

Tier-3 Implementation Committee (T3IC)

**US ATLAS Distributed Computing Workshop
University of Arizona**

December 11, 2013

Mark Neubauer

University of Illinois at Urbana-Champaign

Jason Nielsen

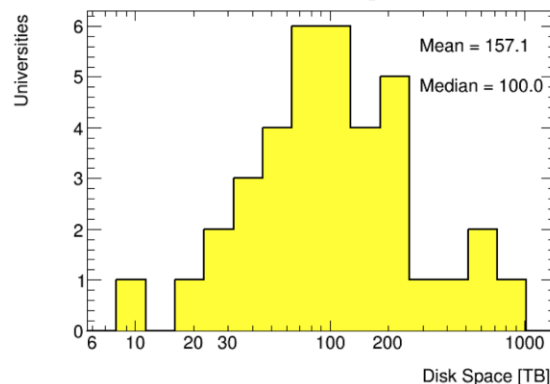
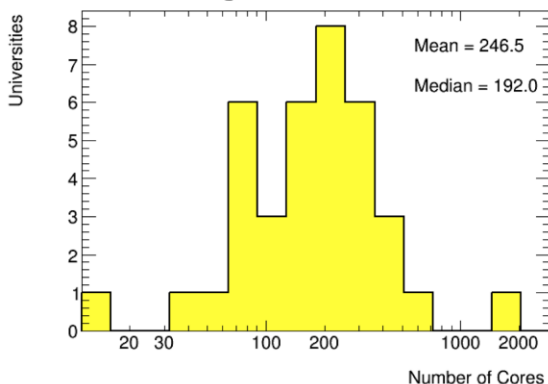
University of California at Santa-Cruz



US ATLAS Tier-3s



- 90% of US universities consider Tier-3 computing as critical to their group's productivity in ATLAS
- The current deployment of the US Tier 3 computing was purchased using ARRA, MRI and university funds



- The majority of equipment was purchased in 2010, w/o replacement
- Two committees reviewed the Tier-3s and produced a report
 - 2009 Tier-3 Task Force. Chair: Chip Brock
 - 2013 Tier-3 Task Force: Chairs: Jianming Qian and Gabriella Sciolla



2013 T3TF Report



- Six recommendations, paraphrased below:
 - ① Support the LT3Cs and make strong case for continued LT3C funding
 - ② Provide mechanisms for Tier-3 jobs to expand onto resources outside of the LT3Cs when the LT3C resources are fully utilized. Invest in technologies that give Tier-3 functionality to institutions w/o an LT3C
 - ③ Support the use of WAN data access for Tier-3s
 - ④ Provide all users a sufficient amount of guaranteed storage space located at Tier-1/Tier-2s (where the batch computing resources are)
 - ⑤ Provide capability for users to direct, upon submission, their output from Grid jobs back to their local storage. Output should be retrievable with minimal delay after the batch job is completed
 - ⑥ Provide documentation and organize comprehensive tutorials to train interested users on how to take advantage of new analysis resource
 - Report also included needs estimates and possible solutions
-



Tier-3 Implementation Committee



- Formed in November 2013 by US ATLAS Operations Program
 - Need for such a committee driven primarily by
 - Desire from the US ATLAS Operations Program to formulate a specific **plan-of-action** to implement the core T3TF report recommendations, with estimates of any Ops-supported resources required
 - ❖ The technological solutions discussion in the T3TF Report was deliberately short and vague, so as to leave that work to our committee where the expertise in the relevant areas is concentrated
 - Desire from the funding agencies to receive a **clear message** from US ATLAS regarding needs and how future Tier-3 funding should be spent
 - Chairs: Mark Neubauer (Illinois) and Jason Nielsen (UCSC)
 - Members: D. Benjamin, K. De, M. Ernst, R. Gardner, G. Sciolla, E. Varnes, T. Wenaus.
 - Ex-officio: J. Cochran, S. Rajagopalan (Ops), C. Brock (US IB)
-



T3IC Charge



- To carry out a comprehensive study that proposes a cost-effective implementation plan to address the T3 challenges. [...] Among the questions and issues you should address are:
 - Provide a best estimate of the computing capacities required to satisfy the physics analysis activities in the U.S. over the next five years
 - Address how far the existing T3 infrastructure goes to accomplish these goals and how to make better use of all existing resources to support the U.S. physics analysis requirements
 - Address the incremental capacities that are needed to provide adequate support for U.S. physics analysis

To be addressed in an interim report by **December 15, 2013**



T3IC Charge (cont.)



