Uso di HPC negli Esperimenti LHC

Daniele Spiga INFN-PG Tommaso Boccali INFN-PI On behalf of

Stefano Piano, Francesco Noferini, Concezio Bozzi, Leonardo Carminati, Alessandra Doria, Alessandro De Salvo, Lorenzo Rinaldi

Outline

Why we talk about HPCs integration @HEP

What are the Integration challenges

LHC Experiment: status and ongoing activity

- A quick look at national landscape

Summary

The HEP (and close friends) Computing circa 2020

- Since ~1980, the HEP world has been facing a steady increase in the computing needs, at least from LEP times
- The increase in needs has seeded most of INFN Computing R&D and operations in the last 20 years
 - The GRID
 - The Tier-1 and Tier-2 centers
 - \circ $\,$ Now the Cloud, the Datalake, ...

	ALEPH	CDF	CMS	
	1995	2004	2007	
Dimensione dei	1 TB = 1000 GB	1 PB = 1000 TB	~10 PB	
dati raccolti		x1000	x10	
Capacita' di	<<100k	1.4 M	>25 M	
calcolo (SI2k)		x50	x20	

Nota: 1 PC attuale ~ 1 SI2k

2021: WLCG

- 161 official sites, in 42 countries
 - "We only miss Antarctica"
- Pledged resources
 - ~ 8 MHS06 (~800kCores)
 - **~700 PB disk**
 - ~1 EB tape
- Other resources (HLT farms, overpledges) count for at least another 50% on CPUs
- Transfers (as seen by FTS)
 - ~ 50 GB/s



What is in front of us?

- LHCb + ALICE upgrade already happened, start data taking in ~ 10 months
 - Ambitious, but manageable in semi-adiabatic mode
- ATLAS and CMS ~ 2028 with Phase-2
 - Currently still unable to match a flat funding profile



Many solutions / ideas under test; using HPC resources is one of them!

The HPC ("Supercomputer") panorama

- The world is literally full of Supercomputers. Why?
 - Real scientific use cases
 - Lattice QCD, Meteo, ...
 - Industrial showcase
 - And hence not 100% utilized, opportunities for smart users. Can we be one of them?
- They are usually very closed systems
 - Few users, with systems designed for their benefit
 - Highly parallel: the cost of networking is a big part of the total
- They are expensive
 - 100M-1B\$ each is the typical range (for comparison, the total "cost" of the WLCG resources is ~ 100 M\$)
- Some hint of global slowing down, but not for top systems where the "war" is on

TOP #1 SYSTEMS

In the last 20 years, the following systems made it to the top of the TOP500 lists:





1994 1996 1998 2000 2002 2004 2006 2008 2010 2012 2014 2016

Supercomputer - The expected future

- The race will go on, at least between major players
- EU wants to enter the game never at the top in the last 25y
- Next big thing is **ExaScale** (10¹⁸ Flops operations per second)
 - Should be well available by HL-LHC times; indeed in principle by this year!
- Somehow difficult to compare, technologies / benchmarks, but
 - LHC needs today the equivalent of ~30 PFlops
 - Not all the flops are equal! This is CPU currently
 - A single Exascale system is ok to process 30 "today" LHC
 - Scaling: a single Exascale system could process the whole HL-LHC with no R&D or model change
- Some FAs/countries are explicitly requesting HEP to use the HPC infrastructure as ~ only funding; it is generally ok IF we are allowed to be part in the planning (to make sure they are usable for us)



Top machines now

- 7.3 MCores, with ARM architectures
- 5 PB of RAM
- RedHat Linux
- 29 MW (~60 M\$/y ??)

The second system is

- 2.4 M computing cores Power9 + NVidia V100
- 2.8 PB ram
- RedHat Linux 7
- 10 MW power use (which alone is ~ 20 M\$/y)

It is difficult to translate this power in the HS06 metric we use, but:

- Intel Xeon Xeon E5-2640v3 is rated at 350 GFlops
- The same CPU is rated at 230 HS06
- \rightarrow 1 HS06 ~ 1.5 GFlops
- Fugaku ~ 250 MHS06 (which is ~ 40-60x WLCG)





A quick look at HPCs main features

HPC design features:

- Performant node-to-node interconnection
- Strict user access policy (used == ID card)
- Scarce / absent local scratch disk
- Limited capability for data access outside the facility, if not absent
- Relying on accelerator to boost performance (today mostly Nvidia GPUs, tomorrow FPGA, Intel Xe, AMD, ... ??)
- Storage systems are optimized for latency and speed, and not for total size.

