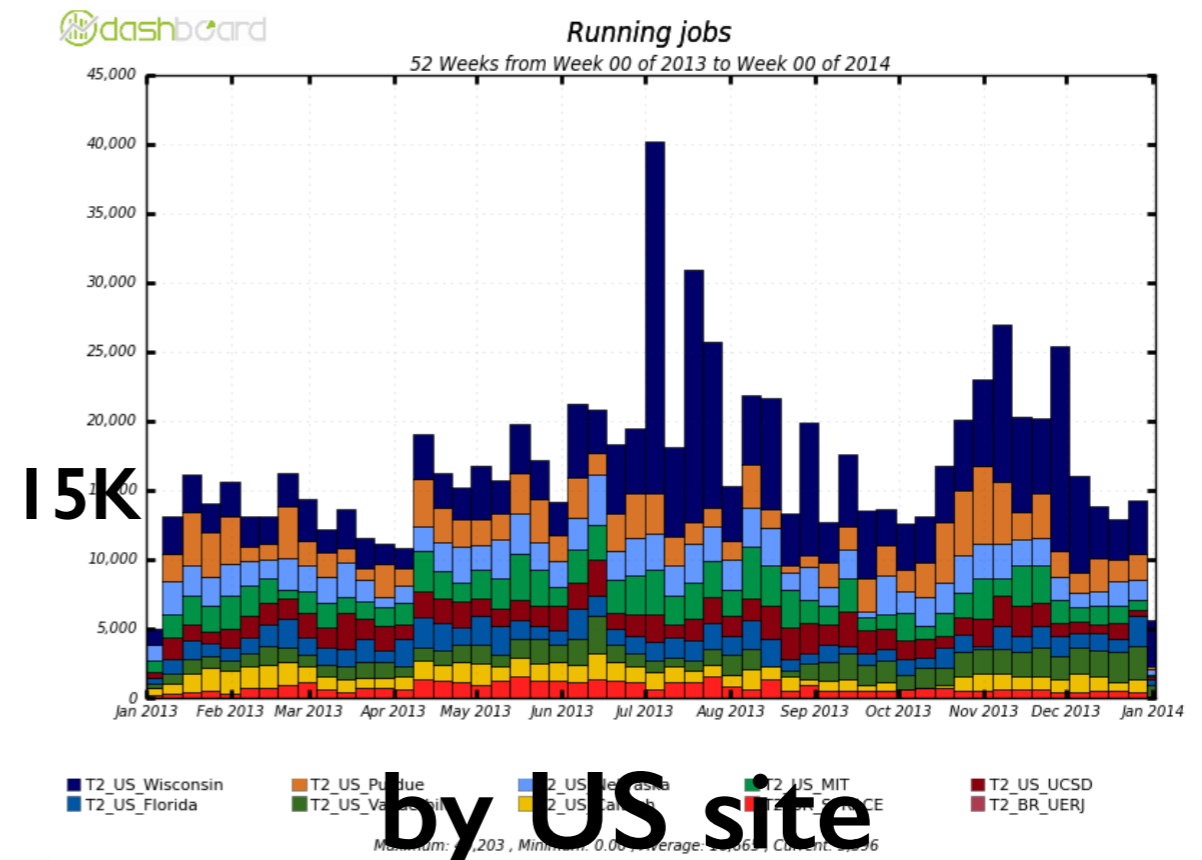
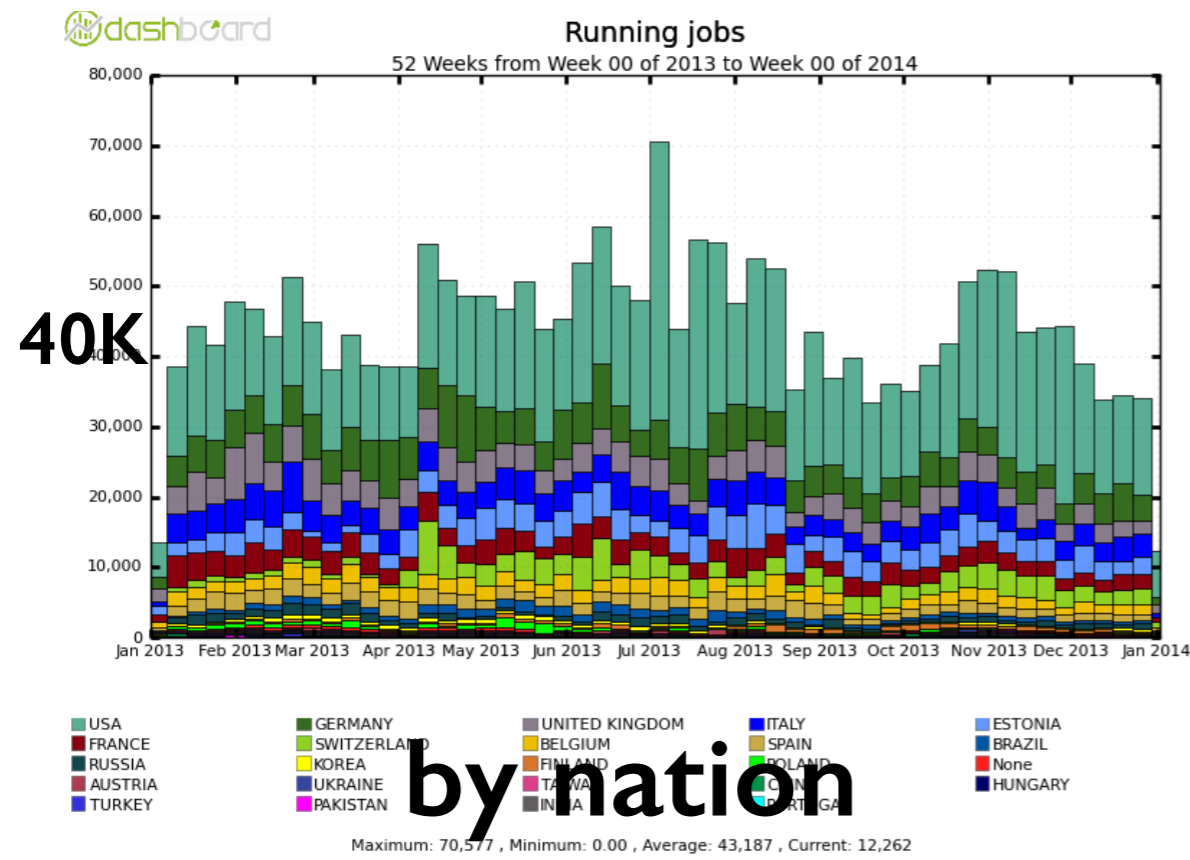
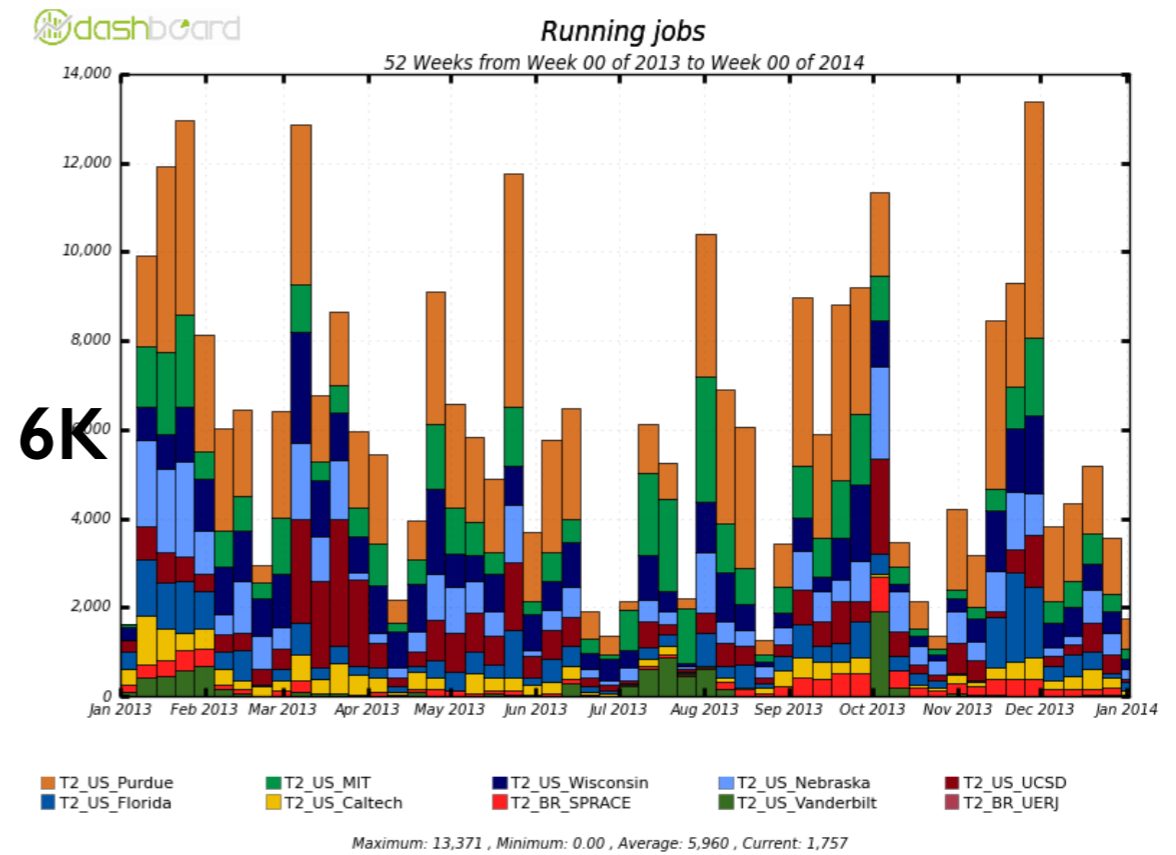
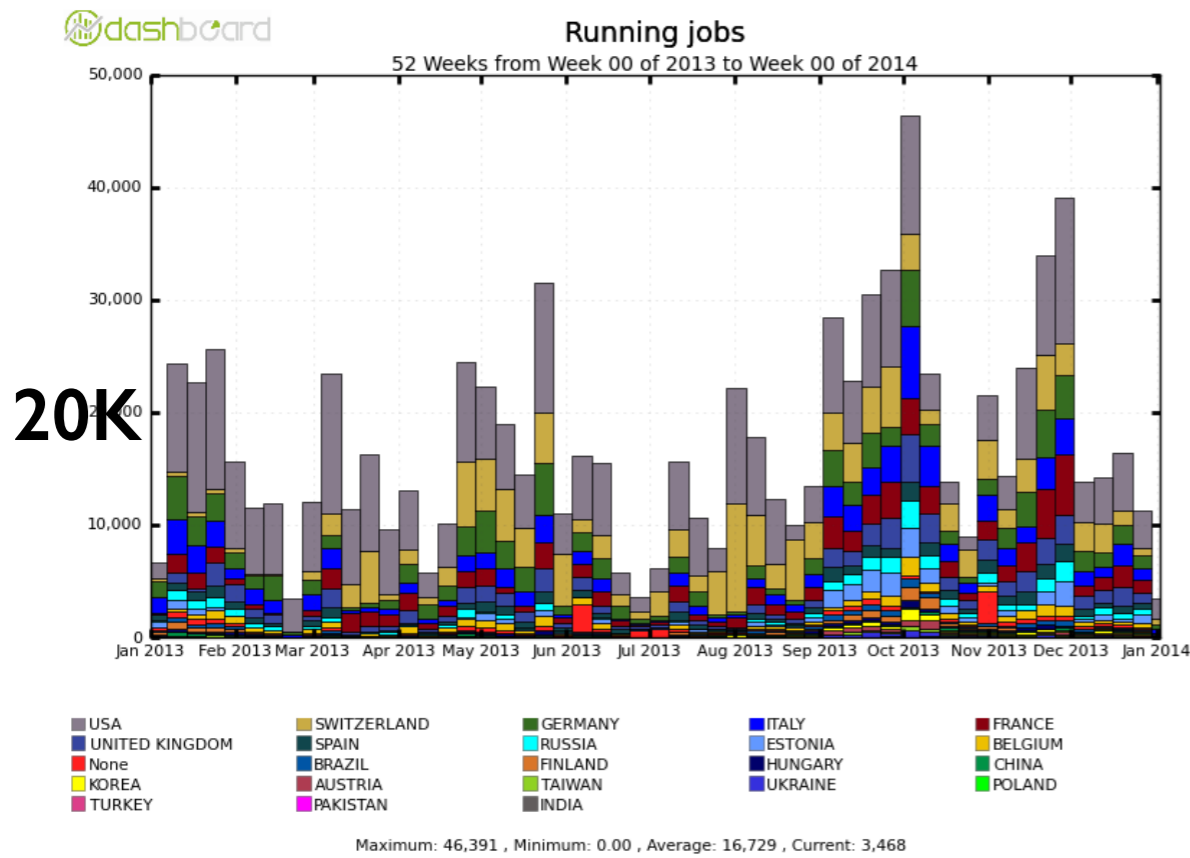