- [...] Among the questions and issues you should address are:
 - Identify potential implementation plans that address the T3 needs of U.S. ATLAS physicists and evaluate their cost-effectiveness. Your evaluation of the cost-effectiveness should take into account and identify any synergies, efficiencies, institutional or laboratory leveraging, potential for existing or additional funding sources, and possibly other intangibles. A final ranked comparison table should summarize the cost effectiveness of these plans.
 - Identify how the T3 computing resources and personnel would be managed for (centrally, institutionally, etc) for each plan. This is important for any solutions that require Operations Program funds.
 - Identify what can be accomplished within the current Operations Program budget guidance and prioritize additional requests if supplemental funding materializes (from Operations or other sources).

Full report due by **February 28, 2014**



T3IC Activities



- The T3IC has met four times (~weekly) since being formed
- It has been difficult at times to get everyone together for the weekly meeting due to busy travel schedules (particularly of late), but very productive discussions at meetings
- Use of Google Drive to allow Committee members to upload material and other members to browse is working well
 - Can easily pull this material for our reports and status talks
- Primary activities thus far have been consideration of
 - Estimate of resources required for physics analysis
 - How extra resources help satisfy the analysis needs
 - Cost-effectiveness of possible solutions
 - Detailed implementation plans

} First three charges



Resource Estimate



- To address the first three points in the charge regarding the required resources for physics analysis, Doug, Kaushik, Mark, and Anyes were asked to provide more information on their typical analysis workflow. In particular,
 - 1) *required (or actual) turnaround time (wall clock time)*
 - 2) *location of input*
 - 3) *size and location of output*
 - 4) *typical number of remote or local job slots used*
 - 5) *additional requirements or constraints on remote vs. local data**and, for each of the following analysis stages:*
 - 1) *skim of group D3PD to secondary D3PD (or small ntuple)*
 - 2) *loop over small ntuple to tune cuts, etc., and make plots*
 - Several of these are now in the T3IC Google Drive area
-



Resource Estimate (cont.)



- To scale to 100 fb^{-1} of 14 TeV data, we need best estimates
 - To be useful, estimates should be accurate to $\sim 20\text{-}50\%$.
 - We guess the scaling will be less than a simple x4 factor, but probably more than x2.
 - What are some things we could do differently with x4 data?
 - We will need a multiplicative factor for the number of workflows, including information on how they are shared among analyzers
 - ❖ A rough way to do this is to count the number of analyses being performed in the US, which is not easy. We have a count of the number of notes with US contacts; maybe this is close enough
 - What about MC generation, toyMC (pseudo-exp's for stats interp?), etc?
 - ❖ MC production on Tier-3 is probably small, given ATLAS central prod for MC
 - ❖ Toy MC may not be small, also others like ME calculation, MVA training, etc
 - ❖ Could lead to non-linear scaling w/ lumi. Decided to do linear estimate 1st
 - Analysis model changes from AMSG need to be considered carefully



Resource Estimate (cont.)



- In Dec 5 meeting, two of these workflows were discussed
 - UC Irvine (Anyes): Detailed breakdown on a SUSY multilepton analysis
 - Analysis workflow shared among 5-6 UCI analyzers who produce different results
 - 1st step is to produce skims from ~200 TB of NTUP_SUSYSKIM input
 - Output of these grid jobs is ~0.8 TB, 60% of which is MC; Output brought to T3
 - 2nd step is to produce a private ntuple for plotting
 - Desire ~2 hr turn-around, so that multiple passes can be made per day, requiring several 100 parallel jobs using 100 CPU-hours to accomplish
 - 3rd and final step is the plotting or limit calculation step, always on the Tier-3
 - Duke (Doug): Overview of Top/SM dilepton analysis (AIDA)
 - This analysis also performs a skim down to private ntuples that are about 6% of the NTUP_COMMON size. This skimming takes 800 CPU-hours on the Grid.
 - Then the derived ntuple analysis on the Tier 3 takes about 6 CPU-hours. Input/output dataset sizes will be provided soon.
 - Kaushik get SUSY jet+MET input; Mark on UIUC analysis flow
-



Resource Estimate (cont.)



