

Table of Contents

Summary of GDB meeting, September 9, 2015 (CERN)	1
Agenda	2
Introduction	3
GDB Future - I. Bird	4
MW Readiness - A. Manzi	5
Information System - M. Alandes Pradilla	6
httpd Deployment - O. Keeble	7
Volunteer Computing - N. Hoimyr	9
Machine/Job Features - S. Roiser	11
CPU Benchmarking	12
LHCb - P. Charpentier.....	12
CMS - D. Abdurachmanov.....	12
ATLAS - A. Filipic.....	13
ALICE - C. Grigoras.....	13
Discussion.....	14
Conclusion.....	15
Offline summary by HEPiX experts.....	15

Summary of GDB meeting, September 9, 2015 (CERN)

Agenda

<https://indico.cern.ch/event/319751/>

Introduction

October GDB in BNL as part of HEPiX (Wednesday 14)

- Register
- Submit abstracts

November GDB: known clash with EGI Conference but confirmed

Next pre-GDBs

- DPM workshop : December 7-8 (pre-GDB)
- Discussions progressing on a workshop about cloud storage around the end of the year
 - ◆ Focus will be running data-intensive workloads in a cloud
- IPv6 F2F meeting: no date yet

ARGUS

- Last collaboration meeting last week: see <http://indico.cern.ch/event/442356/>
- 2 more people at CNAF to work partly on ARGUS (INDIGO)
- EL7/Java8: good progress
 - ◆ All packages available
 - ◆ PAP already updated with last dependencies, PDP ent PEPd being worked on
- pepd issue at CERN: finally found this summer one situation triggering it, should help to fix it quickly

WLCG workshop, Lisbon, February 1-3

- <http://indico.cern.ch/event/433164/>

Actions in progress

- Multicore accounting now in good shape
- wiki with "class2 services" by VO still to be completed
- wiki page Storage for helping diminishing the numebr of protocols: still to be completed

Forthcoming meetings

- HEPiX, BNL, October 12-16
 - ◆ <https://indico.cern.ch/event/384358/>
- EGI Community Forum, Bari, November 10-13
- SuperComputing, Austin, November 15-20
- WLCG Workshop, Lisbon, Feb. 1-3
- HTCondor European Workshop, Barcelona, Feb. 29-March 4
 - ◆ Including ARC CE workshop

GDB Future - I. Bird

MB : supervises the work of the collaboration

- Mandate defined by the MoU
- Well-defined membership, chaired by WLCG project leader

WLCG workshop : present and inform collaboration of strategic directions, opportunity for broad discussions

- Must be seen as collaboration meetings
- Also WLCG partners like Belle2
- Organised by OpsCoord and MB

WLCG Technical Forum (WTF): new board to be proposed at MB next week

- Mandate: Prepare the long-term future of the WLCG infrastructure, complementary to,

and supporting, the evolution of the computing models of the experiments. The timescale of Run 3 and Run 4 are the focus, although improvement should be ongoing.

- Members: experiment nominees (1+1), site representatives (1 per T1, 2 or for T2s) + technical providers as necessary
- Uses GDB for broader discussions
- Chair nominated by WLCG project leader, 2 year term

GDB: mandate as defined by the MoU still valid except the for following the operations

- Operations delegated to OpsCoord
- Attendees: should be open to all WLCG members + partners as approved by MB
- Technical Forum should propose topics in addition to GDB members and MB
- **Time to call for nomination now: election proposed in November**
 - ◆ Ian will send a call for nomination to the GDB list
 - ◆ Send nominations to Ian: self-nomination allowed, Ian will approach the nominees
 - ◆ Eligible: any WLCG site representative
 - ◆ Voters: 1 per country member of the MoU

MW Readiness - A. Manzi

Increased coverage of products: now DPM, CREAM, dCache, StoRM, EOS, ARC CE, VOMS clients

Products not yet tested: WN, Xrootd, ARGUS, FTS3

- ARGUS: agreement to wait for the new version before integrating it

VO participation

- ATLAS and CMS : specific workflow made for this effort
- LHCb: generic document for certification but no active participation since then
 - ◆ S. Roiser: still willing to participate but limited manpower...
- ALICE: not yet involved

Several sites involved but more welcome!

