

WLCG Strategy and Facilities

Overview of strategy towards HL-LHC computing

Introduction	1
Summary of the current understanding of HL-LHC computing requirements	1
Main challenges for HL-LHC computing	2
Hardware Outlook	7
Facilities	12
Traditional WLCG Facilities	12
Opportunistic resources	14
Clouds and HPCs	14
Software Adaptability & Efficiency	15
Networking	16
Authentication Authorization and Identity (AAI)	18
Security & Security Coordination	20
Operational Model and Service Provisioning	21

Introduction

Summary of the current understanding of HL-LHC computing requirements

In elaborating their computing models and strategies for HL-LHC the experiments are working with different baseline assumptions concerning the LHC and data taking parameters. Those assumptions try to be as realistic as possible. At the same time the experiments want to ensure they are prepared for a challenging nominal year of HL-LHC physics. To facilitate comparisons between the models of the two experiments we agreed on a common set of parameters to which we refer as the "LHCC review parameters". Those parameters are summarised in Tab. 1. While they differ from the baseline parameters of the experiments and they might not be the most realistic ones, they offer a common ground for comparison. The experiments therefore provide in their respective documents some estimates based on the LHCC review parameters in addition to their baseline parameters. In the future, we will be happy to use a different common set of parameters under the advice of the LHCC.

Live Time (pp)	7M seconds/year
Pileup	200
Collision Energy	14 TeV
Trigger Rate	10 kHz
1st LHC "nominal" year	2028
Integrated pp luminosity	500 fb ⁻¹ /year
Flat Funding Scenario	+10% hardware / year (disk, tape, CPU)

Table 1: "LHCC review" set of parameters

The ATLAS and CMS documents elaborate all the details of the respective HL-LHC computing requirements. At first glance, the needs for the first production year of HL-LHC are still several factors above what a constant spending for computing hardware can provide. On the positive side, the scenarios where a set of aggressive R&D can be completed provide estimates more compatible with the currently anticipated budgets. Such R&D needs to be properly prioritised, supported and funded and this is the focus of the documents submitted to this review.

Main challenges for HL-LHC computing

At a high level the main challenges that we must face in looking towards HL-LHC are the following:

- Fitting within a restricted cost envelope for computing;
- Managing a new scale of data volumes - at the multi-Exabyte scale;
- Adapting the software investment to a new era of rapidly and continually changing heterogeneous computing hardware;
- Bringing out the potential commonality in software tools and services across experiments and across infrastructures;
- Recognising that LHC is no longer alone in these challenges - this has clear benefits but also potential limitations.

Fitting within the cost envelope

Given the current understanding of the needs of the LHC experiments for a nominal year of running during HL-LHC, outlined and summarised above, there are consequently a number of high-level challenges that must be addressed. It has been clear for several years that the budget outlook for HL-LHC computing will be constrained, with the key message from funding agencies that the LHC community must keep computing costs within a long-term "flat-budget" scenario. In such a scenario, budgets available for software and computing would be fixed at today's level, most-likely without any adjustment for inflation. This is a real challenge as a number of factors conspire against this.

Firstly the overall level of requirements for computing resources is significantly higher in Run 4 than in Run 3, by several factors. While this is a significant improvement up to an order of magnitude since the first estimates were made a few years ago, it nevertheless represents a problem that is not simple to address. The naive cost improvements that we had previously been used to, following Moore's Law, are no longer applicable. The costs of hardware are subject to large fluctuations even on the timescale of months, and the outlook for the future improvement is much less than the 20-25% per year historically seen for capacity/cost, but close to 5-10% per year. This, coupled with large market-driven uncertainties, means that the affordability of computing on the HL-LHC timescale is currently almost impossible to predict with any certainty. In essence, for many components it is not technology that drives the prices, but market forces, lack of availability, and uncertainty over the future of certain technologies. This is typified by the uncertainty over the future of tape, and the technologies for hard disks that may be available.