- Some very preliminary conclusions based on workflows studied
 - **Space at the Tier-3 is not a limitation** now or in the future. Access, latency, and reliability bigger concerns to users than volume
 - **Computing power is the limiting factor**. At the moment, the analysis framework is tuned so that the turnaround time is at the upper limit of convenience.
 - ❖ With 4x the data, it may not be possible to achieve a reasonable turnaround time on the Tier-3 alone. This implies that either the analysis needs to be more efficient or we need to find computing cycles outside of the Tier-3.
 - One concern expressed about moving some jobs off of the Tier-3 into other facilities: losing a lot in **reliability** and workflow **control** (e.g. checking files for failures, following up on lost files, etc.).
- We need a wider range of analysis workflows before drawing preliminary conclusions about resource estimates



Satisfying Analysis Needs



- Need to show how extra resources (beyond current Tier-3 resources) help satisfy the analysis needs.
- The workflows will point to specific cases, but we should also include the laundry list of enabling technologies. It could include some “best practices” examples of how those technologies can be used to enable analysis
- This list will be used to motivate in part the computing requirements, when the technologies imply a change from current resources



Cost-effectiveness



- The “cost-effectiveness” criterion likely boils down to bang/\$\$.
- The baseline expectation is that this would be the cost of extra hardware, with little or no additional personnel cost
- To 1st order, about how much additional resources should go into L3TCs for Run 2 (to 2nd-order, exactly where and when) vs. using the beyond-pledge resources (BPRs) at the facilities
 - Of course, BPRs are already being used for “US physics”, but more for official MC production and not so much for T3-like workflows
 - Since the BPRs exist right now and could be brought into play with ~zero additional cost, a concern is that all plans that maximize cost-effectiveness might favor BPRs exclusively as the “T3 solution” if they do not take into account some less “tangible” aspects of the L3TCs (most notably, in-kind contributions, convenience, and accessibility)
 - ❖ Needs to be defined and quantified before constructing the set of plans



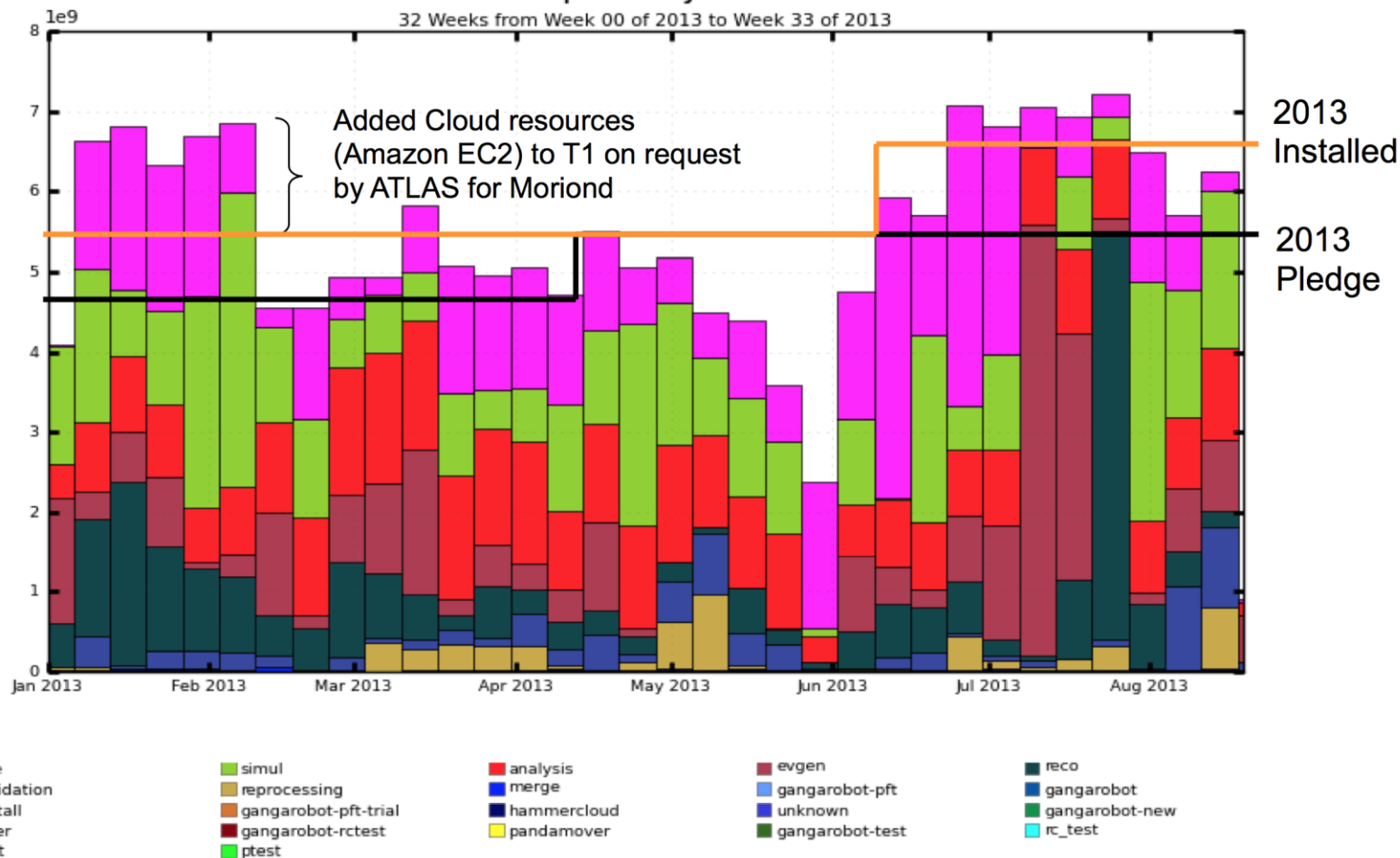
Tier-1 Usage (2013)



