A PROTOTYPE OF A VIRTUAL ANALYSIS FACILITY

Stefano Bagnasco, Dario Berzano, Stefano Lusso, Massimo Masera

Istituto Nazionale di Fisica Nucleare, Torino and Dip. di Fisica Sperimentale, University of Torino

Abstract - Current Grid deployments for LHC computing (namely the WLCG infrastructure) do not allow efficient parallel interactive processing of data. In order to allow physicists to interactively access subsets of data WLOS minister double with the interactively access subsets of data (c.g. for algorithm tuning and debugging before running over a full dataset) parallel Analysis Facilities based on FROOP have been deployed by the ALICE experiment at CERM and elsewhere. Whereas large Tier-1 centres may afford to build such facilities at the expense of their Grid farms, or exploit the large number of jobs finishing at any given time to quickly collect a number of nodes to temporarily allocate for interactive work, this is likely not to be true for smaller Tier-2s centres. Leveraging on the virtualisation of highly performant multi-nore machines, it is possible to build a fully virtual Analysis Facility on the same Worker Nodes that compose an existing LOG Grid Farm. Using the Xen paravirtualisation hypervisor, it is then possible to dynamically move resources from the bach instance to the interactive one when needed, minimizing latencies and wasted resources. We present the status of the prototype being developed, and some experience from the very first users.

Xen Dom()

INTERACTIVE ANALYSIS IN ÷۲ **H** •

• At Tier-1s

- Large Computer Centres with thousands of CPUs
- Feasible to take some out of the Grid infrastructure to build a PROOF-based Analysis Facility • Or may even be possible to "drain" jobs and switch to interactive mode quickly as soon as jobs finish and WNs
- become available

• At Tier-3s

- Very small number of CPUs
- · Probably not even a Grid site, at least with gLite middleware: use PROOF for parallel processing

And Tier-2s?

- Most resources are provided as Grid WNs
- In many experiments' computing models, including ALICE: this is where user analysis runs



• Xen can dynamically allocate resources (both CPU priority and memory) to either machine, no reboot or restart needed

• Normal operation: PROOF slaves are "dormant" (minimal memory allocation, very low CPU priority)

- When needed, resource can be moved to from virtual LCG WNs to virtual PROOF slaves with minimal latency
- Grid batch job on the WN ideally never completely stops, only slows down: non-CPU-intensive I/O operations can go on and do not timeout As the demand for interactive access increases, resources can be added either by shrinking further
- the WNs or by "waking up" more PROOF slaves As soon as everybody goes home, resources can be moved back to the WNs, that resume batch processing at full speed
- WN WN

Had Kill unrei FAM mod orego had, santer may 16/24

LCG CE

THE PROTOTYPE

Hardware

LCG C

WN

WN

- 4x HP ProLiant 360DL, dual quad-core, plus one head node for access, management and monitoring
- Separate physical 146GB SAS disk for each virtual machine performance isolation
- Private network with NAT to outside world (currently including storage)

Software

Linux CentOS 5.1 with kernel 2.6.18-53 on domO





MOVING RESOUR

- Moving resources from the Wn to the PROOF slave, as
 - seen by the WN About 1.5 minutes to complete transition (but the slave
- can be started almost immediately)
- Swap usage continues to increase afterwards
- CPU efficiency (CPU time/wall clock time) for regular ALICE MC production/reconstruction jobs, in three
- different resource configurations Jobs become increasingly I/O-bound as swap activity increases
- No abnormal job terminations observed (above what is seen in regular WNs)



•

• Main advantage: no big development needed, all tools are production-grade already

254 254

204 174 800

256 8 800 100

800

63 671

127964.4

• Prototype deployment did not require any development



→ Contacts: [bagnasco|berzano|masera|lusso]@to.infn.it