- 3 new sites joined in the last months
- Next meeting: Sept 16, 4 pm CEST

In the last 6 months, several problems identified during MW readiness testing

Plan to start CentOS 7 readiness testing as soon as PTs release their first compliant version

- A few sites available to test it
- Discussing with experiments the priorities

Collaboration with EGI working well

- Participating in URT
 - ◆ Good collaboration around UMD4
- Several WLCG sites reporting their results also to EGI

MW Package Reporter

- Site information collected in a DB
- An application is mining this information and filing it into the SSB: accessible to MW officer and site admins
 - ◆ Also allow to define the baseline for each product
 - ◆ Currently being added: ability for a site admin to see products used at his site, based on packages (similar to Pakiti dashboard)
- Using Pakiti v3 for the reporting: support "tags" to flag in the DB that the site is participating to MWR

Based on testing activities, will decide wider deployment in WLCG

Information System - M. Alandes Pradillo

BDII : no new version during the last year

- Working on EL7 version: no code change
 - ◆ Done by EGI
- Main issue currently: slapd crashing, seems related to a problem with 2.4.40-5
 - ◆ Sites advised not to upgrade, ticket opened at RH
- Most sites are running the last version (5.2.23) or the one before
- Increased number of GLUE2 endpoints

GLUE2 validation proved useful: very low number of error/warning since it is part of SAM tests

- Still not validating the storage capacity: to be discussed with EGI

OSG announced their plan to stop BDII support in the future

- Partly related to well known limitations in our IS architecture: need to move forward
- Need to define the authoritative source of information for all WLCG infrastructures and implement a clear distinction between static, dynamic and mutable information
- An IS TF meeting planned on Sept. 24: OSG, NDGF and EGI will present their plans
- A first document presenting use cases is being prepared: will be presented at MB next week

Discussion

- I. Bird : this document should summarise what the experiments need. It is important that we get driven by what is needed by the LHC experiments and not by OSG, EGI and NDGF specific needs (that can be a subset or a superset of what LCG needs)
 - ◆ Maria: this is what the document does. But it is also important to hear their plans to ensure a smooth transition and a good coordination where it is needed.
 - ◆ Oxana: this document will also be useful for sites as it is often difficult for them to understand where the experiments is getting some information from. This is documented and it will be easier for sites to fix the problems.
- Michel: do we really believe that we can develop a new IS? We have a manpower problem as exhibited by the difficulty to progress on projects we had in the last 2 years (GSR, AGIS)
 - ◆ Markus : we have more and more dispersed source of information, coping with them in all experiments is a waste of resources, would be better to promote a common way of aggregating it.

httpd Deployment - O. Keeble

Reason for creating the TF

- All pieces in place now for use of http by HEP
- ATLAS and LHCb have started to develop http-based infrastructure
- All storage implementations supporting http protocol
- ROOT support now as performant as xroot, thanks to Davix
- Cloud storage based on S3 which is an http-based protocol

Mandate

- Define the minimum set of functionality
 - ◆ Also document further functionalities desirable and their use cases
 - ◆ For example space monitoring is important for storage offering only http access
- Deploy monitoring/validation tools
- Summarize deployment advices for sites, including baseline versions
- Track and support deployment until it can be integrated into the standard experiment operations

Membership

- ATLAS, CMS, LHCb
- Several sites
- All storage implementations (including EOS and xrootd)
- WLCG monitoring

3 meetings since inception

- Functionality: described in a Twiki page, HTTPTFStorageRecommendations
 - ◆ Covering all operations required
 - ◆ includes 3rd party transfers, possibly with S3 endpoints
- Access monitoring, as for gridftp and xrootd: 2 approaches discussed
 - ◆ Xrootd f-stream like UDP summary packets
 - ◆ Summary information according to a JSON schema
- Monitoring/Validation: a shared SAM probe almost ready, working on integration into the new SAM/Nagios scheduler and discussing visualisation into SAM3

