



WLCG

Worldwide LHC Computing Grid

WLCG Ops. Coord. Job Allocation and Handling WG

Antonio Pérez-Calero Yzquierdo (CIEMAT & PIC)

WG Kick-off meeting, November 14th 2024





Welcome to the new WLCG Ops. Coord. WM

Welcome to the first meeting of the WG on Job Allocation and Handling!

- This new WG aims at providing a discussion platform for experts from LHC VOs and WLCG sites on matters of compute resource management and workflow to resource matchmaking
- WG membership is (of course) open
 - please join the *wlcg-ops-coord-wg-job-alloc* [e-group](#)
- [Twiki](#) for the WG
- Meeting notes [document](#)



WG Mandate

- **Mandate:** “The goal of the WG is to define the most efficient strategy for job allocation and handling, namely, the most efficient **core count** for batch job allocation and a more efficient handling of different **job classes** with **diverse requirements**. This includes **analysis of the experience accumulated so far** for running multi-core jobs, whole job scheduling, different approach for handling jobs with specific requirements, as for example, high-memory jobs. Based on these studies, the **recommendations** for the sites and LHC VOs should be developed and **presented to the community**. The WG includes experts from the WLCG sites and LHC VOs.”



The meeting today

Kick off meeting objectives and proposed agenda:

- WG presentation
- Summary from the 2024 WLCG Workshop at DESY
- Next steps (next meetings)



A summary of summaries (I)

From the [slides and discussion in Hamburg](#), the main messages were (my personal take):

- ATLAS:
 - Himem jobs exist while other jobs consume much less than 2 GB/core. Objective is to mix them on the WNs relying on CE+BS
 - When all jobs are himem, cores will be left idle, this needs to be accounted for.
 - 16 core slots seems a good choice for evolution of standard common size
 - sites could opt to continue with 8
 - single core still needed
 - whole node scheduling
 - a requirement when using HPCs, but not all payloads can scale
 - not clear benefit on the grid
 - to be tested at a small number of volunteer sites



A summary of summaries (II)

- CMS:
 - Single pilot type multicore, can run on 8, 16, etc fixed size, or whole nodes in the Grid.
 - Internal partitioning deals with diversity of payloads jobs
 - No need for dedicated himem slots
 - CMS can benefit from increasing standard slot from 8 to 16
 - Whole nodes already in use in multiple sites (exclusive sites in the US)
 - Beneficial for CMS due to providing increased flexibility
 - Efficient use possibly linked to extended pilot lifetime
 - Also, whole nodes when running on HPCs, even thinking about multi-node pilots...
- ALICE:
 - Transition from score to mcore tasks achieved at SW level, but support for score must remain
 - Evolution towards combination of multi-sized multi-user payloads to fit the 8-core standard slot
 - CMS-like?
 - In favor of higher than 8-core slots, such as 16, and even whole node, to improve flexibility and ease sysadmins operations.
 - CPU efficiency is to be kept high by means of extended lifetime, also CPU pinning and oversubscription
 - Whole node already in use in HPC and exclusive sites
 - Wider adoption depending on cgroups v2 (EL9, HTCondor 23)



A summary of summaries (III)

- Sites' perspective: Survey **submitted and answered by 57 sites**
 - Support for himem jobs:
 - Sites in general in favor, some doing it already for some VOs
 - Collaborative spirit in general but need for discussion, improvement of practices and accounting
 - Mem requests need to be explicitly specified: scheduling of the mix becomes crucial for effective use of resources
 - Himem slot usage must be properly accounted for:
 - Evolve the accounting metric to reflect use of CPU and mem
 - Including effect of CPU cores left idle (charge for them?)
 - Whole-node scheduling: Sites are open to tests but express concerns on:
 - Slot utilization efficiency:
 - draining
 - Resource management
 - less flexibility on the site side to allocate resources
 - insufficient turnover of resources in multi-VO sites?
 - 8 to 16 core slots: Some sites express need for further study
 - In general, similar concerns to whole node (flexibility in resource allocation, resource utilization efficiency)
 - But also concerns on higher fragmentation if need to also deal with single core, multiple core counts
 - Compatibility with some WN types (not multiple of 16) also architecture (tasks running over multiple NUMA processors?)

Discussion in DESY

Proposals for discussion

- **Allow some jobs to request memory over physical amount per core**
 - up to site configured max RSS
 - unlimited number of cores if mix means no idle cores, or capped at low(10%) level?
 - VO responsible for this. Monitored and enforced by the site - how? Trust but verify?
 - accounting for any idle cores important for pledges
 - important enough to develop memory dimension or rescaling in APEL?
 - HS23 does not scale linearly with HT cores. Accounting to reflect this.
- **Move to 16 core as new standard(currently 8) where a VO requests it,**
 - if then CMS would want 16 everywhere. ATLAS ok with mix, especially when HT-off.
 - major/all VOs send 16 core jobs as standard, to ease slot re-use
 - must allow 1 core jobs
 - ALICE want then 72hr walltime. CMS also prefer longer walltimes for more cores.
 - to do with draining inefficiency at the end

Discussion in DESY

Whole node scheduling

- Wholenode only where advantageous, e.g. GPU, numa pinning, many-core job(non-parallel)
 - oversubscription of cpu can improve efficiency when some jobs not cpu-bound
 - useful for packing gaps and draining, i.e. better slot efficiency
 - should be allowed also in MCORE if cgroup contained
 - otherwise let the Batch System schedule
- Pledged resources must be able to continue with S/MCORE jobs
 - combining whole nodes with S/MCORE jobs is problematic for many sites
 - If (CMS & ALICE & expert BS admin & volunteer) then do it
 - need at least 2 VOs with reliable whole node jobs, to keep slots
 - for this subset ATLAS would submit some wholenode(but also S/MCORE)

~~Relaunch Multicore TF to address this between exps, sites and WLCG~~

Job Allocation and Handling WG



Next steps

- Discussion on next steps
 - Dedicated sessions for each VO detailed view?
 - Report on tests from VOs with dedicated sites?
 - Etc?

- Next meetings of the WG
 - Fixed period, e.g. bi-weekly?
 - When new material is available?