

# UnifiedPandaQueue

Motivation, deployment and consequences

# What's a unified queue?

- Single Panda queue for S/MCORE, LO and HIMEM per physical resource
  - Actual job requirements are set for pilot
    - passed via CE to batch system

File  
ATLAS  
←

ATLAS Grid Information System

/C=DE/O=GermanGrid/OU=LMU/CN=Rodney Walker | Logout

RC Site ATLASite DDMEndpoint PANDA Queue Service Central Services DDM Groups PandaQueue combined resources Docs TWiki OLD JSON

Show 200 entries

First Previous 1 Next Last

give me url of this page **hold shift + click column for Multi-column ordering** VO ATLAS Site PanDA Site Template PanDA Resource PanDA Queue state Final Status Manual HC Switcher type capability rtype CLOUD TIER use newmover deprecate oldmover

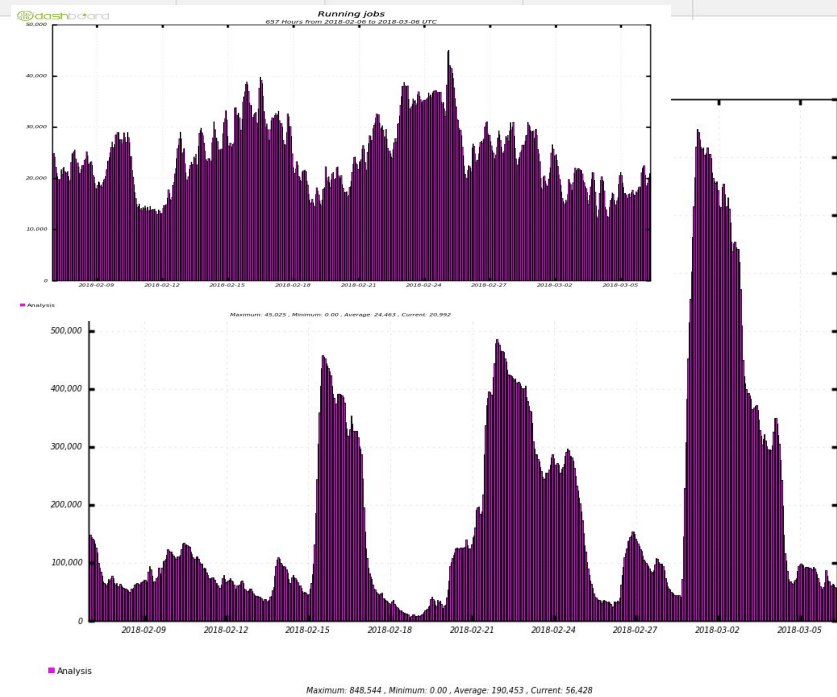
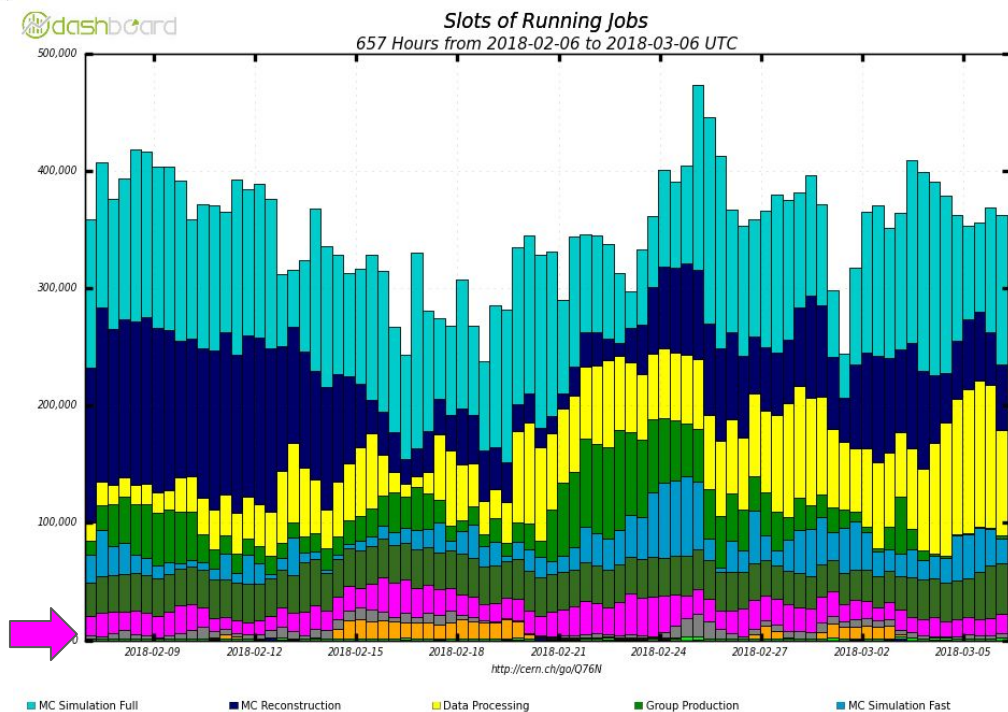
VO	ATLAS Site	PanDA Site	Template	PanDA Resource	PanDA Queue	state	Final Status	type	capability	CLOUD	TIER	use newmover	deprecate oldmover