- | | | | | |
|----------------------|---------------------|------------------|------------------|-------------------|
| T1_US_FNAL_Buffer | T1_IT_CNAF_Buffer | T1_DE_KIT_Buffer | T1_UK_RAL_Buffer | T0_CH_CERN_Export |
| T1_FR_CCIN2P3_Buffer | T1_ES_PIC_Buffer | T2_DE_DESY | T2_CH_CERN | T2_US_Purdue |
| T2_UK_London_IC | T2_DE_RWTH | T2_ES_CIEMAT | T3_US_FNALLPC | T2_US_Wisconsin |
| T1_TW_ASGC_Buffer | T2_US_Nebraska | T2_BR_SPRACE | T3_US_Minnesota | T2_US_Florida |
| T2_US_MIT | T2_UK_SGrid_RALPP | T2_US_UCSD | T2_EE_Estonia | T2_US_Vanderbilt |
| T2_IT_Pisa | T2_US_UCSD | T1_IT_CNAF_Disk | T1_UK_RAL_Disk | T2_IT_Rome |
| T2_PT_NCG_Lisbon | T2_ES_IFCA | T2_FR_GRIF_IRFU | T2_CH_CSCS | T2_IT_Legnaro |
| T2_FR_CCIN2P3 | T2_IT_Bari | T2_FR_GRIF_LLJ | T2_FR_GRIF_LLJ | T2_US_Colorado |
| T2_BE_UCL | T2_UK_London_Brunel | T1_DE_KIT_Disk | T2_TW_Taiwan | T2_US_Caltech |
| T2_FR_IPHC | T2_CN_Beijing | T2_CN_Beijing | T2_RU_INJ | T2_BR_UERJ |
- Total: 5,441 TB, Average Rate: 0.00 TB/s

by destination
by source

- ▶ The site readiness metric measures the usability of each site on any given day, on the basis of ability to execute basic functions, including data transfers
- ▶ Makes allowances for short-term problems, scheduled downtimes, weekends
- ▶ CMS goal for this is >80%, all US sites are above