Moreover

- HPC centers differ on a variety of specialized hardware setups and policy wise strict usage rules can apply mostly for security reasons

All experiments are working on integration of HPC resources with varying levels of technical difficulties because the HEP systems are built using different technical needs in mind:

 HEP workflows are typically data intensive, and systems deploy large storage systems close to the computing, balanced with CPU deployment. Most of the software is still designed and optimized for the x86_64 architecture

Let's see a "good" scenario [an example]

A typical Marconi A2 node was configured with	A typical WLCG node has
A KNL CPU: 68 or 272(HT) cores, x86_64, rated at ~¼ the HS06 of a typical Xeon	1/2 Xeon-level x86_64 CPUs: typically 32-64 cores, O(10 HS06/thread) with HT on
96 GB RAM, with ~10 to be reserved for the OS: 1.3-0.3 GB/thread	2GB/thread, even if setups with 3 or 4 are more and more typical
No external connectivity	Full outgoing external connectivity, with sw accessed via CVMFS mounts
No local disk (large scratch areas via GPFS/Omnipath)	O(20 GB/thread) local scratch space
Access to batch nodes via SLURM; Only Whole nodes can be provisioned, with 24 h lease time	Access via a CE. Single thread and 8 thread slots are the most typical; 48+ hours lease time
Access granted to individuals	Access via pilots and late binding; VOMS AAI for end-user access

A detailed analysis by category

(a longer list..)