dashboard

CPU consumption All Jobs in seconds

32 Weeks from Week 00 of 2013 to Week 33 of 2013

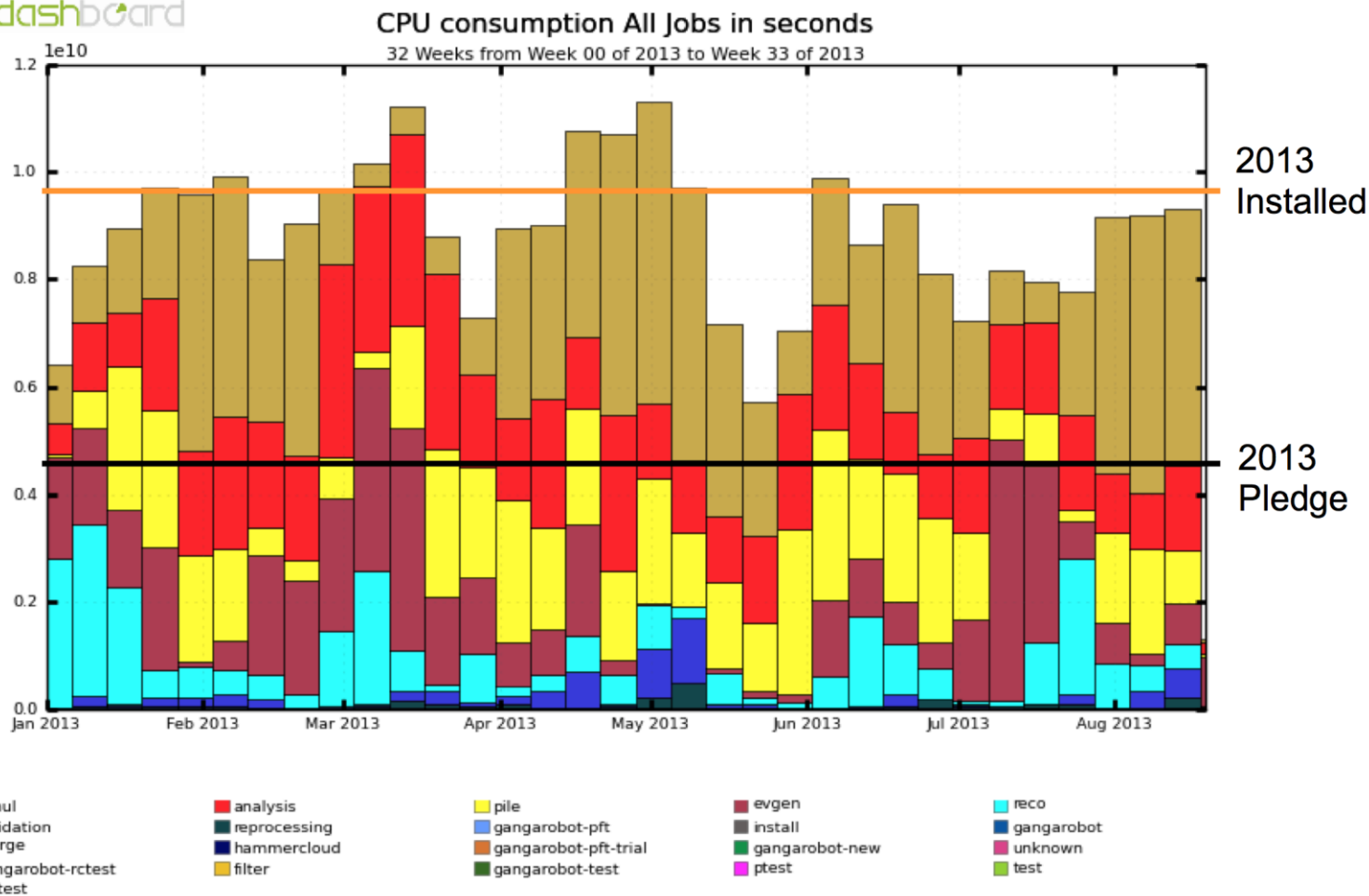




Tier-2 Utilization (Jan-Aug 2013)



dashboard





Implementation Plans



- There are a few contributed implementation plans so far in the T3IC Google Drive area
 - Agile infrastructure
 - ATLAS Connect / Condor flocking / FAXBox
 - PandaAtTier3
 - More detailed plans for future Tier-3 functionality forthcoming
 - Fortunately, we have a couple of very promising enabling technologies for future Tier-3 workflows which should not require an large amount of additional manpower
 - Suggest that we include these in our planning; ultimately, users will settle into using the ones that work best for their physics output
 - It was agreed that we would also look at the resources available at the FNAL LPC-CAF
-



Summary



- The T3IC is now very active
- We need a wider range of analysis workflows before drawing preliminary conclusions about resource estimates.
- The first workflow studies indicate that CPU will be a scarcer resource than disk. The CPU availability will probably drive the implementation plans.
- Anecdotal evidence from users indicates that job failures or larger-than-anticipated latencies could push people away from distributed computing in favor of Tier-3.
- Development of implementation plans well progressed; Documentation of these plans progressing in parallel with the resource/cost estimates
 - Want use of BPRs and user feedback to be available by Agency review