by nation

by US site

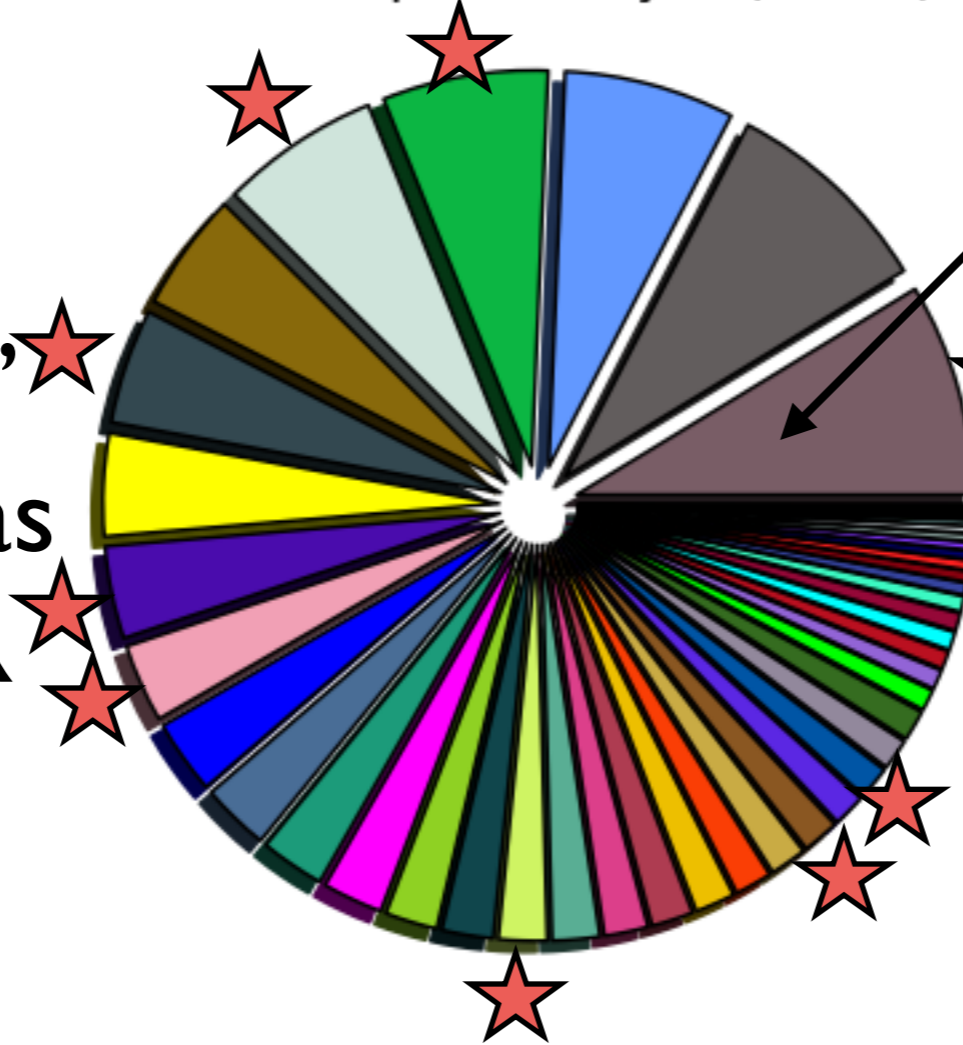
Analysis at T2s Production at T2s



weeks: CPU consumption Good Jobs (Sum: 1,327,224)

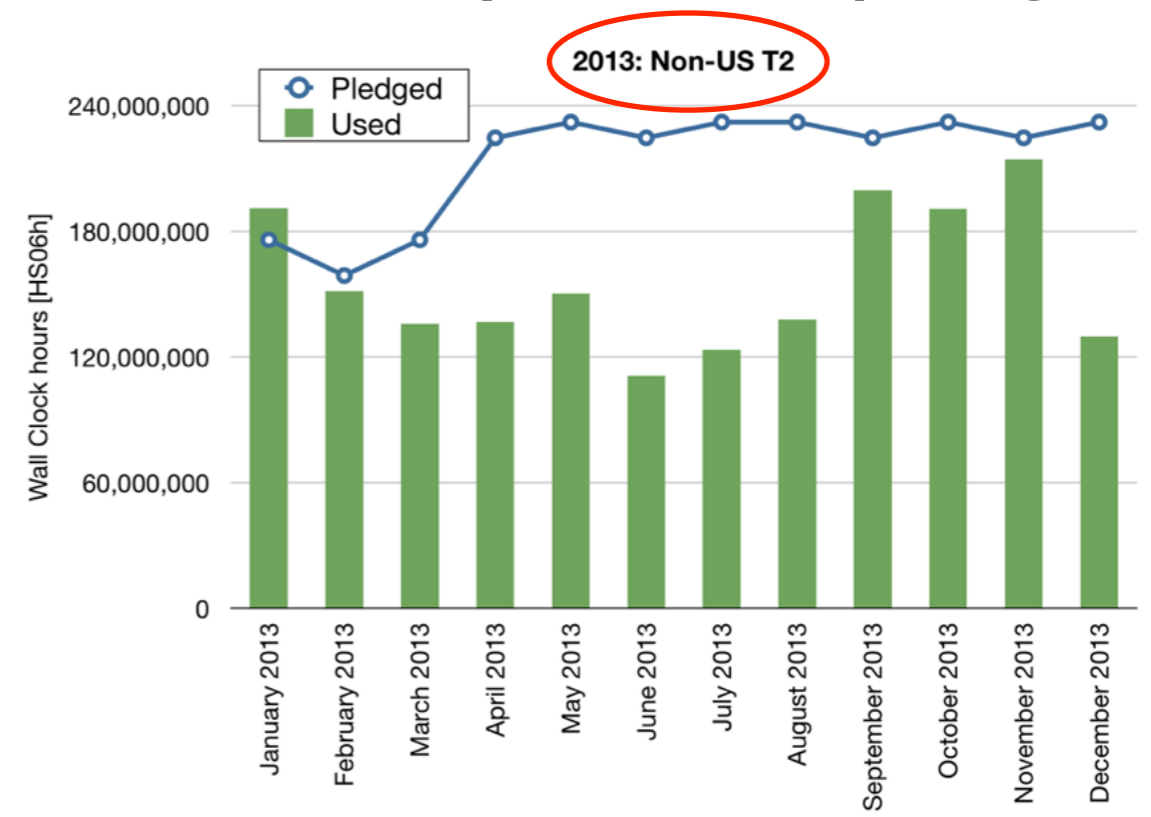
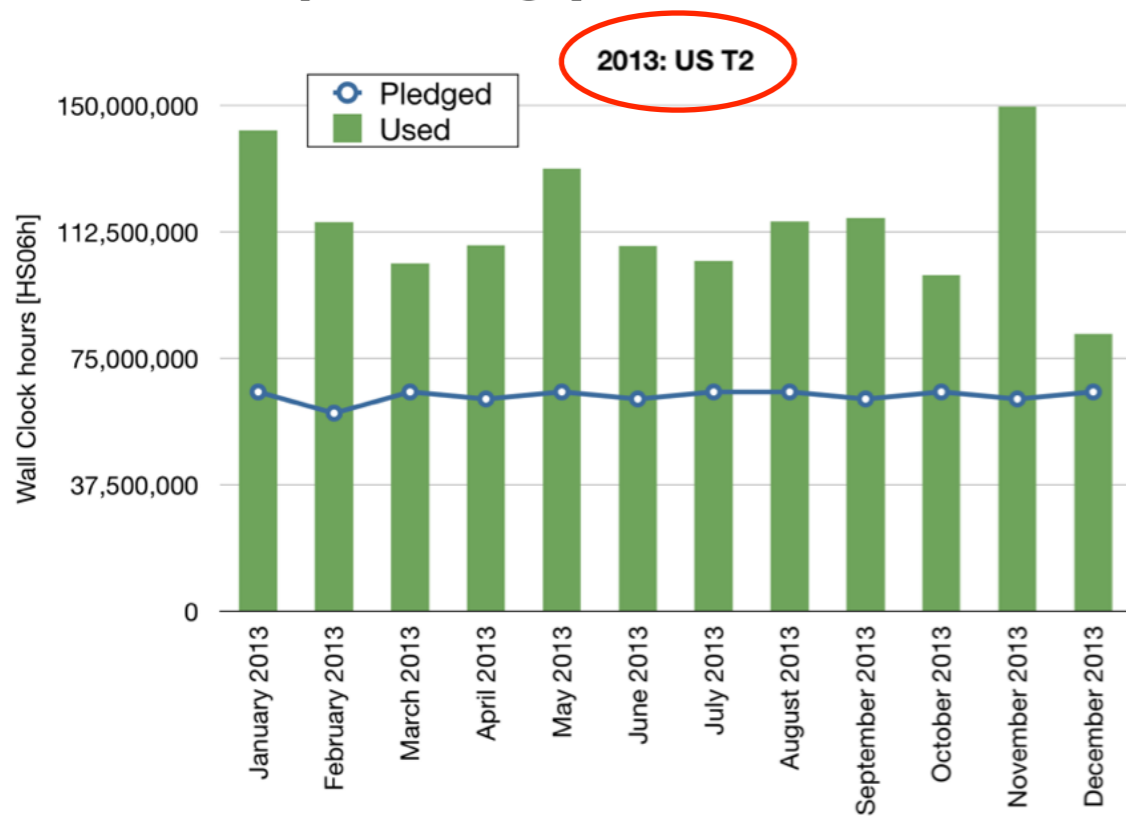
38.4% of “good”
T2 CPU time was
from the US/BR