https://cds.cern.ch/record/2707936 /files/NOTE2020_002.pdf

Explanation	CMS standard solution	CMS preferred solution for HPC	CMS fallback workable solution (full utilizability)	CMS fallback solution (for a fraction of workflows)	CMS no-go scenario	Possible CMS devels to solve the no-go
Base system architecture	x86_64	x86_64	x86_64 + accelerators (with partial utilization)		Currently, OpenPower, ARM, they could be used but at the price of physics validation	QEMU? Recompiling + physics validation?
Memory available to each thread / process	2 GB/Thread	2 GB/Thread	Down to 0.5 GB/thread needs heavy multithreading, at the expenses of CPU efficiency	GEN and SIM workflows need less than 2 GB/Thread (0.5GB/Thread would be a limit) in order to run efficiently	Less than 0.5 GB/thread	
I/O demand per process	5 MB/s/core	5MB/s/core		GEN and SIM workflows are mostly CPU bound still ok with 0.1 MB/s/core	Less than 0.1 MB/s/core	
Local space per production job	20 GB/Thread	20 GB/thread local	Less than 20 GB/thread ok if a shared high performance FS is available on all the machines Large multithreading lowers 20 GB/thread requirement to ~ 10	Some CMS workflows run for hours without creating huge local disk areas (GEN, SIM)	No sizeable local space and no shared usable FS	
Needed on WNs in order to access remote data, conditions, and to speak to the CMS Global Pool	Full outgoing connectivity	Full outgoing connectivity	Connectivity to only a subset of the IP ranges (for example, to CERN, and to a close xrootd proxy cache) And to everywhere we have condor services?	NAT with a very limited bandwidth via an edge service	No outgoing connectivity from the compute nodes and no NAT available	Edge service running Harvester or HTCondor? Prepare a single container to be deployed at the edge and doing: NAT for Condor, Squid, Xroot proxy cache,?
	Explanation Explanation Base system architecture Memory available to each thread / process I/O demand per process Local space per production job Needed on WNs in order to access remote data, conditions, and to speak to the CMS Global Pool	ExplanationCMS standard solutionBase system architecturex86_64Memory available to each thread / process2 GB/ThreadI/O demand per process5 MB/s/coreI/O demand per process20 GB/ThreadLocal space per production job20 GB/ThreadNeeded on WNs in order to access remote data, conditions, and to speak to the CMS Global PoolFull outgoing connectivity	ExplanationCMS standard solutionCMS preferred solution for HPCBase system architecturex86_64x86_64Memory available to each thread / process2 GB/Thread2 GB/ThreadI/O demand per process5 MB/s/core5 MB/s/coreLocal space per production job20 GB/Thread20 GB/thread localNeeded on WNs in order to access remote data_conditions, and to speak to the CMS Global PoolFull outgoing connectivityFull outgoing connectivity	ExplanationCMS standard solutionCMS preferred solution for HPCCMS fallback workable solution (full utilizability)Base system architecturex86_64x86_64x86_64 + accelerators (with partial utilization)Memory available to each thread / process2 GB/Thread2 GB/ThreadDown to 0.5 GB/ThreadI/O demand per process5 MB/s/coreDown to 0.5 GB/ThreadDown to 0.5 GB/ThreadLocal space per production job5 MB/s/coreLess than 20 GB/Thread ok if a shared high performance FS is available to - 10Needed on WNs in order to access remote data, conditions, and to speak to the CMS GlobalFull outgoing connectivityConnectivity full outgoing connectivityConnectivity onnectivityNeeded on WNs in order to access remote data, conditions, and to speak to PoolFull outgoing connectivityConnectivity onnectivityConnectivity onnectivity	ExplanationCMS standard solutionCMS preferred solution for HPCCMS fallback workable solution (full utilizability)CMS fallback solution (for a fraction of workflows)Base system architecturex86_64x86_64x86_64 + accelerators (with partial utilization)CMS fill outpartielyMemory available to each thread / process2 GB/Thread2 GB/ThreadDown to 0.5 GB/Thread needs heavy multithreading, at the expenses of CPU efficiencyGEN and SIM workflows need (5.68/Thread needs heavy multithreading, at the optimization of cPU efficiencyGEN and SIM workflows need (9.658/Thread uotiflows need (9.658/Thread) (9.658/Thread) workflows are mostly CPU bound still ok with 0.1 MB/s/coreI/O demand per process5 MB/s/coreSMB/s/coreGEN and SIM workflows are mostly CPU bound still ok with 0.1 MB/s/coreLocal space per production job20 GB/Thread localLess than 20 GB/thread ok if a shared high performance FS is available on all the machines Large multithreading lowers 20 GB/thread requirement to ~ 10Some CMS workflows run for hours without creating harea (GEN, sith)Needed on WNs in order to and to speak to the CMS GlobalFull outgoing connectivityConnectivity coring a subset of the IP ranges (for example, oces reword and to speak to the CMS GlobalFull outgoing connectivityConnectivity coring a subset of the IP ranges (for example, oces reword and to speak to the AS GlobalFull outgoing connectivityNAT with a a cose xeroid orx	ExplanationCMS standard solutionCMS preferred solution for HPCCMS fallback workable solution (full utilizability)CMS fallback solution (for a fraction of workflows)CMS no-go scenarioBase system architecturex86_64 scenariox86_64 scenarioX86_64 accelerators (with partial utilization)Currently, OpenPower, APM, out they out be used but at the price of physics validationMemory available to each thread / process2 GB/ThreadDown to 0.5 GR/Ihread needs heavy multithreading, at the expenses of CPU efficiencyGEN and SIM workflows need less than 2 GB/Ihread (0.5GR/Thread 0.1 MB/s/coreLess than 0.1 MB/s/coreI/O demand per process5 MB/s/coreSMB/s/coreGEN and SIM workflows and the expenses of CPU efficiencyLess than 0.1 MB/s/coreI/O demand per production job20 GB/Thread coll20 GB/thread local satial bio and thigh performance FS is satial big accel at the satial big accel at the performance FS is SIM)No sizeable local space satial big accel satial big accel at the satial big accel at the satial big accel at the performance FS is SIM)No support satial big accel satial big accel satial big accel satial big accel at the performance FS is satial big accel at the performance FS is SIM)No support satial big accel at the and an and an and an accel at the and and an and an accel at the accel at the accel at the price at the accel at the price

HPCs integration challenges in summary

Integrating HEP experiment workloads on HPC systems poses technical issues and challenges in two distinct areas:

- Data management and submission of jobs
 - how to schedule jobs on HPCs, how to connect them with experiments services, how to access experiment software and calibrations ..
- Development of the software applications
 - HEP software has so far almost exclusively been developed for traditional x86 CPUs, while HPC systems provide (and will provide more and more) their computing capacity using Accelerators (GPUs /FPGAs) and a variety of technologies such as PowerPC processors

https://zenodo.org/record/3647548#.YKzxqqIzZH5

https://cds.cern.ch/record/2016011/files/actchep2015-linennumbers.pdf

ATLAS: Grid Like execution @HPC

HPC provides the external connectivity on the nodes and enables access to cymfs software area

- Any ATLAS workflow can be executed. Push and pull mode supported although push mode is preferred on HPCs without close Storage Element. Dedicated agents (ARC-CE or Harvester) are used
- The data on the grid is typically accessed directly with xrootd protocol or copied to/from local node scratch space.