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	ANALY_MWT2_HIMEM	ANALY_MWT2_HIMEM	ACTIVE	ONLINE	analysis	himem	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	ANALY_MWT2_HIMEM_MCORE	ANALY_MWT2_HIMEM_MCORE	ACTIVE	ONLINE	analysis	mcورهimem	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	ANALY_MWT2_MCORE	ANALY_MWT2_MCORE	ACTIVE	ONLINE	analysis	mcوره	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	ANALY_MWT2_SL6	ANALY_MWT2_SL6	ACTIVE	ONLINE	analysis	score	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	MWT2_HIMEM	MWT2_HIMEM	ACTIVE	ONLINE	production	himem	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	MWT2_HIMEM_MCORE	MWT2_HIMEM_MCORE	ACTIVE	ONLINE	production	mcورهimem	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	MWT2_MCORE	MWT2_MCORE-condor	ACTIVE	ONLINE	production	mcوره	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	MWT2_SL6	MWT2_SL6	ACTIVE	ONLINE	production	score	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	MWT2_VHIMEM	MWT2_VHIMEM	ACTIVE	ONLINE	production	himem	US	T2D	True	true
atlas	MWT2	MidwestT2	MidwestT2_VIRTUAL	MWT2_VHIMEM_MCORE	MWT2_VHIMEM_MCORE	ACTIVE	ONLINE	production	mcورهimem	US	T2D	True	true