UCSD is biggest
because of special
opportunistic use
of SDSC resources
this year!

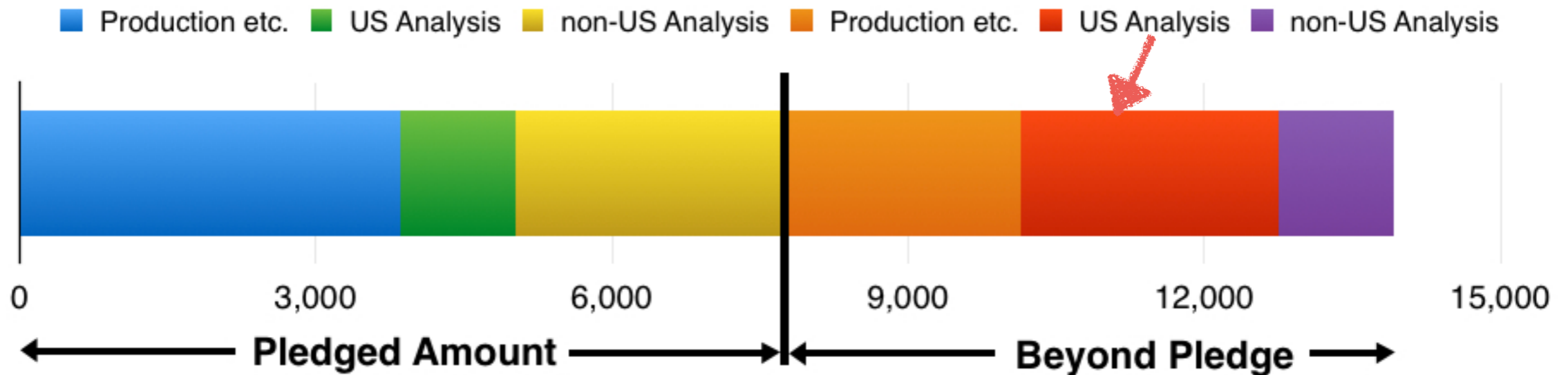


- T2_US_UCSD - 8.69% (115,360)
- T2_US_Wisconsin - 6.76% (89,778)
- T2_US_Nebraska - 4.70% (62,316)
- T2_US_Florida - 3.19% (42,400)
- T2_IT_Legnaro - 2.63% (34,895)
- T2_FR_CCIN2P3 - 2.10% (27,844)
- T2_FR_IPHC - 1.78% (23,671)
- T2_UK_SGrid_RALPP - 1.51% (20,089)
- T2_IT_Rome - 1.35% (17,854)
- T2_CH_CSCS - 1.25% (16,551)
- T2_DE_DESY - 8.69% (115,299)
- T2_US_Purdue - 6.51% (86,419)
- T2_EE_Estonia - 4.18% (55,433)
- T2_DE_RWTH - 3.03% (40,214)
- T2_ES_CIEMAT - 2.47% (32,720)
- T2_US_Caltech - 1.97% (26,190)
- T2_BE_IHHE - 1.71% (22,652)
- T2_FR_GRIF_IRFU - 1.50% (19,929)
- T2_BR_SPRACE - 1.34% (17,830)
- T2_FR_IFCA - 0.87% (11,488)
- T2_CH_CERN - 6.96% (92,381)
- T2_UK_London_IC - 4.84% (64,282)
- T2_US_MIT - 3.78% (50,111)
- T2_IT_Pisa - 2.78% (36,944)
- T2_FR_GRIF_LLR - 2.13% (28,286)
- T2_IT_Bari - 1.93% (25,653)
- T2_BE_UCL - 1.63% (21,641)
- T2_US_Vanderbilt - 1.50% (19,928)
- T2_UK_London_Brunel - 1.27% (16,829)
- plus 27 more

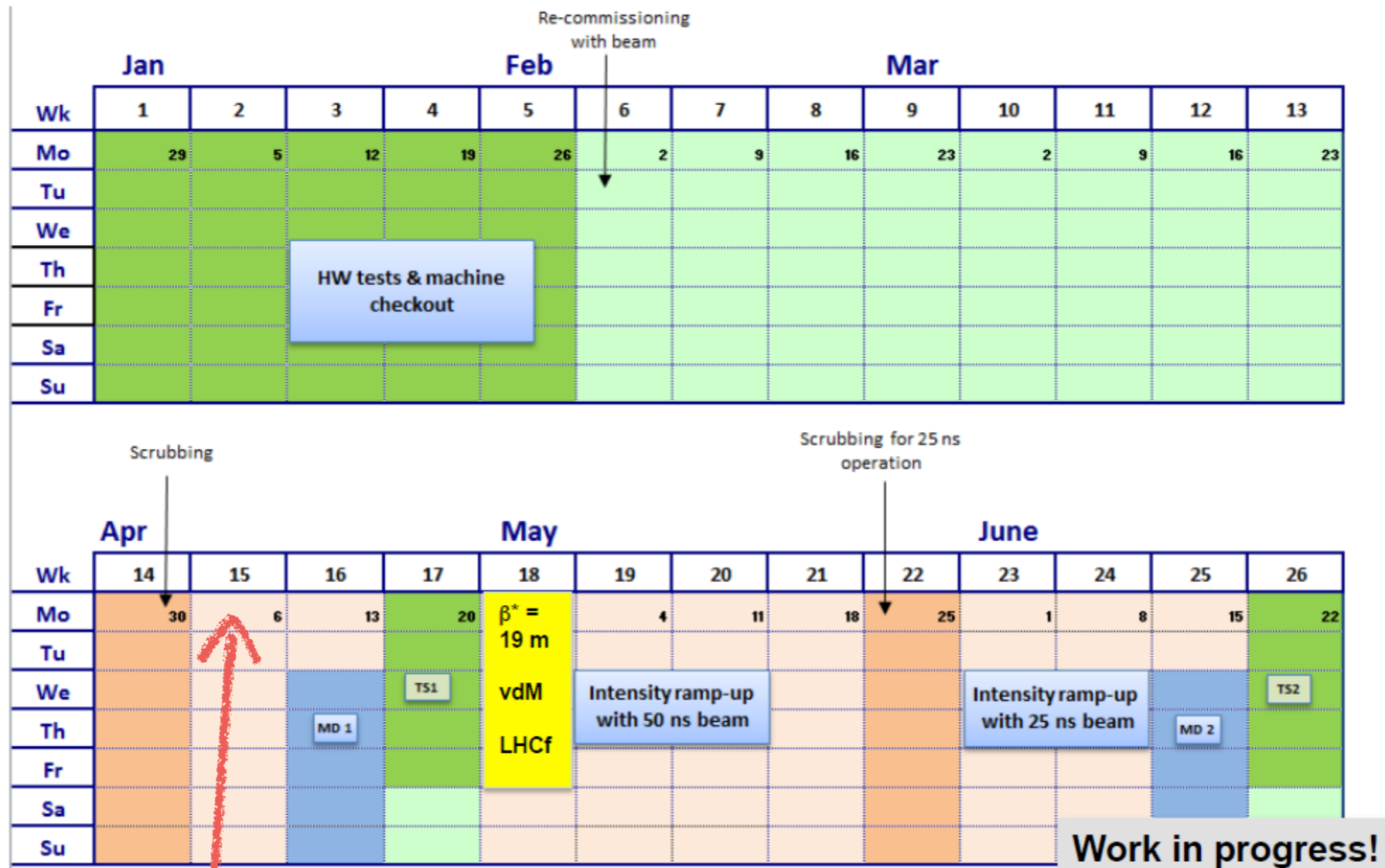
Unsurprisingly, CMS uses US resources well beyond the pledge:



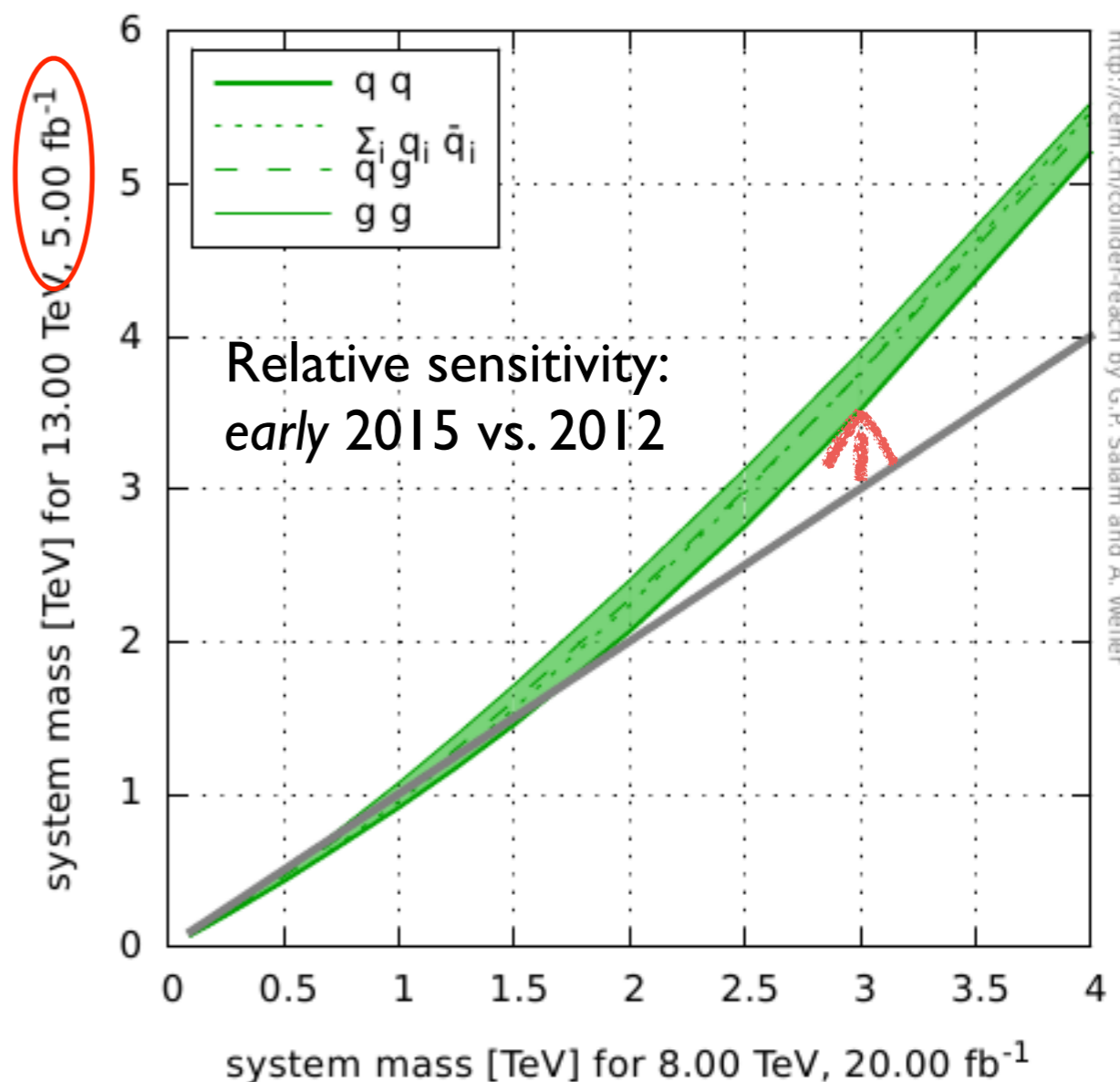
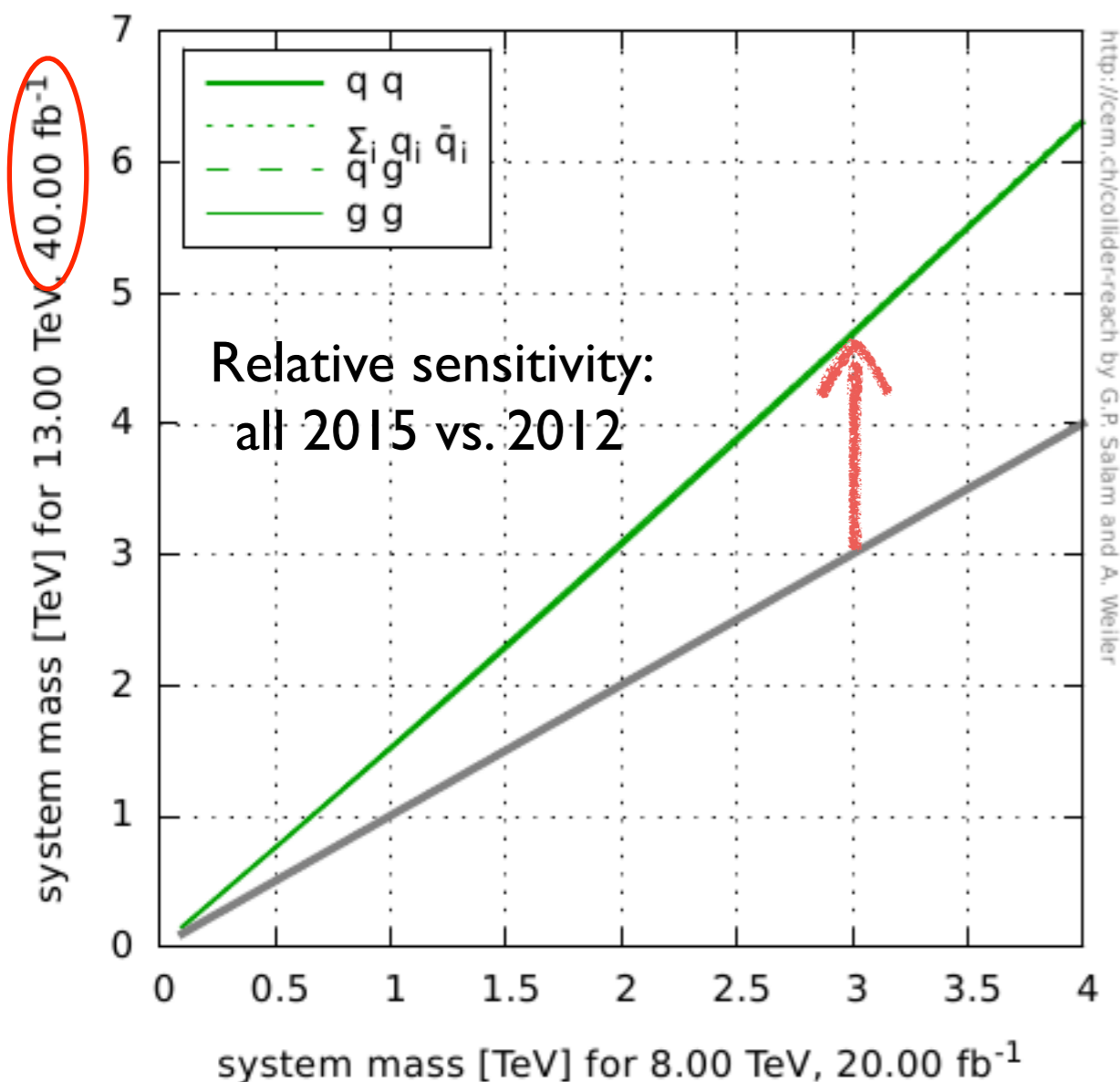
How are these resources used? Best estimate (by J.Letts):



