



# Machine/Job Features TF

**Andrew McNab**  
University of Manchester,  
GridPP, and LHCb



# Overview

- MJF aims
- Task Force mandate
- Specification
- Implementations
- SAM probe
- Rollout
- Interruptible jobs
- Benchmarking and performance studies

<https://twiki.cern.ch/twiki/bin/view/LCG/MachineJobFeatures>

## Aims of Machine/Job Features

- A common API that jobs can use to discover the parameters of their environment
  - eg wall clock time limit
- Otherwise requires a patchwork of environment variables and command call-outs
  - Different for each batch system: qstat etc
  - Not available in VM-based environments
- `$MACHINEFEATURES` and `$JOBFEATURES` given locations of the key/value files
  - Local directories or HTTP(S) URLs

# Machine/Job Features Task Force mandate

*“The Scope of this task force is divided into several subsequent steps*

- Check the completeness of the proposal for machine/job features*
- Coordinate implementations used in WLCG and an interface for its usage to the VOs*
- Provide means to monitor the correctness of the provided information*
- Plan and execute the deployment of those implementations at all WLCG resources”*

<https://wlcg-ops.web.cern.ch/ops-coordination-task-forces/machinejob-features>

# Specification

- Several iterations, starting from the HEPiX virtual machines working group, then into this WLCG TF, as talks/Twiki pages
- Task force agreed a specification at the start of 2016 MJF
  - The set of key/value pairs to publish
  - How jobs can get the key/value pairs
- Published as HEP Software Foundation technical note (HSF-TN-2016-02)
- Consistent definitions with APEL and Infosys TF

## Key / values

- The technical note has the full list with definitions
- Sites should supply them if they know the value (eg HS06)
- Values can typically be discovered from batch system, with OS values as a fall-back
- shutdowntime allows sites to declare a cut-off when draining

### \$MACHINEFEATURES

total\_cpu  
hs06  
shutdowntime  
grace\_secs

### \$JOBFEATURES

allocated\_cpu  
hs06\_job  
shutdowntime\_job  
grace\_secs\_job  
jobstart\_secs  
job\_id  
wall\_limit\_secs  
cpu\_limit\_secs  
max\_rss\_bytes  
max\_swap\_bytes  
scratch\_limit\_bytes

# Machine/Job Features Task Force mandate

*“The Scope of this task force is divided into several subsequent steps*

- *Check the completeness of the proposal for machine/job features ✓*
- *Coordinate implementations used in WLCG and an interface for its usage to the VOs*
- *Provide means to monitor the correctness of the provided information*
- *Plan and execute the deployment of those implementations at all WLCG resources”*

<https://wlcg-ops.web.cern.ch/ops-coordination-task-forces/machinejob-features>

# Implementations

- PBS/Torque, HTCondor, and GridEngine scripts exist in GitLab and as RPMs
  - Common code where possible; same ideas
  - Basic mjf-onlymf RPM works for any site but only supplies \$MACHINEFEATURES
  - All use /etc/sysconfig/mjf and /var/run/mjf for fine tuning
  - But will work sensibly out of the box
  - MJF RPMs in WLCG SL6 YUM repo
- Vac/Vcycle always supply MJF directories to their VMs
- See <https://twiki.cern.ch/twiki/bin/view/LCG/MachineJobFeaturesImplementations>



# Machine/Job Features Task Force mandate

*“The Scope of this task force is divided into several subsequent steps*

- *Check the completeness of the proposal for machine/job features ✓*
- *Coordinate implementations used in WLCG and an interface for its usage to the VOs ✓*
- *Provide means to monitor the correctness of the provided information*
- *Plan and execute the deployment of those implementations at all WLCG resources”*

<https://wlcg-ops.web.cern.ch/ops-coordination-task-forces/machinejob-features>

# SAM probe

- WN-mjf.py script in the MJF GitLab repo
  - Runs inside jobs to look for MJF keys/values
- Sites with MJF pass tests when script is run
  - By hand inside jobs
  - Or in the LHCb ETF (SAM) service
    - Viewable in check\_MK  
[https://etf-lhcb-prod.cern.ch/etf/check\\_mk/view.py?view\\_name=hostsbygroup](https://etf-lhcb-prod.cern.ch/etf/check_mk/view.py?view_name=hostsbygroup)
- Key to rollout and debugging
  - Provides an objective test as to whether the site has “got there”
- Same script ok for any batch system or experiment
  - That’s the point of MJF!

# Machine/Job Features Task Force mandate

*“The Scope of this task force is divided into several subsequent steps*

- *Check the completeness of the proposal for machine/job features ✓*
- *Coordinate implementations used in WLCG and an interface for its usage to the VOs ✓*
- *Provide means to monitor the correctness of the provided information ✓*
- *Plan and execute the deployment of those implementations at all WLCG resources”*

<https://wlcg-ops.web.cern.ch/ops-coordination-task-forces/machinejob-features>

# Rollout

- This is the last part of the process and has become the hardest
  - Today 14/65 LHCb sites have at least one cluster/CE passing ETF test of MJF
- Circular problem
  - Sites knows experiments can gather (most) MJF information in other ways
  - Experiments won't spend effort on supporting features until most sites deploy
- Working with some sites directly (eg CERN, Tier-1s, Tier-2 volunteers) to get critical mass
- “Killer app” features for sites/experiments?

# Shutdowntime mechanism

- Sites can set `$MACHINEFEATURES/shutdowntime` to state the Unix time at which the job must stop
- This provides a way for interruptible jobs to be interrupted
- Can avoid the need for draining
  - As reboot approaches, only run interruptible jobs
  - Run those jobs up to the last few minutes before reboot
- LHCb Interruptible Monte Carlo has been demonstrated in production (in VMs) and will go mainstream soon
  - Non-interruptible jobs won't match pilot if too little time left
  - We have already had queries from large sites facing mass reboots whether they can use this mechanism yet

## shutdowntime (2)

- Very applicable to similar job types from other experiments
  - eg ATLAS event service jobs?
- Also for non-traditional quality of service scenarios
  - 15 minutes warning lets me finish current event and upload?
- Adjustable job length allows job masonry with single/multi processor mixes
  - Make single processor jobs finish in groups of eight ...

## Benchmarking and performance

- `$MACHINEFEATURES/hs06` and `$JOBFEATURES/hs06_job` have been key to recent studies
  - Helped identify step change due to Haswell
- DB12-at-boot scenario developed to run an Open Source fast benchmark at boot time
  - Communicated to jobs via MJF `db12/db12_job`
  - Provides a reliable (< 1%) worst-case, fully loaded measure of CPU power available to jobs at the individual host level
  - Distributed in `mjf-db12` drop-in RPM



# Summary

- Long and winding road from original HEPiX VMs proposal
- TF has addressed specification, implementations, and monitoring aspects of the mandate
- Rollout started but uphill struggle
- Break out of circular problem with new uses?
- Availability of interruptible LHCb Monte Carlo jobs in near future: no need for draining?
- DB12-at-boot and other benchmarking/performance studies