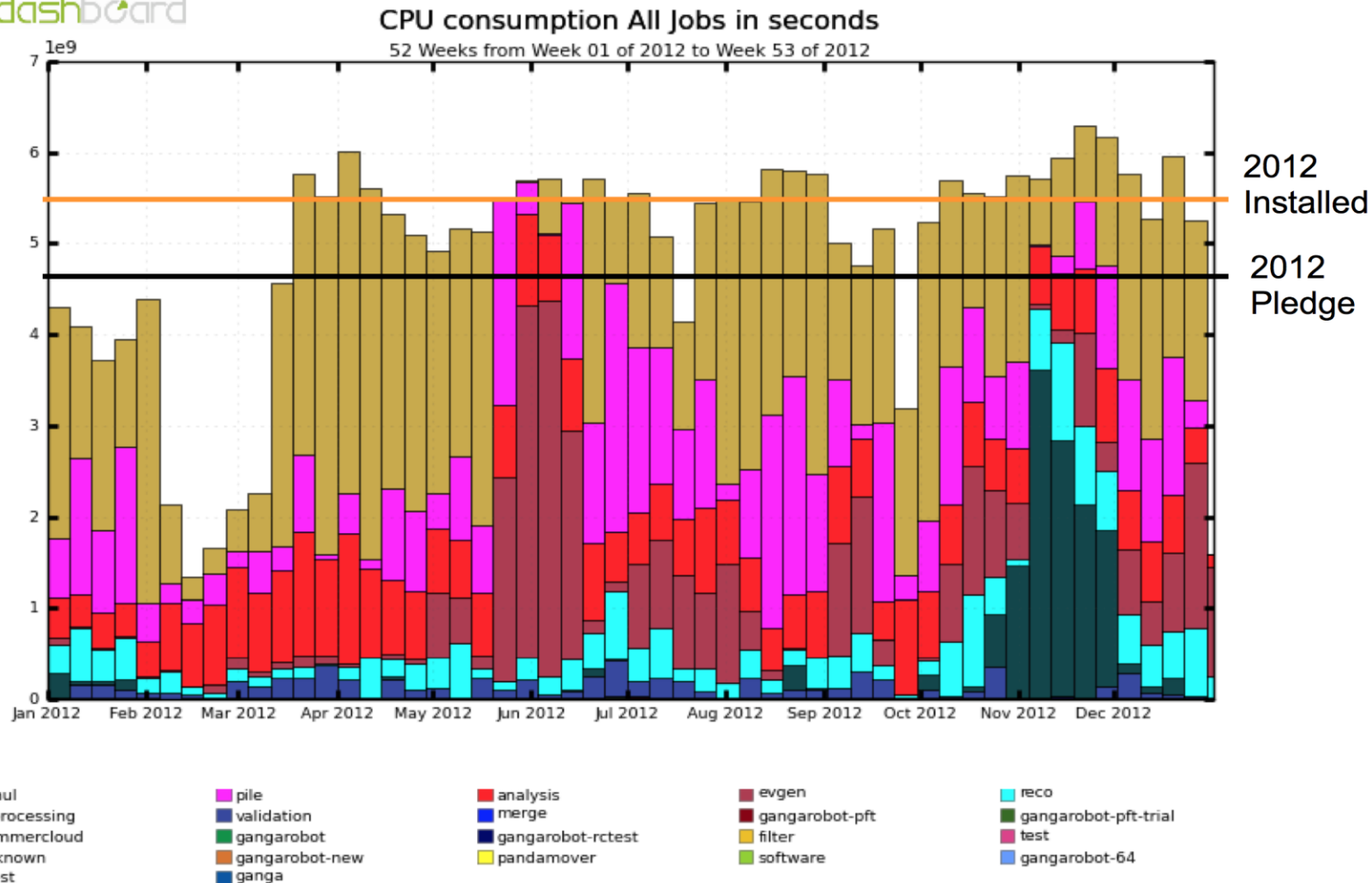
Bonus Material



Tier-1 Usage (2012)



dashb^oard





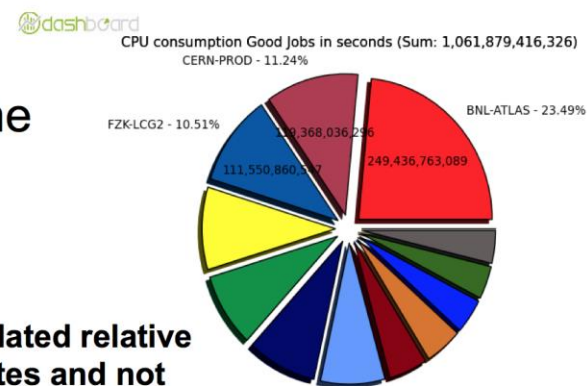
Tier-1 Comments



- ❖ U.S. ATLAS pledge is based on MoU share = 23% of total capacity requested by ATLAS.
- ❖ U.S. ATLAS installed capacity = pledge + 20%. Additional capacity targeted for U.S. physicists.
 - Beyond pledge capacities are managed and allocated by U.S.
- ❖ Tier 1 runs at full load despite additional capacity.
- ❖ 2013 has not seen a decrease in usage despite shutdown

Reprocessing, ongoing analysis, simulation (incl. pile-up) in preparation for the upcoming run at ~14 TeV and pileup ~50.