Early 2015 schedule → “Nominal” 25 ns Physics from July’15

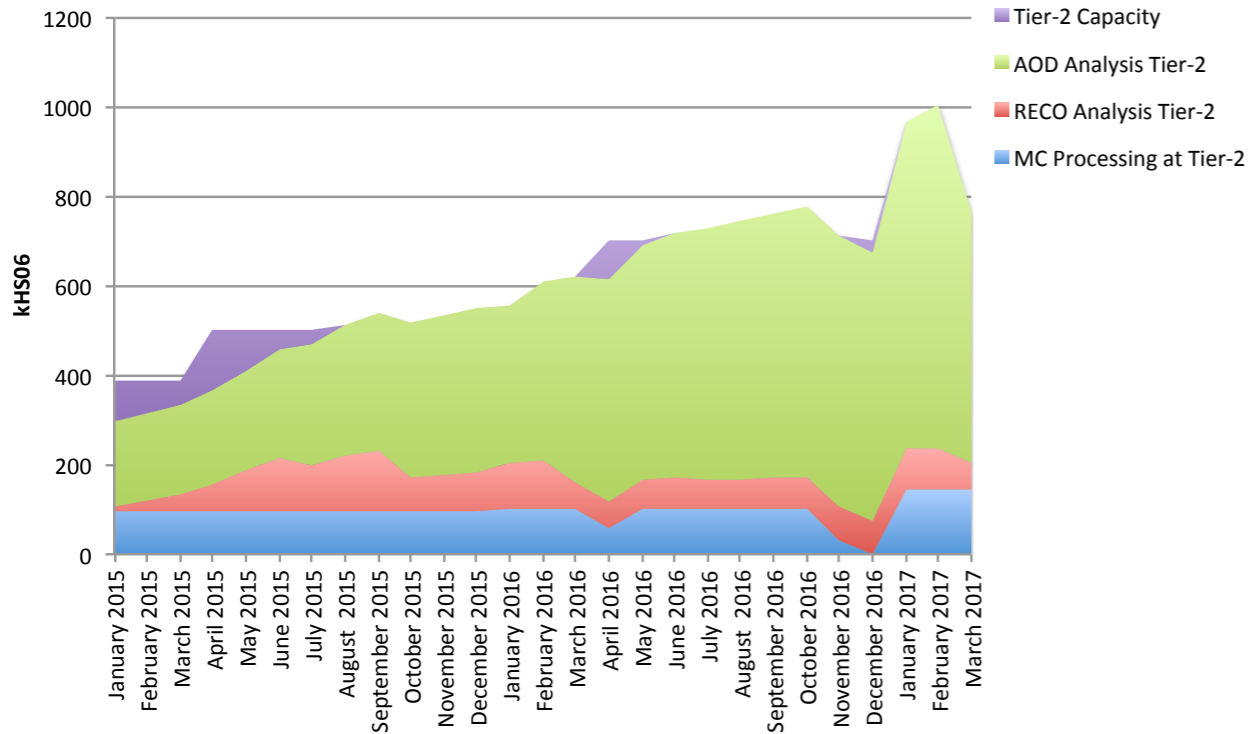


- ▶ First Monday of April: LHC collisions resume
- ▶ Higher beam energy, higher beam intensity, higher pileup rate
- ▶ Larger trigger rates, larger events and processing times

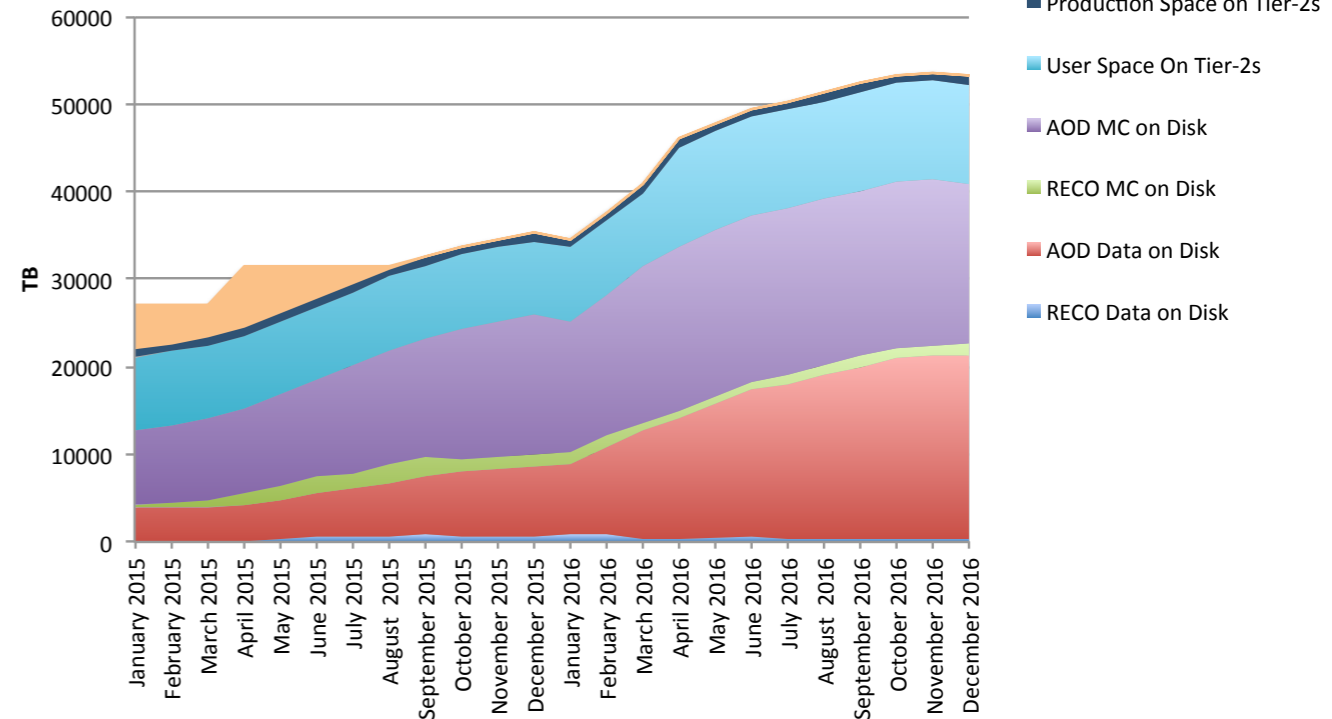


- ▶ We could see new physics in the first year — or first months
- ▶ There is no next major step in collider energy on the horizon
- ▶ We must be ready and we must execute!