Experiment computing models

The experiments' computing models, experiment-specific and common application software, all clearly have a direct impact on the total amount of computing resources required. The models are described in the separate ATLAS, CMS and software documents. Among the model parameters and assumptions, we can identify those which have a larger impact on the resource needs.

Firstly, simulation has always been a significant fraction of the overall need for CPU, up to 50%, with a large variation among the experiments. In the past the event generation part has been insignificant in terms of CPU needs, but as the need for greater precision grows, the generators must calculate at higher orders for many processes, leading to a huge growth in the computing requirement for event generation. This problem must be addressed together with the generator communities, offering them the computing expertise of the experiments and WLCG community in a collaborative way. This effort has begun, with several workshops held; but there is still a significant effort required in order to achieve a realistic cost for the accuracy requested. This requires good investment from the generator community and additional support for technical and physics improvements required.

The other aspect of simulation is that of the detector response, mainly through Geant4¹. As noted, overall Geant4 accounts for around half of the overall WLCG computing budget. Geant4 is common to all of the experiments, so it is clear that performance improvements will have a significant and common impact on all experiments. This is an area that the HSF and simulation communities are addressing, but it is clear that this must have a significant and ongoing investment in effort and expertise. This may require long-term and deep changes in the Geant4 code-base in order that it can make appropriate use of modern compute facilities, and to adapt to the ongoing evolution in processor technologies.

Simulation is clearly also an area where new types of resource (HPC, GPU clusters) could have a significant impact, but also requires appropriate software modernisation. In addition, other

¹ <https://geant4.web.cern.ch>

strategies, such as so-called “fast simulation” must be integrated with the full simulation of Geant4, where appropriate.

The Geant4 collaboration is planning for how the software should evolve, including several R&D efforts. However, it is clear that this is a specific area where long-term investment from funding agencies would have a clear and beneficial impact on the overall resource needs of HL-LHC, and beyond as well as revitalising the development community with a younger generation of physicists.

The second main area of concern for compute requirements was in event reconstruction. In the regime of pile-up between 140 and 200, the use of Run 2 and Run 3 versions of the reconstruction software leads to exponential increases in reconstruction time. The experiments have already invested significant effort in this area, and reconstruction is a more manageable problem, while still a major resource consumer. The progress here is described in the experiment computing model documents. This is one of the ways in which the overall compute requirements have been so reduced during the past few years.

In terms of analysis there are several strategies being pursued. Since analysis data is often accessed many times it can potentially lead to very high I/O rates, and resource contention. Both ATLAS and CMS are building very reduced data formats for analysis that can be used in the majority of analyses. This will allow the data sets to be cacheable close to the analysis compute resources and avoid frequent downloads. The aggressive size reduction of these formats means that the scale of the overall data management problem can be reduced significantly, as long as their use is mandated wherever possible. There are also ongoing discussions about the potential use of dedicated analysis facilities, for example with a high IOPS capability and/or high performance internal networking. It is not yet clear whether this is something that is really needed, or desirable for cost and complexity.

All experiments are investigating machine learning applications to augment or replace some traditional workflows. ML applications often concentrate the bulk of the processing time in training, which can introduce latency-critical steps. The IO needs for training may be large and chaotic and may impact the facility storage requirements. Industry and the open source community have invested heavily in hardware and software to accelerate machine learning applications. Facility planning will need to consider the role of ML applications and the storage systems and accelerated processors needed to support them.

Operational costs at sites

One of the trends seen in recent years is the consolidation of effort across WLCG sites, particularly at Tier 2s. The strategy towards HL-LHC must take into account the cost of staff at the sites needed to keep the operation going. One area is the need for managed storage at smaller sites, that usually requires skilled staff to operate. Part of the data management strategy (see the DOMA document) is to remove as far as possible the need to have a managed storage system at many sites, replacing it with a simple cache where possible, thus allowing the site to operate with fewer staff. Another option is to enable the federation of

sites (for example within a country) in order to optimise both hardware and personnel costs across the federated sites. Again, this is part of the DOMA strategy.