Discussion

- S. Luders: http or https
 - ◆ Oliver: this the all http protocol group (http, https, WebDav, S3...). Authentication is always https but data transfer is generally http for performance reason (using a session token to assert the user has been authenticated/authorized).. * Maria D.: now that HEP experiments are using http, should we try to establish closer contact with W3C to discuss the extensions we needed/made?
 - ◆ Oliver: HEP experiments are not really using http, they are only capable of using it. But contacts with W3C would certainly be good and having the extension we made part of the standard would certainly be a good thing.
- Michel : plans in experiments, in particular ATLAS and LHCb, to use http for data access in addition to metadata operations ? * Stefan: LHCb, with the help of Fabrizio, started a dynamic federation of all LHCb endpoints with http enabled to check data consistency against the catalogs. We want to explore other possible usages but we have no concrete plans yet.
 - ◆ P. Charpentier: for data access, we need to demonstrate a clear benefit of using http rather than xroot

- ◆ Alessandro D.G.: same situation for ATLAS. http currently used for metadata operations but doesn't exclude other usages in the future. For example, currently discussing using http for accessing log files.
- ◆ Oliver: pressure to move to http for data access may come from sites rather than from experiments. Sites supporting non HEP communities may tell that they want to consolidate all accesses over one standard protocol.

Volunteer Computing - N. Hoimyr

Volunteer computing can provide substantial resources and outreach benefits

- Cost of the resources is the effort to attract and assist volunteers
- Virtualization support in BOINC allows an easy integration of HEP applications relying on CERNVM and CVMFS
- Share commonalities with the VAC model: possible to share some underlying services

LHC@Home

- SixTracks (beam simulations): the initial BOINC project, revived for HL-LHC studies
 - ◆ 12K volunteers, 19K hosts
 - ◆ Peak of 150k simultaneous jobs
- Test4Theory (theoretical fitting of past exp data): since 2011
- ATLAS@Home: in production since summer 2014
 - ◆ Interface to Panda through an ARC CE
 - ◆ ~6K jobs continuously running: 2nd largest simulation site for ATLAS
- Beauty@Home (LHCb), since 2012
 - ◆ Required x509 credentials in the VM: soon a new proxy layer will be available
- CMS: prototype since 2014, progressing well, based on data bridge

Service infrastructure at CERN: Drupal portal (lhathome.cern.ch), BOINC instances in OpenStack

- vLHC@Home: instance to host projects related to LHC relying on VMs and connecting to an experiment framework
- Other projects are hosted on separate instances
- Databridge solution available from IT/SDC

why volunteer computing: free and substantial resources plus outreach channel

BOINC suitable for hundreds of applications

volunteer is similar to cloud and the vacuum approach (CernVM) can use common pieces

LHC@Home examples : - LHC SixTracks, since 10 years, high CPU application (Fortran) running on Linux, Mac and Windows renewed activity for High Luminosity LHC - Test4Theory : launched in 2011, pioneered virtualisation more than 2 trillion events since 2011

- ATLAS@Home: pilot beginning in 2014, open to public since 1 year CernVM ARC-CE currently more than 6000 jobs currently (in the top 5 ATLAS T2 for simu) - Beauty@Home: pilot since 2012 required x509 credential -> only LHCb members develop new proxy to be ready soon - CMS@HOME: since summer 2014, prototype service

Boinc at CERN - service layers : drupal portal + openstack VM from the boinc server

- MySQL DB + NFS back-end for storage
- vLHC and other projects on separate instances to avoid I/O bottlenecks

Would be good to have another dedicated meeting (pre-GDB?) in the winter when CMS and LHCb are in prod, like we did one year ago.