Showing 1 to 10 of 10 entries

# Motivation

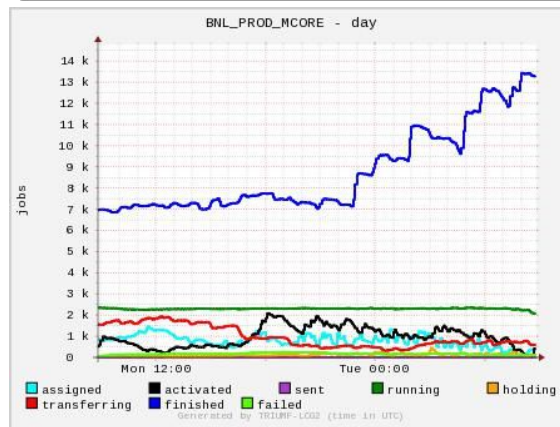
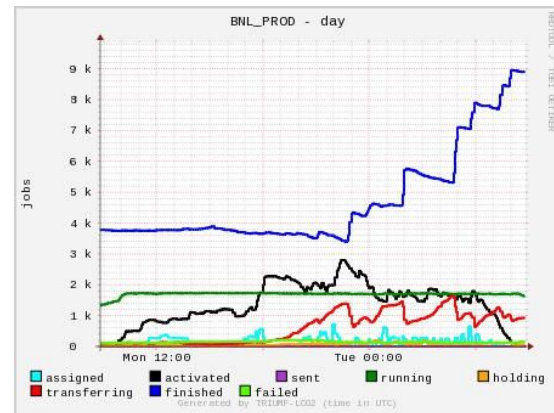
- Gshare and prio decide which job starts next (getJob)
  - is the case now BUT only on single PQ
  - shares between S/MCORE/HIMEM undefined or static
- UnifiedQueue: site runs pilots in order of priority
  - e.g. only submit MSCORE, then only runs MSCORE
    - no more low prio SCORE evgen using resources - follow gshare
- Evgen can run anywhere and pushes out MSCORE
  - currently cap evgen running globally
    - leave some resources empty, if that is all they can run
    - leave many resources empty, if no other job types activated
- Include ANALY
  - Could fill cluster with ANALY, when it has popular data
    - just push prod elsewhere or delay it
    - no need to make replicas or do inefficient remoteio

L1 Share	L2 Share	L3 Share	Actual HS06	Target HS06	HS06 ratio	Queued HS06	Actual share	Target share
Analysis [20.0%]			244,347.76	1,085,905.82	22.50 %	551,660.18	4.50 %	20.00%
Express [3.0%]			140,706.55	162,885.87	86.38 %	36,833.32	2.59 %	3.00%
Production [75.0%]			4,965,947.54	4,072,146.84	121.95 %	22,572,243.52	91.46 %	75.00%



# Challenges

- Need unpartitioned cluster - sites probably like this(Xin)
  - no hard partitioning or intra-VO shares
    - both S/Mcore and prod/analy
  - no loss of resources by not submitting SCORE
- Some local partitioning/shares/limits are to protect site
  - Directio to storage, bandwidth, #connections
  - WN scratch disk io, space
- PandaQueue is unit in monitoring, tests, switcher, ...
  - might need to switch capabilities individually, eg. stop analysis
- ANALY inclusion is tricky
  - pilot or prod proxy
  - need gshare for users



# Brokerage & Monitoring

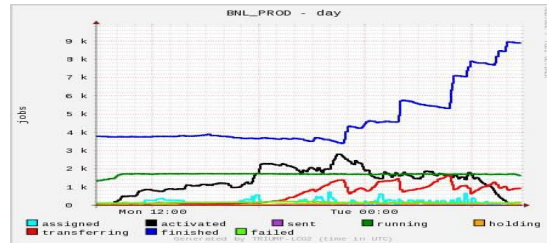
- Unified queue has internal sub-resources used for brokerage
  - see these in the brokerage logs
    - skip site=DESY-HH\_UCORE/MCORE due to core mismatch site:8 <> task:1
- Not exposed on bigpanda
  - maybe could expand UQ to show sub-resources
  - panglia #running is mix of S/MCORE
    - dashboard uses job.corecount, but panglia is so convenient!
- HammerCloud PFTs
  - can run both PFT and PFT\_MCORE
  - can be is\_default, and all works fine - all switched at once
- Accounting ok
  - all job based, not PQ

# Deployment

- Initially only possible in push-mode
  - pre-loaded pilot from aCT has requirements passed to batch via ARC CE
  - aCT can now submit to HTCondor, so ARC CE not required - tested at CERN
- Also possible in pull-mode - any CE, including Cream
  - strict control of pilot streams
  - next talk.....