In addition there are a number of R&D efforts under way to explore the possibility of remote operation of certain site services. The feasibility of such an approach depends on local policies, but is being investigated as another potential mechanism to reduce the need for operations staff on the spot. In a similar vein, there are R&D efforts into further automation of operations tasks, and the use of AI technologies in order to facilitate and automate operational aspects. The results of these efforts are likely to be deployed in production as soon as they are demonstrated to bring operational efficiencies.

Managing Exabyte scale data

The expectation from both ATLAS and CMS during Run 4 is close to 1 EB of data to be collected each year, with a few times that in derived and simulated data. Thus each experiment will be in a multi-Exabyte per year regime of data to be managed. This represents a significant challenge compared to the current situation.

There are a number of explicit strategies focussed on addressing this challenge. These are centred around the following areas:

- Efficient and more cost effective data management;
- Reduction in the size of derived data that is to be distributed;
- Management of operations costs.

The primary strategy to have cost-effective and efficient data management services is addressed via the Data Organisation, Management, and Access (DOMA) project. This work is treated in the separate DOMA document, and contains a number of R&D and prototyping activities that update the underlying data management tools and protocols, as well as addressing at a high level the policies for data management. One of the most significant aspects is that in the future data management must be able to serve data to heterogeneous compute resources, from a pool of well managed large scale data stores, with automated and policy driven replication and load balancing within and between the stores (the so-called "Data Lake"). This allows the optimisation of managed data at the larger sites with large storage systems and experienced staff, and the more effective provision of compute at both traditional grid sites, or opportunistic use of clouds and HPC sites without the need for the installation of special data services in them. This should result in a better use of available funding for compute services at many sites without the need for expensive storage systems or additional staff.

This change of strategy that moves away from replicating data "everywhere" to serving data when and where it is needed is more efficient as only data that is actually being processed will be transferred. In the past a lot of data was moved in advance to ensure availability and accessibility. However in practice a lot of that pre-placed data was never (or rarely) accessed, and over time it became clear that the network was perfectly capable of supporting a remote access model.

The strategy is also aided when considering the trend towards very small analysis data sets, based on nano-AOD (or other experiment-specific names). For a large fraction of analyses only a very small subset of event data is required; these reduced data sets will obviate the need for large transfers of analysis data, and could even be cached in various points on the network for access. Again this helps reduce the amount of data transfers and management, and the need for large managed storage systems.

These strategies go towards managing the overall cost of operations at sites and with the data transfer services. Of course they rely upon efficient and reliable networking with appropriate bandwidth available. As we have seen in the past decade, this is achievable and we have been able to rely on the networks fully, and early strategies to manage potential network unreliability have been removed or greatly simplified.. The developments around data management are being done together with the networking communities across the world. There is concrete planning for network provisioning that is required, discussed later.

Heterogeneous computing and portability

In the last several years, the era of "x86-only" processor dominance has started to change. Even within the "x86" processors the introduction of vector units, multi-threading, as well as other forms of parallelism support has meant that most of HEP software is not using all of the available processing at maximum efficiency. In addition, the introduction of new processor types (e.g. non-x86-64 architectures), and the ubiquity of coprocessors, such as GPUs, has made it clear that we must be able to both adapt our software to make use of the available processing, and also to be able to port it to a more parallel environment. To do that effectively requires re-engineering of the software and algorithms in most cases.

The resources will no longer be mainly provisioned through dedicated grid sites, with HPC systems and cloud compute may also be part of the resources that we must be able to take advantage of. Each of these types of facilities may well deploy a heterogeneous set of processors and accelerators. While the access to and operation of these facilities may be quite different, there is a common software portability challenge that must be faced no matter which type of facility is used.

Such a re-engineering and porting effort is also not a one-off enterprise; the future is one where there will be a continual evolution of processor types and capabilities. Our community must equip itself to address that challenge in an ongoing and efficient way. The HEP Software Foundation (HSF) was set up to help address these challenges. The strategies for re-engineering and porting are outlined in the software document.