Contribution (2012) is calculated relative to resources provided by sites and not relative to the MoU share)





Tier-2 Utilization (cont)

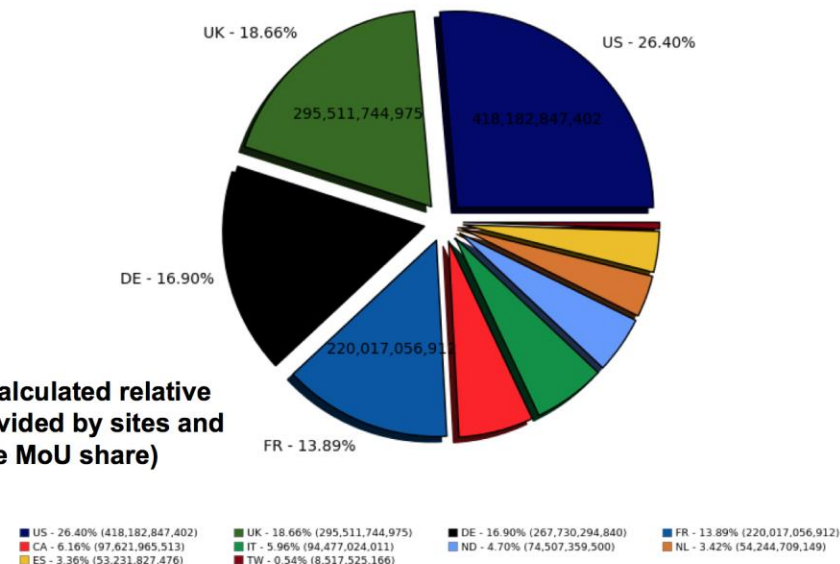


dashboard

CPU consumption Good Jobs in seconds (Sum: 1,584,042,354,944)

Baseline: Pledge based on MoU share (23% of total capacities requested by ATLAS + capacities from uncommitted funds leading to significant resources beyond pledge.

Contribution is calculated relative to resources provided by sites and not relative to the MoU share)



- ❖ Despite additional capacities, Tier 2s are running at full load, even during the shutdown period.
- ❖ The U.S. Tier 2 are amongst the best performers amongst all international Tier 2 centers.



Comments on US ATLAS BPRs



- While they are $\sim 100\%$ utilized, a question of optimizing future hardware deployment vs. software development (see Torre's talk)
- Their existence is very defensible (see utilization plots), but they are in constant need of defending
 - US has large BPR, but CPU consumption by US is in line in MOU share
 - Regular talking point of funding agencies
 - ❖ They are very sensitive to the prospect that the US is subsidizing computing for ATLAS beyond their "fair share" (not the case)
 - There is some risk of reduction if we do not make sufficient use of these resources for US physicists **and** provide clear accounting
- I've focused on CPU, but disk over-pledge is also substantial



Usage of US ATLAS BPRs



- Dedicated requests (mostly for MC) to RAC (handled by E. Varnes) from US physicists
 - Seems that requests well accommodated but are sporadic (by the nature of the requests)
 - Q: For production requests, can we do any more than encourage and reassure the community?
- Panda brokerage for jobs identified as being from US people
 - Q: How well has this been working? Can we improve this? How is fair-share among US physicists done?
 - Q: Is there monitoring of this to allow for transparency and feedback for tuning?



Usage of US ATLAS BPRs (cont)



- US Tier-3s
 - At most institutions, T3 hardware was purchased ~3 years ago with ARRA and MRI funds (T3TFv2 Observation 4) → at a level not likely to repeated any time soon
 - ❖ Hardware-wise, US ATLAS T3s are aging w/o funds to refresh
 - ❖ Physics-wise, T3s are critical resources for effective participation by the US in ATLAS (T3TFv2 Observation 1)
 - CPU BPRs could augment the T3 capacity given challenging T3 funding
 - ❖ Panda dedicated queues to access these resources
 - ❖ Using flocking capability in HTCondor
 - ❖ Both of these have been demonstrated to work. Users would like to retain their T3 “look-and-feel”
 - ❖ Q: What would be the fair-share policy? How would the accounting work?
 - Disk BPRs could provide T3/WAN-accessible read/write scratch space (e.g. via Faxbox). Archiving service?