Discussion

- Markus: last year there were concerns about BOINC funding/future
 - ◆ Nils: BOINC missed one round of NSF funding and there were some worries then. But turned out to be no as dramatic: BOINC was turned into a community project, core developers from Berkeley University are still contributing, project active. Should not worry to much.
- Michel: last year the data bridge was presented. What is its current production usage today?
 - ◆ Nils: used in production by theory apps, LHCb is planning to use it. ATLAS wanted to test it but in fact the ARC CE solution is working well for them so no real activity. Need a wider usage to evaluate exact scalability.

Machine/Job Features - S. Roiser

MJF: a way for resource providers to provide information on WN/job slot to users

- Works with both standard WNs (file-based) and virtual machines (Apache)
- Fine granularity compared to BDII
- 2 sets of information : machine features and job features
- Possible usage: discover limits on WN, compute time left for a job, announce shutdown of a WN

Reference implementation for all batch systems (HTCondor, LSF, SGE, SLURM, Torque/PBS) + a pervasive one based on Apache

- Apache implementation originally for VMs but could be used for batch systems too. Allow to publish information when there is no way to do it on a local file system (ex: commercial clouds). Can be hosted separately from the cloud/resource.
- Actual deployment at site may require some adaptation
- Querying all the implementation can be done the same way
- Probe to monitor that the required information is published with meaningful values
 - ◆ No checking of correctness of the values currently

Deployment on the infrastructure not progressing quickly...

- According to the probe, only 4 grid sites having it deployed and running: CERN, GridKA, Imperial College, GRIF/LPNHE
- Also several cloud sites
- LHCb would like to have MJF deployed at each site before the end of the year

Assessing/validating values published is outside the mandate of the TF but is important for experiments to be confident they can rely on it.

- Publishing CPU power is strongly related to the benchmark discussion

Discussion

- Tim: on a public cloud, how does it work ?
 - ◆ Stefan: this would be through the Apache-based implementation, that could be hosted at CERN for example
 - ◆ Tim Bell: but if the VM is moved, the CPU power will change
 - ◆ Stefan: clearly a use case for a fast benchmark run on the fly
- Andrea S.: what is the interest in experiments other than LHCb?
 - ◆ Stefan: this project was started as a WLCG effort, with all the VOs represented
 - ◆ Latchezar: ALICE has nothing against MJF deployment but has concerns about the correctness of the values published.
- A. Sciaba: is the doc for a deployment ready ?
 - ◆ Stefan: doc is ready and deployment is easy. LPNHE, just did it and it pop up
- D. Abdurachmanov (CMS): can we add new features? More details on the machine could help to build a global view of the grid.
 - ◆ Stefan: this can be flexibly extended. Currently we focussed on basic features. Need to agree as it makes sense only if this is published by all sites.
 - ◆ Michel: I'm not sure this mechanism is suited for what you seem to have in mind. There is no way to access this information outside the WN, so it cannot be used as a source of information to build a global view. It's really intended to provide fine granularity information to a payload running on a particular machine. For a global view, use the information system but it cannot provide the same granularity.

CPU Benchmarking

LHCb - P. Charpentier

2 reasons for benchmarking

- Job matching / masonry : again 2 different use cases
 - ◆ Can a pilot match a job of a certain known duration with the resources left?
 - ◆ Adjust job duration to the resources left (for example, the number of simulated events)
 - ◆ Precision: CPUWorkLeft should not be overestimated, for masonry as precise as possible (worst 20%)
- Accounting

Main features required

- Time left: expect it to be defined as a duration (normalize what is provided by a batch system)
- CPU power available to the job (CPUPower)
 - ◆ Can be obtained either by running a fast benchmark or through a WN lookup table but this option proved to be difficult to implement (see CERN experience)
 - ◆ Power per allocated core : must be based on the number of jobs slots rather than the number of logical cores when hyperthreading is enabled
- $\text{CPUWorkLeft} = \text{CPUTimeLeft} * \text{CPUPower}$