Common software tools and services

The second aspect of the software challenge is to be able to continue to build on the common software solutions in a far more aggressive way than in previous years. During Run 1 and Run 2 we have seen commonalities arising, in addition to the long standing examples of Geant4 and ROOT. As well as common event generators, and run-time libraries such as TBB, all experiments now use CVMFS for software delivery, and with the DOMA project we see the potential for a real common data management service.

It is essential to continue to identify the commonalities and benefit from them, as the available effort for many competing but similar solutions is expensive. This is particularly true in the environment where a continual effort to port and adapt software to new facilities will now need to be made. It is therefore essential that commonalities must be taken advantage of wherever practical. It is also important to note that other related sciences now are facing similar challenges to HL-LHC and commonalities between LHC experiments can also benefit those other communities, and eventually build a stronger overall community. It is also very important to recognise that many of those other sciences will share data centres with the LHC experiments, and it is not realistic to expect those data centres to operate different solutions for different experiments where it is not necessary.

Hardware Outlook

The worldwide revenues for the general semiconductor market were recovering from some weak demands during the last 12 months.

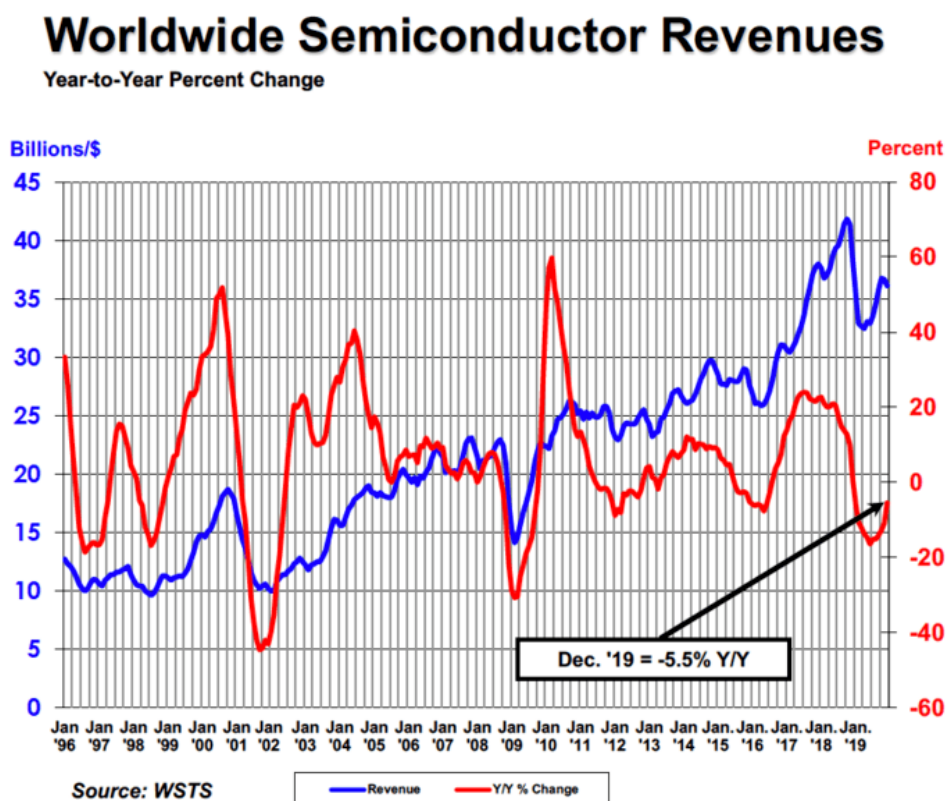


Fig. 1: Worldwide semiconductor revenues and the corresponding yearly fluctuations

This situation will of course now change due to the COVID-19 crisis, which has started to create a major worldwide recession. All depends on the structure and timing of this downturn. A high unemployment rate and large-scale business bankruptcies would severely affect

the demand for PCs, notebooks, smartphones, etc., and then indirectly the demand for server computing and storage in the cloud providers (plus telecom industry). R&D activity may stall or be reduced. A short crisis might create a rebound effect. More digitalisation, more teleworking would cause large scale investments in the cloud infrastructure (server, storage, networking) and trigger even more technology improvement. The evaluation of these possible developments can only be done in a credible manner at the earliest in the middle of 2021.