# Deployment - aCT, ARC CE

- LRZ-LMU\_MCORE, DESY-HH\_UCORE(EL7), DESY-HH\_MCORE, FZK-LCG2\_UCORE, RAL-LCG2\_UCORE, SiGNET-NSC\_MCORE, TRIUMF\_DOCKER\_UCORE, CA-SFU-T2\_UCORE
  - mix of names because Ivan likes UCORE, but making new site was painful
    - clone now attaches DDM endpoints and sets pr/pw protocols
- In use with old PQs removed or long-term brokeroff
  - some care needed to remove T1 or nucleus PQ
- Step towards combined ANALY/PROD queue
  - ANALY\_FZK-LCG2\_UCORE, ANALY\_DESY-HH\_UCORE have sub-resources for MCORE, HIMEM
    - may be use use-cases for MCORE, e.g. proof-lite, deep learning, ART
    - need to address accounting

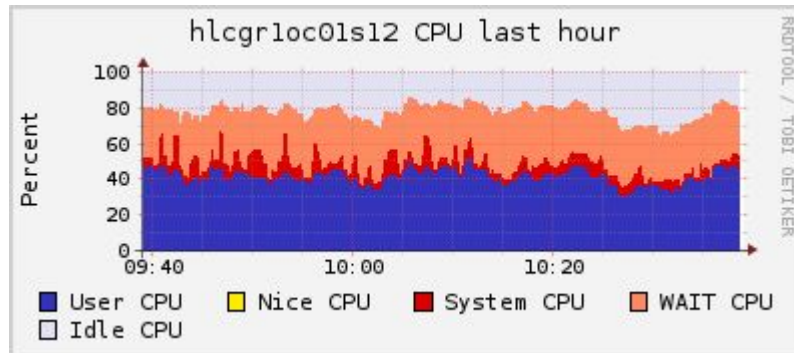


CA-SFU-T2\_UCORE ([862](#))  
corecount (2) 1 ([18](#)) 4 ([844](#))  
MC 16 simul ([757](#)) Reprocessing default ([105](#))  
Single core are repro  
AODMerge\_tf.py ([12](#)) ESDMerge\_tf.py ([6](#))



# Need to limit jobs

- Cluster can fill with single job type
  - high disk io, or ANALY hitting storage directio
  - currently might have hard limit on ANALY jobs
    - better to have limits on physical properties
    - sum quantity over running jobs
- Have iointensity, maxrss, Frontierload
- Now added DISKIO NUMBER(9)
  - "Local disk access measured by scouts  
(totWBytes+totRBytes)/(endTime-startTime)"
  - E.g. PowHeg jobs write/read O(100MB) files continuously



# Optimize job mix for a site

- Batch schedule with RAM as consumable resource
  - can run very himem jobs but will leave cores idle
- LRZ 3200 logical cores, 5.5TB -> ~1500MB/core -> PQ.rsspercore
  - 3GB per physical core, on average, but HT gives 10-20% more HS06
- How to achieve a mix of jobs
  - run some himem jobs, if gshare wants that, but take lomem to optimize core usage
  - up to now have multiple PQs and rely on entropy
- Sum over running jobs:  $\text{RAM/cores} = \text{rsspercore\_run}$ 
  - broker job with  $\text{job.minramcount} > \text{PQ.rsspercore}$  only when  $\text{rsspercore\_run} < \text{PQ.rsspercore}$
  - overshoot to start, then may stabilize - need it more reactive?
  - mean will be correct, but no control over what BS does - hope it is sensible.
- Same story with scratch disk io, directio, Frontier
  - E.g. stop brokering PowHeg when  $\text{sum io} > \text{\#hosts} * (\text{disk rate})$

# Conclusions

- Several large sites with prod UQ
- Already following shares better, but not enough to remove evgen caps
  - need to deploy to more sites
- Including ANALY brings many advantages, but is tricky
  - no immediate need to make replicas or do inefficient remote io
    - only after the obvious improvement of pushing out prod
  - ready for MCORE/HIMEM
  - discuss this week
- Will need to replace site limits and entropy 'protection' with proper limits