Validation of values published in MJF

- With a real job
 - ◆ The job is a MC production job (CPU bound) with a large number of events (100+)
 - ◆ job/MJF ratio: matches pretty well, both at CERN and GridKA (after taking into account properly the slot/phys. core ratio at each site)
 - ◆ Still some subtle effects giving wrong results on some machines : same HW model may give different results (memory speed? turbo boost?)
 - ◆ AMD better
- LHCb fast benchmark / MJF : rather good matching, around 1 but still some corner cases. Anyway sigma better than 10%

Conclusions

- MJF important for LHCb
- Important to agree on how the formula to compute the power per slot from the machine power
- Need to agree on conditions as the one used by jobs: in particular the same number of benchmarks as the number of job slots rather than the number of logical cores
- Fast benchmark is a compromise but less precise than MJF (that can publish the actual HS06 rating on a machine)

CMS - D. Abdurachmanov

Too many different workloads: cannot expect one benchmark to allow to predict the power for all workloads

CMS using several benchmarks: CMSSW simulation, Geant4 simulation

- Recently did a wide data collection campaign: difficult to correlate to actual power of a machine
 - ◆ May want to look at the LHCb methodology to see if some useful correlation can be extracted from the data

- Also started to look at ARM architecture

ATLAS - A. Filipic

Reasons for benchmarking similar to LHCb (and other experiments)

- A precision of 10% is good enough
- New job infrastructure (JEDI) is using site HS06 for brokering
 - ◆ Site HS06 in AGIS, based on BDII information (site average). This is a starting point, aware that this is not representing any real machine at the site...

Atlas benchmark infrastructure is in CVMFS

- HTCondor benchmarks: promising a few % accuracy but not consistent with current results, need to be checked
- Whetstone+ Dhrystone: used by BOINC
 - ◆ Whetstone quite reliable at a few % level
 - ◆ Runs in ~10s
- KitValidation (muon MC simulation): a few minutes, pretty accurate, good candidate for high precision benchmark but too long to be run as part of each job
- HS06 unusable directly by ATLAS

Would support WLCG defining one common fast benchmark shared by all experiments

- Target : ~10s for a fast benchmark run as part of each job
- Accuracy target: 1%
 - ◆ Many comments that this is totally unrealistic... and probably excessive when currently relying on a site average value...

ALICE - C. Grigoras

Rely on 2 fast benchmark : ROOT /test/stress (~30s), condor_kflops

- Every benchmark run twice after the simulation job: only the second execution is recorded
- ROOT benchmarks well correlated with real simulation jobs, not the case for kflops
 - ◆ kflops issues may be related to a different integer/float ratio
 - ◆ Is ROOT benchmark scaling with workload of other VOs?

Fast benchmark results don't scale with HS06 results found

- During discussion, it appeared that the HS06 value was picked up based on standard results found for the matching CPUModel rather than a value actually computed for the specific node: clearly not a reliable method...

Planning to build a per node database of benchmark results

- Currently based on CPU model

HS06 considered useless, even for accounting, as it is a unit not correlated to the real workload performance.

Discussion

Latchezar to Philippe: do you have correlation between nb of events and your benchmark value ?

- Philippe : yes, good correlation
- Latchezar: why are you interested in MJF ?
- Philippe: if MJF can provide more precise values than a fast benchmark, we'll use it
- Charpentier : if MJF would be more precise we would use it. Usage encouraged by a few sites we talked to.

Latchezar: HS06 is not working any more in Intel, it is probably still working on AMD: on Intel machine we have lost the correlation between number of events processed and HS06. We have 4 different fast benchmarks, we could consolidate to 1 instead of having 4 of them.

- Philippe: not true for LHCb according to the study presented
- Michel: thanks for your detailed study, would have been good to see the same kind of numbers/plots presented by other experiments
- Dirk: want to mention that CERN did quite an extensive analysis of both ATLAS and CMS jobs in term of event processing rate and found a good correlation between them. In fact, were able to predict an ATLAS job duration based on a CMS benchmark with an error less than 20%.

Philippe: hyperthreading seems responsible for the spread of the benchmark results for a throughput gain around 20%. If this results in MJF value being wrong and leading to 5% of jobs crashing, is it really worth activating it?