For the future of HEP computing and especially for the HL-LHC computing predictions one can today say that the possible technology evolutions are in principle doing well. Credible technical roadmaps in the area of processors, storage and networking exist. However, these roadmaps are currently thrown into doubt by the current Covid-19 crisis, and we will need to closely monitor the evolution.

A few more details in the areas of Processors, Disk storage, Tape storage, Accelerators and Networking follow.

Processors and Servers

The server market is still dominated by processor lines from Intel, with a market share of about 95%. In the last 12-month AMD was able to grab a market share of ~5% with their new EPYC processor line and analysts expect a further increase during the next 18 months. This strong competition has led to large price variations in the server market. At the same time, we see a resurrection of the anticipated ARM servers (Ampere and Marvel). The AMD-Intel competition and the anticipated economic down-turn will make it very unlikely that an ARM server model will succeed in the near future; pushing a new architecture at scale in the market requires a very large amount of up-front investments.

Processor technology evolution itself looks very healthy. The foundries are using the 7nm process for the current generation of server chips (the Intel 10nm process is equivalent to the 7nm process of Taiwan Semiconductors (TSMC) used for AMD chips). Roadmaps down to the 1.4nm level exist, but one expects the first issues at the 3nm level, as one has to include new materials (e.g. Ge) and new methods (e.g. nanosheets, nanowires); these technology changes require very large and ever-increasing investments.

In addition to the processor price fluctuations, we have seen during the last 2 years an increasing level of price volatility in other server components, specifically memory (DRAM) and SSDs (NAND). The following diagram shows the evolution of server price/performance (CHF per HepSpec) at CERN during the last 15 years and some possible extrapolations.