Tier-2 CPU



Tier-2 Disk



▶ Resource needs increase as LHC resumes operations and data accumulates once again

▶ Year-to-year increases relatively moderate at T2, as we assume they function as data caches, older data doesn't stay around long

▶ Note: very different at T1 — ~70% CPU increase in 2015, ~30% disk increases in 2015 and 2016

	2014	2015	2016
CPU	+14%	+25%	+40%
Disk	+4%	+16%	+20%

- ▶ Currently expect to maintain build-to-cost model with \$250K/site/year. This should be more than sufficient to meet CMS needs.
- ▶ We continue to be CPU-rich, already have deployed far beyond the ~20 kHS06/site needed by CMS in 2016
- ▶ Cost of disk drives what we can deploy
 - ▶ Have used the LHC shutdown to invest in disk, with the goal of getting ahead of the CMS resource needs
 - ▶ Sites added “extra” 0.5 PB last year and will add 0.5 PB this year
 - ▶ This will allow us to deploy 1 PB/site that is not pledged to CMS but is available to US CMS users through /store/user etc.
 - ▶ Can be accessed via AAA, space management needs development effort
 - ▶ To maintain the extra 1 PB/site and meet CMS request, need to deploy *total* of 2.2 PB/site in 2015 and 2.4 PB/site in 2016 — quite do-able!
- ▶ Likely deployment strategy is then to add disk as described above and beyond that buy as much CPU as possible

- ▶ But: the CMS computing resource request is designed to fit into a flat-ish budget, under the assumption that we can take the operational measures to get the job done with resources available
- ▶ This leads to an operational plan with ~no contingency!
 - ▶ Example: end-of-year reprocessing of that year's data can only be completed in time with HLT, which is only available in a particular time window → all reprocessing ingredients must be on time
 - ▶ HLT is as big as the entirety of CMS TI!
- ▶ CMS will be eager to have access to opportunistic resources
- ▶ Your site can make a difference!
 - ▶ Now more than ever, we need your help in finding extra resources on your campus that can be incorporated when CMS needs them
 - ▶ Have had some notable successes here already, but we need more

- ▶ For Run 2, we will be operating a much more flexible system of distributed facilities that will help meet the computing needs
 - ▶ Different tiers are becoming more similar in functionality
 - ▶ Wide-area data access — AAA — is a key enabling technology
- ▶ Many examples of new operational modes and capabilities
 - ▶ Prompt reconstruction at T1, use of HLT during shutdowns (and between fills?), global scheduling/prioritization through glideinWMS, improved user analysis tools (CRAB3), dynamic data management
 - ▶ Oli, Tapas and Christoph will discuss some of these today
- ▶ Thanks to tools such as AAA, Parrot, and CVMFS, CMS is in a much stronger position than ever to make use of any opportunistic resources that sites can provide to us, and we might well need them
 - ▶ E.g. if we miss the HLT window, do re-processing at T2 sites?

- ▶ CMS computing has an extensive set of milestones for this year; US T2 milestones are aligned with the major ones of CMS
- ▶ 1 April: Deployment of T2 computing resources sufficient to meet CMS resource request for 2014
 - ▶ Timed to start of WLCG resource year; we have done this
- ▶ 30 April: Achieve full participation of T2 sites in CMS data federation
 - ▶ Also done, of course
- ▶ 30 June: Demonstrate appropriate fraction of full production scale of CRAB3 distributed analysis tools at T2 sites: 5K cores, 50K jobs/day, average 80 users/week
 - ▶ The numbers are 25% of the total CMS goal
- ▶ 31 October: Complete participation in exercise of full system for organized production: 5K cores at T2 sites for simulation
 - ▶ Also 25% of the total CMS goal
- ▶ All of these milestones should be readily achievable

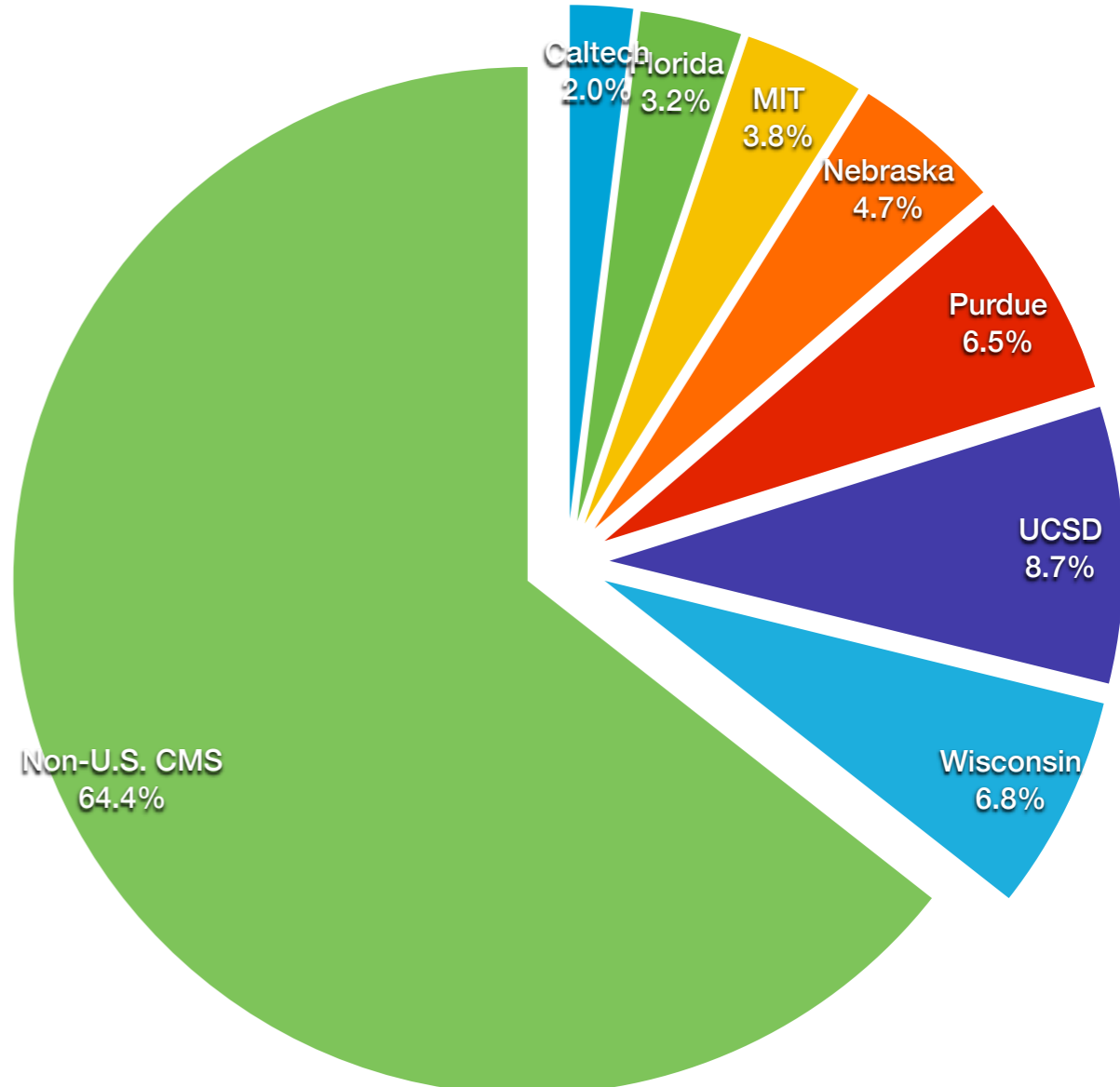
- ▶ One of the reasons the US T2 sites are so successful in CMS is that we try to stay at the leading edge of technology
- ▶ US T2's are always the first to try anything new in CMS computing
- ▶ Today each site will be discussing how they are getting ready for the future, in terms of technologies and otherwise
- ▶ Brian will be offering us some challenges for us to meet over the next few months
- ▶ And we'll hear about specific projects underway at sites that will be of interest to everyone:
 - ▶ Jeff on automatic failover inside FUSE
 - ▶ Manoj on experience with whole-node scheduling
- ▶ We're always happy to put such a presentation into a monthly US T2 meeting, too