Ale: we have seen that all 4 VO have 2 cases in mind: accounting and job duration estimate for brokering. Both can safely be separated. Accounting : a posteriori; Job length: must be fast and +/- 20% would be fine.

- Costin: you are doing it because you do not trust HS06
- Ale: no, but because the actual slot power is not available currently: MJF may help.
- Philippe: why insist for a few seconds rather than a minute: one benchmark per pilot.
- Ale: in ATLAS, only one job executed per pilot...

Philippe: the values in the EGI accounting are random... and don't correlate with experiment accounting.

- Ian: this is not true! Nobody is claiming it is 1% accurate but each time we have compared to the VOs, it was quite all right. We clearly rely on somebody putting a value (the CPU power) by hand... Room for improvements. But not related to the benchmark accuracy.
- Michel: emphasize the need for clearly defining how to run the benchmark and/or compute the value. Having MJF + fast benchmark may help to identify problems by studying correlations (see Philippe's presentation).

Ian to all experiments: what about scaling between different types of workload? can you scale from MC to reco ?

- All VOs answering yes: can predict with a reasonable accuracy the duration of any workload based on MC (CPU intensive) and the average CPU efficiency of a given workload (known).

Accounting: long discussion on whether a fast benchmark should also be used for accounting.

- Latchezar pushes to explore this direction
- Ian: insists that for accounting we need a solution common to all experiments, usable by vendors (benchmark is used for procurements) and understandable by funding agencies. Must also be easy to adapt to new architectures as new ones are coming every year...

Eric Lancon : ATLAS needs to do its own accounting because EGI accounting was not available quickly enough to follow the everyday operation. This was the first reason for the fast benchmark but 2 use cases have been added over time: job length prediction and computation of CPU power at (commercial) clouds and resources outside EGI/OSG.

Conclusion

I. Bird: experiments invited to sit together to:

- Discuss how to improve the accounting accuracy with the benchmark we have today
- Discuss the feasibility and work on a common fast benchmark to address the use cases mentioned during the discussion, independently of the work done in HEPiX about next generation HS

Next step: discussion of the proposal at MB next week

Offline summary by HEPiX experts

"Note: This is the content of an email received after the discussion and published with the authorization of the authors (Manfred/Michele). Note that these persons left before the final discussion where Ian made his proposal to move forward."

Scaling of HS06 with HEP code

Philippe has presented the results of very detailed investigations which are showing good scaling (in average) of LHCb jobs with HS06. Although GridKa has enabled all these "nasty features" like HT and Turbo Boost, Philippe has found out perfect scaling also at GridKa.

Alice has reported about bad scaling of Alice jobs "with the few HS06 results they found".

Atlas' concern is that they cannot run HS06 inside of jobs to estimate the performance of the provided job slot.

Conclusion: There are currently still no strong arguments that HS06 doesn't reflect the performance of the typical workloads of LHC VOs. So there is absolutely no pressure to migrate to another benchmark.

Performance of the provided job slot

Users want to know about the performance of the provided job slot (on batch farms as well as in cloud environments). There are 2 possible solutions:

- Using MJF to get the performance details (HS06, number of slots, ...).
- Using fast benchmark tools to estimate slot performance.

Conclusion: This is indeed an important topic. We must tackle it with high priority.

How to proceed (next months)

- We are ready to assist also Alice, Atlas, and CMS to monitor scaling of their jobs with HS06 using the already existing MJF implementations.
- We have already started some investigation of the LHCb fast benchmark and found quite good scaling with HS06 (in average, less precise than long running benchmark). Other VOs have mentioned other fast tools like Atlas KitValidator, Drystone+Wetstone, and so on. We should have a look at these tools also.

-- MichelJouvin - 2015-09-11

This topic: LCG > GDBMeetingNotes20150909

Topic revision: r2 - 2015-09-14 - MichelJouvin



Copyright &© 2008-2020 by the contributing authors. All material on this collaboration platform is the property of the contributing authors.

Ideas, requests, problems regarding TWiki? Send feedback