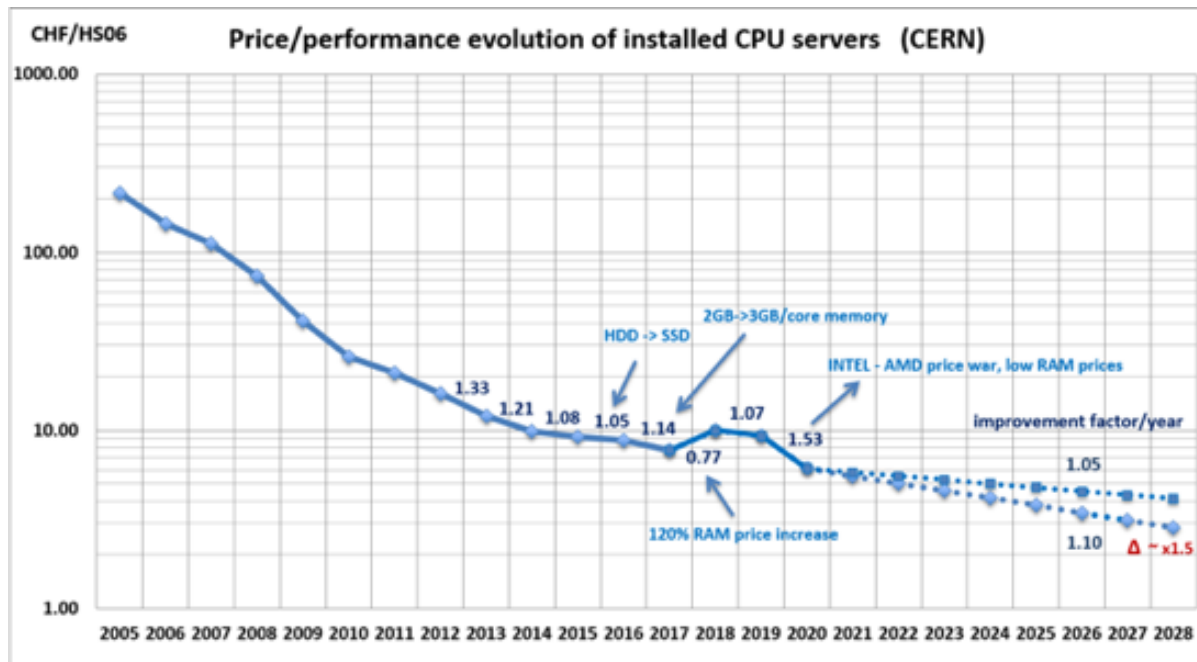


Fig.2 The price/performance evolution as a function of time of CPU servers at CERN

Again, how the current crisis and its after effects will affect the price/performance evolution of servers during the next years is far from clear.

Disk Storage and Servers

During the last 5 years the yearly amount of HDD unit shipments have reduced by a factor 2, while the amount of storage capacity (EB) increased by nearly a factor 2. The low end of the HDD market is essentially disappearing, as they are being replaced by SSDs: at the end of 2020 all new notebooks and the majority of the desktop PC shipments will likely have only SSDs. All three remaining HDD manufacturers (Western Digital, Seagate, Toshiba) are investing heavily into SSDs, e.g. today WD has already an SSD market share of 17%. In the HDD sector the companies are investing into a still growing nearline high capacity drive market (cloud provider, backup, video surveillance, etc.). They have a 20% share of the total HDD shipments while providing already 45% of the revenues.

The technology roadmaps for the next 6 years are credible with a variety of new density improving methods (HAMR, MAMR). However, these new approaches have already been delayed by several years in the past, due to their sophistication and cost.

The cost difference between SSDs and HDDs in terms of price/TB has been diminishing, with SSDs still about a factor 5 more expensive. This will not change during the next few years as the NAND fabrication units would need to increase by factors with estimated investments of more than 100 B\$. Only 10% of today's disk storage capacity is provided by SSD shipments.

The storage cost calculations need of course also to include the infrastructure around the bare HDDs. The following diagram shows the storage (HDD) server cost evolution during the last 15 years at CERN and some possible extrapolation (the price/GB numbers include server mirroring = 2 copies of the data).