VEGA alone could be able to cover the full **ATLAS pledge !**



CINECA_HPC covered up to 10% of the total CNAF

ATLAS: Workflows for 'closed' HPCs (the hardest case)

Limited connectivity execution: when the nodes do not have outbound network access.

- The ATLAS software can either be installed locally on a shared file system or provided through fat containers on HPCs that support singularity or shifter. Such HPCs typically run ATLAS Event Generation or Geant 4
 Simulation where conditions data can be stored in a local SQLite file.
- The ATLAS agent (ARC-CE or Harvester resources broker) receives the payload description including a list of input files. This list is then processed by the agent's data-transfer subsystem. Dedicated data-transfer nodes can be transparently used and tuned for transfer performance. The completed payloads are processed by the ATLAS agent and the output files are then transferred to the pre-assigned remote Storage Elements.



Harvester has been interfaced to multiple Cloud resources (Google Compute Engine, OpenStack), several US DOE High ranking HPCs (Titan, Theta, Cori, US NSF HPC Frontera)

HPC and cloud resources

* O2 = Online-Offline Project

(ALICE Run-3 software)



- Thanks to the new O2 simulation and reconstruction code (Run 3) possible to fully exploit the multi process features
- Significant progress has been made to incorporate HPC and cloud resources in the standard ALICE Grid workflows
 - Multicore queues at CERN used to test and benchmark the O2 MC code
 - Intel based HPCs (Marconi @ CINECA, Cori and Lawrencium @ LBNL) have been used for the O2 MC challenge
- Future steps:
 - porting the O2 code to Power 9 and ARM platform

LHCb: distributed computing integration

- Support HPC with external connectivity (multi-core allocations)
- Tools to exploit HPC with no external connectivity are in progress

••• Tackling the distributed computing challenges

Overview

- 1 main variable directly affects the chosen solution (push, pull):
- + Do WNs have an external connectivity? yes (or only via the edge node), no
- Other variables generate variations that can be added up to the proposed solution:
- $+\,$ Is CVMFS mounted on the WNs? yes, no
- $+ \$ Is LRMS accessible from outside? yes, no
- + What type of allocations can we make? single-core, multi-core, multi-node
- $\Rightarrow~$ We will go through different cases: from the easiest to the hardest one

••• Tackling the distributed computing challenges

PeePush model: No Ext. connectivity, No CVMFS

In Progress: v7r2? VO action

As it was already said, SC do not provide CVMFS by default. CVMFS-exec cannot be used in this context.

Subset-CVMFS-Builder

- Run & extract CVMFS dependencies of given jobs
- Use CVMFS-Shrinkwrapper [1] to make a subset of CVMFS
- Test it & deploy it on the SC shared FS







HCb (F. Stagni et al. @CHEP 2019) arxiv:2006.13603



LHCb multi-core Grid submission

- DIRAC Matcher service
- One pilot per node, partitions the node for optimal job allocation



LHCb MC Simulation software

- RAM enough for ~80 single-process jobs
- Using a multi-process application, may use 136 threads with 8 MPx17 jobs, but with minimal gain in event/sec throughput





From 2019 to 2020 HPCs contribution had x3 increases

From 33.6 to 108 kHS06/y. _

9.5

0.6

0.9

0.4

0.5

0.1

0.4

0.3

0.0

0.1

1.8

1.9

17.0

1.3

0.0

0.3

0.1

35.3

GPUs and heterogeneous computing

- Since 2019 (Patatrack) CMS has started woking the porting the code to CUDA
 - Currently offloading about 25%
- Currently investigating the adoption of High Level Frameworks such as Alpaka, Kokkos and SYCL
 - allows writing single common code basis
 - EuroHPC/TEXTAROSSA Project
 - CNAF and PISA involved

25% CPU time reduction when offloading pixels, ECAL and HCAL to GPU

(GPU used at 25%)

OBJ



See S.Dal Pra @ ISGC

A few words on our national (INFN) landscape

Integrating Tier1-CNAF and CINECA in the most transparent way to the end user/ experiment