- ▶ Had hoped for a presentation on experiences with 100 Gbit networking today, but we are not quite ready yet
 - ▶ Thanks to sites that are starting to explore this
 - ▶ To my knowledge, almost every site has a plan to get to 100 Gbit
 - ▶ This may soon become our baseline network connectivity!
- ▶ We now have the technology to give priority to “national groups” of users on suitable portions of sites
 - ▶ We first talked about this idea at our 2010 workshop at FNAL!
 - ▶ US users should have priority on compute resources above our pledge to CMS
 - ▶ However, this apparently needs a lot of work on the part of sites to partition the cluster appropriately
 - ▶ I am not sure that it is worth it, but it might be necessary in case of a resource crunch

- ▶ The US CMS Operations Program had its annual external review three weeks ago, and as usual it went well for the program as a whole and for software and computing in particular
- ▶ This is the agencies' annual opportunity to reveal budget guidance for the next few years — may not be the actual funding, but a sense of what we need to plan for
- ▶ Scenarios offered were either flat-flat or a 5% cut
 - ▶ Neither one has pretty outcomes, and a budget reduction could endanger the program as a whole
 - ▶ We could see a shifting of priorities away from the LHC this year....
- ▶ Even though it is agreed that T2 sites are a good investment, reducing the T2 hardware budget is an easy cut to make
 - ▶ We have strongly articulated the consequences of reducing site staff or reducing the number of sites that we operate

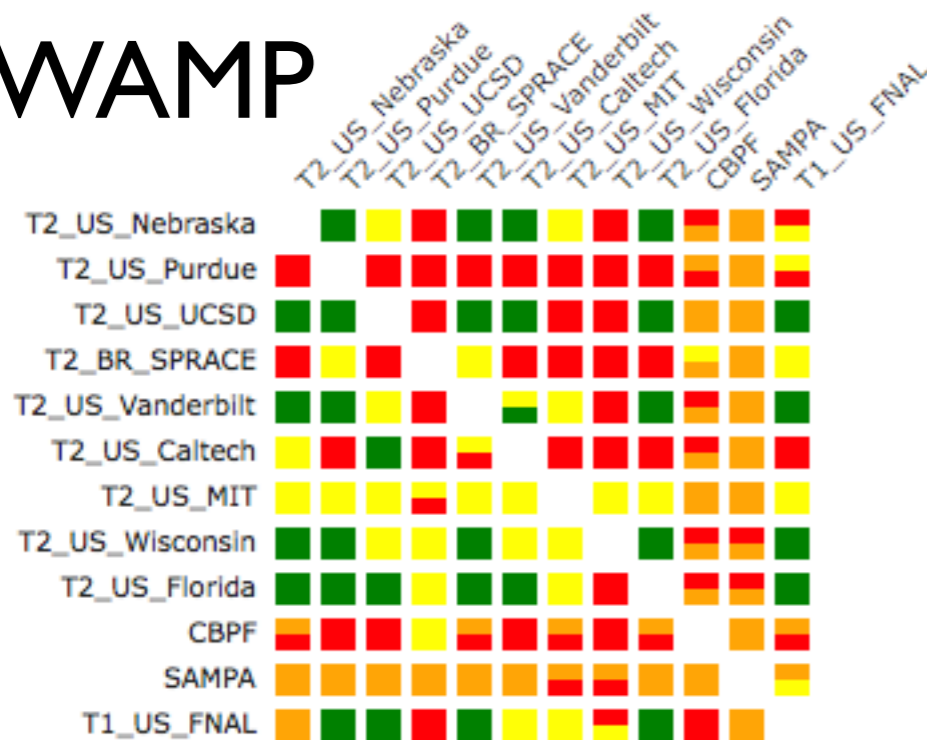
- ▶ The 2012-13 T2 cost/benefit review gave a set of metrics that could be used to evaluate the value delivered by sites
- ▶ Sites in fact gave a substantial return on investment thanks to the participating universities
- ▶ Should we need to reduce hardware purchases, we would use those metrics to decide how to allocate funds to sites that provided the most value
- ▶ It goes without saying that we want every site to deliver as much value as possible to CMS

2013 Tier-2 CPU Usage

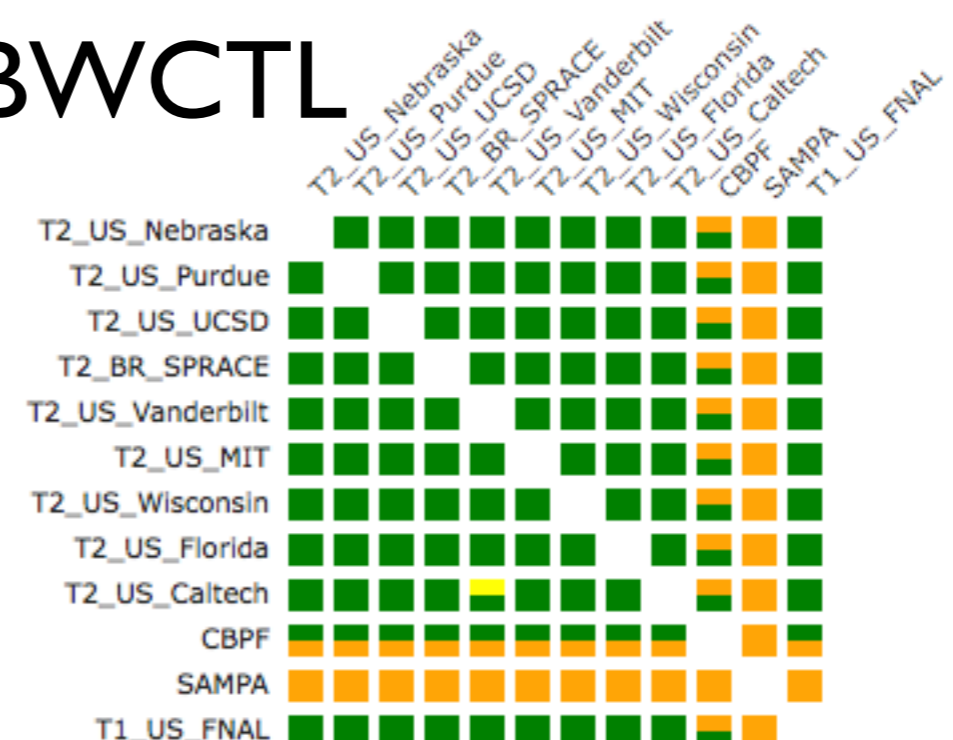


One measure of value; others include quality of site support, contributions to development/operations....

OWAMP



BWCTL



- ▶ These perfSONAR tests look a lot better — thanks!
- ▶ Remember, orange (no data) is worse than red (< 0.5 Gbps)
- ▶ Garhan will give us some pointers on pS tests today
- ▶ (Need to get rid of those sites that aren't ours...)
- ▶ We are looking for some sites to test out more pervasive monitoring of CMS disk usage, see [here](#) for details

- ▶ The excellent quality of work performed at the US CMS Tier-2 sites — over the past nine years! — is a key driver of the successful CMS physics program
- ▶ Discovering the Higgs was exciting, but there is a lot of physics potential in run that starts a year from today and we must be ready for it:
 - ▶ CMS will need all the computing that we can throw at it
 - ▶ We need to respond to operational changes
 - ▶ We must continue to modernize our facilities
 - ▶ All while operating the sites at the highest level of service
- ▶ I'm looking forward to working with you as we get ready for the next era of the LHC
 - ▶ Even if it means that I have 9 PM meetings — only three more of those until I come home!