TRACK SUMMARY: DISTRIBUTED PROCESSING AND ANALYSIS ON GRIDS AND CLOUDS

Oliver Gutsche Philippe Canal Johannes Elmsheuser

FNAL Ludwig-Maximilians-Universität München

25 May 2012/CHEP 2012, New York

- **2** LHC EXPERIMENT OVERVIEW AND STATISTICS
- **3** MultiCore and Job Scheduling
- **(1)** CLOUDS AND VIRTUALIZATION
- **5** FUTURE EXPERIMENTS
- **6** FUTURE SOFTWARE



- Merged track from previous CHEPs: Grid and Cloud Middleware and Distributed Processing and Analysis
- 174 abstracts after merging and reassignments to/from other tracks
- 31 talks in 7 parallel sessions 2 no-shows
- 143 posters accepted
- 27 papers already submitted to the journal
- \Rightarrow Largest Track very difficult to make everybody happy
- Broad variety of Grid and Cloud related topics



Topics of the different session:

- LHC experiment DDM and WM overviews
- Multi-Core and parallel scheduling
- Virtualization and Clouds
- Monitoring and Security
- Storage
- Grid middleware
- Dynamic data placement
- Experiment computing models

POSTER SESSIONS



Lively discussions during the Poster sessions 142 posters are close to impossible to mention during the summary, apologies for being so brief about the posters

J. Elmsheuser (LMU München)

2 LHC EXPERIMENT OVERVIEW AND STATISTICS

- **3** MultiCore and Job Scheduling
- CLOUDS AND VIRTUALIZATION
- **5** FUTURE EXPERIMENTS
- **6** FUTURE SOFTWARE

- All experiments have built their customized workload management systems for production and analysis and data management system on top of the existing grid middleware
 Very successful in delivering physics results
- But experiments are trying to streamline systems, remove unneccessary components, ease operations with limited person-power, find commonalities, scale to higher needs, adapt to new technologies

Workload Management Systems



J. Elmsheuser (LMU München)

Workload Management Systems



Workload Management Systems



New WM system

- · Consolidate analysis and organized activity
 - Same components but different instances
 - Prevent interference
- Single entry point for requests.
- Permanently recorded reproducible
- Requester can view status.
- · Prioritization between requests.
- Approval chain.
- Work distributed automatically & optimally to resources.
- · Reduced manpower needs.
- Adapt to new features / requirements:
- Pilot jobs.
- Multi-core processes.
- · Some use cases require all events to be processed.
- Cope with intermittent problems.

S. Wakefield, CMS

Imperial College London

25/05/2012 8 / 22

DATA MANAGEMENT SYSTEMS



DATA MANAGEMENT SYSTEMS



Large amounts of data replicated - but sometimes never touched - ATLAS pioneered popularity based automatic data replication with Panda/DQ2

DATA MANAGEMENT SYSTEMS



Similar large amounts of data replicated - now also moving to popularity system and replica deletion

EXPERIMENTS STATISTICS



Several 20k-100k concurrent running production and analysis jobs for each

Maeno

EXPERIMENT GRID IMPROVEMENTS, COMMONALITIES

The next \mathcal{DDM} version: \mathcal{R} ucio

Why a new major version ?

- · New high-level use cases and workflows
- · New technologies, paradigms and middleware
- · Difficult to extend the existing system with new concepts
- Old design (2006) with some conceptual limitations and heavy operational burden

High Level Roadmap

- 2011: Technical meetings with other LHC experiments, user surveys, collection of use cases
 - \Rightarrow Conceptual model document
- 2012: Parallel and incremental development track, incubator projects, preparatory steps
- 2013: Rucio in production





ATLAS is building new DDM system - ALICE tries out bit-torrent based SW distribution

EXPERIMENT GRID IMPROVEMENTS, COMMONALITIES



Study if ATLAS Panda is suitable for analysis in CMS - HammerCloud is used among 3 experiments for grid site validation

J. Elmsheuser (LMU München)

2 LHC EXPERIMENT OVERVIEW AND STATISTICS

3 MultiCore and Job Scheduling

4 CLOUDS AND VIRTUALIZATION

5 FUTURE EXPERIMENTS

6 FUTURE SOFTWARE

MultiCore and Job Scheduling



- Number of cores per node is increasing
- Physical memory per core is decreasing
- Different multi-core job scheduling options possible full node scheduling with slightly different resource usage

J. Elmsheuser (LMU München)

MultiCore and Job Scheduling



Alternative methods for parallel job scheduling on local cluster or online system

- **2** LHC EXPERIMENT OVERVIEW AND STATISTICS
- **3** MultiCore and Job Scheduling
- **4** CLOUDS AND VIRTUALIZATION
- **5** FUTURE EXPERIMENTS
- **6** FUTURE SOFTWARE

CLOUD COMPUTING IN ATLAS



EFFICIENCY, ELASTICITY

ATLAS Cloud Computing R&D

ATLAS Cloud Computing R&D is a young initiative

- Active participation, almost 10 persons working part time on various topics
- Goal: How we can integrate cloud resources with our current grid resources?

Data processing and workload management

- PanDA queues in the cloud.
 - Centrally managed, non-trivial deployment but scalable
 - Benefits ATLAS & sites, transparent to users
- Tier3 analysis clusters: instant cloud sites
 - Institute managed, low/medium complexity
- Personal analysis queue: one click, run my jobs User managed, low complexity (almost transparent).