External networking enabled to CNAF and CERN

A **HTCondorCE/Slurm** allowed on one CINECA login nodes

Partial routing to CNAF storage / squids over the dark
 fiber @ 40 Gbit/s (Will improve with next machine)

Enough to guarantee access to CNAF storage and conditions for experiments' workflows + HTCondor pilot communication

Moreover: We use it also for accessing external data via proxies (xrootd)



Looking at ahead: Marconi 100

Started working on the integration of the Marconi100 at CINECA (Power9+GPU), INFN got 3.5 MCoreH for 2021

- Enable **multi-arch support** for CMS production/analyses jobs
- Perform the **physics validation** on Power9 for its utilization in physics analyses

Currently:

- CMS software stack fully ported to POWER9 -
- Established a complete integration stack of the CMS Workload -Management (WMAgent and CRAB)
- Successfully (technically) tested analyses and Release -Validation workflows

GPUs:

"Performance" in A.U. of the HLT application (the higher the better)

Thanks to LHCb for lending the node!	Platform	CPU only	GPU Type	CPU + 1 GPU	CPU + 2 GPU	CPU + 4 GPU	HS06 CPU
	2x Intel Xeon 6130	24	T4	33 (+37%)			865
	2x EPYC 7502	58	T4	75 (+29%)	75 (+29%)		1832
	2x EPYC 7742	100	T4	127 (+27%)	129 (+29%)		3170
In contact with IBM to optimise compilation flags	2x POWER9 (Marconi 100 Node @ CINECA)	18	V100	23 (+28%)	23 (+28%)	23 (+28%)	need new compiler flags

	Rank	System
MARCONI - 100	11	Marconi-100
Nodes: 980		
Processors: 2x16 cores IBM POWERS	9 AC922 a	at 3.1 GHz
Accelerators: 4 x NVIDIA Volta V100 0	GPUs, Nv	link 2.0, 16GB
Cores: 32 cores/node		
RAM: 256 GB/node		
Peak Performance: ~32 PFlop/s		
Quick startup guide		

A first "scale" test performed as well



This first attempt is really promising, the focus now is on

- setup consolidation
- larger production



A few more words on HPCs with no-network

- Negotiate with the site for at least minimal connectivity (as for CINECA: ok for CERN and CNAF)
- Encapsulate network on something else (for example BSC(*): via FS, but it needs a service by service approach...)
- Deploy large containers to avoid needing CVMFS/Conditions (but condor management traffic? Stageout? Input data?)
- Create edge services to transmit / encapsulate / translate accesses
 - \circ $\,$ A Squid is (also) such a thing
 - ATLAS Harvester (**) is (also) such a thing
 - \circ ArcCE , DIRAC
 - \circ \qquad An xrootd proxy is (also) such a thing
 - How many services do we need to bridge one by one? (xrootd, srm, webdav, condor connections, dashboard reporting, DBS, DAS, Frontier, voms connections ...)
- We need a lower level solution, below the application layer
- *: Exploiting network restricted compute resources with HTCondor: a CMS experiment experience
- **: Harvester : an edge service harvesting heterogeneous resources for ATLAS

- Having routing by kernel to a NAT is clearly a solution (but it needs the site admins to agree/implement)
- Having a VPN from the site to "somewhere else" is a solution (but it needs the site admins to agree/implement)
- If the sysadmins agree with deploying the previous two, the problem is not existing...

... but difficult to find the agreement in some (known and future) cases



See M.Mariotti @ISGC



Summary

- HPC will play a key role in the future computing of the Experiments
 - Requires both technical and political effort
- Although many progresses have been done by all the experiments, not all the workloads can't be efficiently run on HPCs
 - Suitable for opportunistic computing but we are still not ready for replacing grid-based resources (pledges) with HPCs
- A lot of work still needed to deploy experiment software to fully profit from the accelerators and variety of processors at HPCs
- At national level, INFN the investment done at MarconiA2 with (Grant LHC) is starting paying back with ongoing Marconi100 system
 - in turn will be of fundamental importance for Leonardo, the pre-exascale system
 - Moreover it is a valuable experience in view of the new technopole datacenter under construction

Crediti per il lavoro di integrazione CNAF-CINECA, contributi ad R&D, supporto etc.

- Stefano Dal Pra
- Stefano Zani
- Gaetano Maron
- Luca dell'Agnello
- Lucia Morganti
- Daniele Cesini
- Vladimir Sapunenko
- Mirko Mariotti