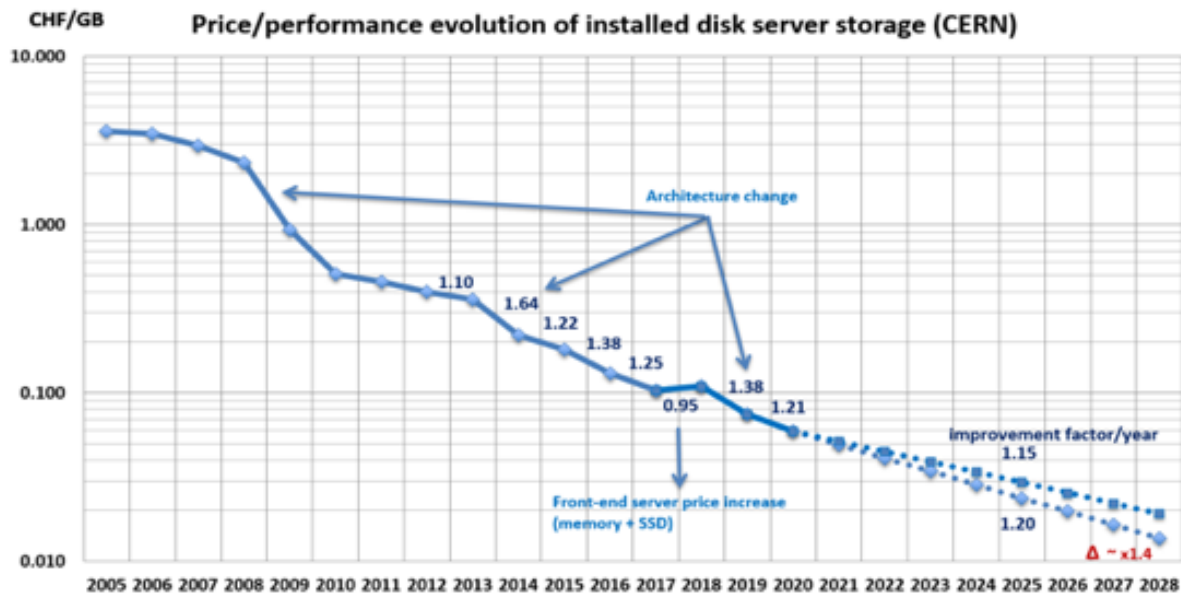


Fig.3 The price/performance evolution as a function of time of disk servers at CERN

Likely evolution of this is extremely uncertain as noted above.

Tape Storage

The technology roadmap for tape storage is actually in good shape, i.e. the LTO consortium has published a preview for generation 12 (up to 192 TB/cartridge) and prototypes in the lab have reached magnetic densities corresponding to 300 TB cartridges. The possible problems are again due to the market situation. Overall revenues in the tape market are estimated to be below 1 B\$, and the LTO consortium has stopped reporting sales numbers (#tapes sold, EB shipped) 2 years ago. Only IBM is evolving the tape-head technology for tape drives. The two companies (Sony, Fujifilm) manufacturing the tape media have been entangled in a severe patent-war in 2018-2019; for many months it was impossible to buy LTO-8 tapes. The current COVID-19 crisis and the following worldwide recession might increase the possibility of a complete break-down of the tape market. The HEP storage systems would be able to cope with a possible disappearance of tape technology, by replacing it with HDD storage. Software and management wise this would not cause a big issue, but of course the costing difference will be a major problem. In terms of \$/TB cost, an HDD storage system is about a factor 5 more expensive than a tape-based system.

Accelerators

The GPU market is healthy with some small percentage growth rates driven by gaming (PC and consoles), engineering and the various machine learning (ML) activities. About 40 million discrete graphics cards were sold in 2019, of which only less than 1 million are in the high-end area (expensive, high profits). There are only two companies left sharing this market (Nvidia 73%, AMD 27%). This has to be compared with about 400 million GPU shipments for PCs and notebooks (low-end) and about 2 billion embedded GPUs in smartphones and tablets. The important part of the GPU market are the discrete graphic card sales and not the embedded ones (e.g. Intel CPU+integrated-GPU or the smartphone GPUs), because only the high-end systems are driving the scientific computing area.

The used foundry technologies are following very closely the ones for processors. The current generation of Nvidia GPUs is manufactured by Samsung with a 12nm process, and it is moving

to 7nm this year. TSMC is already using 7nm for the current generation of AMD chips. Intel will introduce their new graphics engines (Xe) this year possibly using their 10 nm process.

There are two more types of accelerators available: FPGA and ASIC/TPU. FPGAs have a large, but more low-level market in embedded systems, while TPUs are aiming for a new market (ML). High-end FPGAs are very expensive and have a complex programming model. Market Research companies (e.g. Tractica) forecast only very small market shares for FPGAs in the Machine Learning area. GPU's have a large advantage here, as they are able to provide an advanced software ecosystem (e.g. CUDA) and have a very well established and large market (Gaming) as a driving force. There is currently a flurry of small companies with new chip designs for ML (e.g. Graphcore, Nuvia, Tenstorrent, Kneron, Cerebras, etc.) and the large cloud providers have started to offer specific ML services with their own chip design (e.g. Google TPU, Alibaba HanGuang ASIC, Microsoft FPGA). Many designs will not survive or will be bought by companies like Intel.

The most reliable/credible accelerator architecture for general HEP scientific computing in the near future seems to be GPUs.

Networking

The total international fibre link connectivity between the continents (e.g. transpacific, transatlantic, etc.) this year will reach about 3500 Tbps, with an estimated growth rate of 35% per year in the future. Cloud providers are the biggest driving force. Only 20% of this capacity is used today: E.g. the currently available transatlantic link capacity is capable of transferring 40 EB/month.