Data storage

- · Short term data caching to accelerate above data processing use cases Transient data
- · Object storage and archival in the cloud
 - Integrate with DDM

Fernando H. Barreiro Megino (CERN IT-ES) CHEP- New York May 2012



Data Access Tests

- Evaluate the different storage abstraction implementations that cloud platforms pròvide
- Amazon EC2 provides at least three storage options
 - Simple Storage Service (S3)
 - Elastic Block Store (EBS)
 - Ephemeral store associated with a VM
 - · Different cost-performance benefits for each layout that need to be analyzed
- Cloud storage performance on 3-node PROOF farm
 - · EBS volume performs better than ephemeral disk
 - But ephemeral disk comes free with EC2 instances
 - Scaling of storage space and performance with the size of the analysis farm





Results

- · 100 nodes/200 CPUs at Cloud Sigma used for production tasks
- · Smooth running with very few failures
- Finished 6 x 1000-job MC tasks over ~2 weeks
- · We ran 1 identical task at CERN to get reference numbers

	HELIX	CERN 36 failed, 1000 succeeded 8136.6s ± 765.5s	
Success Rates	265 failed, 6000 succeeded		
Mean Running Times	$16267s \pm 7038s$		

 Wall clock performance cannot be compared directly, since we don't have the same hardware on both sites

CloudSigma has ~1.5Ghz of AMD Opteron 6174 per jobslot, CERN has a ~2.3GHz Xeon L5640

· Best comparison would be CHF/event, which is presently unknown







ATLAS with extensive Cloud R&D, tested production on commercial cloud

F. Barreiro Megino, ATLAS

FURTHER CLOUD EXAMPLES



FURTHER CLOUD EXAMPLES



Summary

- The HEPiX Virtualisation Working Group has shown how <u>trusted</u> user generated images can be safely instantiated at Grid sites
 - compatible with site security policies and obligations, and
 - with a guaranteed environment for the experiments.
- Three options now for following this up
 - Ignore and hope it goes away
 - Integrate trusted virtual images into a traditional batch environment
 - Exploit the proposal to deliver a "Grid of Clouds" that enables sites to easily and simply present to experiments a seamless—but dynamically changing—set of resources that the experiments can exploit to schedule work according to physics priorities.

Tony.Cass@<u>CERN</u>.ch



CERNVM co-pilot on opportunistic resources - Virtualization offers coherent environment

FURTHER CLOUD EXAMPLES



- Setup local Virtualization or Cloud cluster with ROCED
- WNoDeS Mixed Mode lets a resource center to progressively introduce virtualized services without disrupting existing setups and maximizing resource utilization

J. Elmsheuser (LMU München)

- **2** LHC EXPERIMENT OVERVIEW AND STATISTICS
- **3** MultiCore and Job Scheduling
- CLOUDS AND VIRTUALIZATION
- **5** FUTURE EXPERIMENTS
- **6** FUTURE SOFTWARE

Computing at Belle II, SuperB

Estimated Data Rates

Experiment	Event Size [kB]	Rate [Hz]	Rate [MB/s]		
High rate scenario for Belle II DAQ:					
Belle II	300	6,000	1,800		
LCG TDR (2005):					
ALICE (HI)	12,500	100	1,250		
ALICE (pp)	1,000	100	100		
ATLAS	1,600	200	320		
CMS	1,500	150	225		
LHCb	25	2,000	50		
MIT Thomas Kuhr	CHEP 22.05	T. Kuł	r, Belle		

SuperB distributed resources

- The distributed computing infrastructure, as of May 2012, includes several sites in Europe and North America
- · EGI and OSG Grid flavours have been enabled
- The LHC Computing Grid architecture was adopted to provide the minimum set of services and applications upon which the SuperB distributed model could be built.
- Computing resources needed in a typical year of SuperB data taking are of the same order as corresponding ATLAS and CMS estimation for 2011



New experiments adopting LHC grid software: DIRAC, AMGA, HappyFace, CVMFS

- **2** LHC EXPERIMENT OVERVIEW AND STATISTICS
- **3** MultiCore and Job Scheduling
- CLOUDS AND VIRTUALIZATION
- **5** FUTURE EXPERIMENTS
- **6** FUTURE SOFTWARE

Software and Middleware



Local and Remote HTTP data access

Software and Middleware



Computing Models have evolved from hierarchical system to mesh - FTS3 will adapt to these changes

Software and Middleware

News: dCache & pNFS



- NFS v4.1 / pNFS has been supported since 2009.
- Deployed **in production** (at DESY) for over a year.
- Fermilab's REX dept. evaluated dCache NFSv4.1 for their Intensity Frontier experiments:

"**Results look promising**, throughput scales well with number of pool nodes"

- Supports:
 - authn: trusted-host and Kerberos
 - all three GSS security modes.

dCache, agile adoption of storage technology | Paul M D. Milla



Summary

- Storage systems based on traditional spinning disks have a limited I/O performance.
- Adding SSDs, as an additional cache layer in currently deployed storage solutions (DPM, dCache..) can boost I/O performance significantly.
- An analysis of CMS user activities has shown that only 10% of cache, results in 70% of random reads from the cache.
- A storage system with 3 Tiers is transparent for the user and cost effective for sites.

	5/22/12	CHEP 2012, The need of 3 Tier storage	17
r			D.Ozerov

News in several areas from large scale storage system dCache - 3 tiered system with SSD cache allow better performance

Projects



Both EMI and OSG are looking into the future, with EMI planning to hand over software projects to software teams and OSG looking to continue support and integrate more non-HEP users

- Many thanks to the Organizers, Speakers and Poster Presenters for their interesting contributions
- We think the sessions were a big success and well received
- In general the session were on time although we had lively discussions
- Proceedings:
 - Due to the large number of possible proceeding contributions we are asking all track participant to help
 - Please volunteer to review 1-2 papers
 - Don't worry: we would not expect detailed technical comments about the content

Thanks a lot for an interesting and fun time in NYC