The worldwide IP traffic will reach this year a value of 200 EB/month, of which 80% is used for commercial video streaming (Netflix, YouTube, etc.). The expected growth rate is 25% per year. The worldwide IP traffic in 2028 (first full year of HL-LHC) is expected to be about 15000 EB per year.

Inside WLCG we see a yearly total network traffic corresponding to a data movement of about 2 EB/y. For HL-LHC we assume at least one order of magnitude of increase, which would yield an amount of 20 EB/y worldwide transfers in 2028. But the planned storage/transfer activities (e.g. DOMA) are aiming for a large scale optimizing and reduction in transfers. The DOMA activity will test the hypothesis that data delivery when needed requires overall less bandwidth than the current strategy of pre-placement of data sets.

In the data centre itself (taking CERN as an example) we have seen considerable price reductions for the network equipment (switches and routers). The TCO contribution of the network (NIC+switch share+router share) to the total cost of CPU and disk servers is below 10%.

Total I/O in the CERN centre has reached about 3 EB data movement per year. The vast majority of activities (MC, processing) is low I/O on average. High I/O jobs (analysis) can be spread evenly across the system without creating network bottlenecks. CPU servers are still using 10 Gbit NICs, while PB scale disk servers have moved to 25 Gbit NICs. The backbone router infrastructure uses 100 Gbit connectivity. For HL-LHC one will only need a slow evolution of the centre network infrastructure.

The network technology driving forces are the big cloud providers, with a pace which is actually faster than the HEP community really requires. Overall the network developments in terms of technology and markets are going very well and will not be a problem for HEP

computing. More the contrary, they will provide an opportunity. While this is true in general, there are several aspects that need to be closely watched and considered: the transatlantic capacity, the network organisation in some regions such as Asia, and the relatively long time needed to plan an increase in capacity in several regions.

Facilities

Traditional WLCG Facilities

The WLCG Memorandum of Understanding (MoU)² and in particular its Annex-3³ defines the role of sites and federations in the organisation and classifies them according to a hierarchical structure: Tier 0 (host laboratory), Tier 1s and Tier 2s. The roles and functions of sites have evolved since the time these definitions were made and today several roles are covered more flexibly by tiers of different kinds. The Tier 1s continue to hold a special function both at the level of operations and organisation. Those large national centres are responsible for the long term curation of the LHC physics data and therefore we expect a shorter response time to problems with respect to Tier 2s, 24/7 in the limit of their capabilities, and special attention to data loss and disaster recovery. Critical incidents or unavailabilities of the Tier 1 services are followed by Service Incident Reports, discussed at the WLCG Management Board and possibly escalated if needed. The Tier 1 centres might be asked to run production services beyond storage and compute upon negotiation. Finally, Tier 1s cover a special role in their region/country as they are the prime point of contact and support for the other centres in that region, primarily the Tier 2s. We do not foresee reviewing the definition of Tier 1 or Tier 2 centres for HL-LHC, while we might decide to update Annex-3 to reflect the current situation. We do not expect to re-sign the MoUs, according to the possibly updated Annex-3, as it would bring little value to the collaboration. We strongly encourage new sites and federations joining WLCG, discussing their role with the WLCG and CERN management and the experiments they would support.

The model we foresee for facilities at the time of HL-LHC is a flexible one, where particular workflows are not bound to be supported at a particular tier type. On the contrary, we expect sites offering different capabilities and the experiments organising the activities to best exploit those capabilities. For example, some sites could offer storage for long term data retention at possibly high latency - what today is tape storage - others might offer online storage in addition or instead. We expect a good fraction of the sites to run stateless storage - a.k.a. caches - for the purpose of supporting data processing, while others may offer no storage at all, relying on storage services of a close and well connected site, and focus on providing larger processing farms. We are investing in tools and services to expose those capabilities and monitor them effectively. The experiments will then be able to organize different workflows taking advantage of those capabilities and the support level offered by the sites. Some facilities, because of a well defined set of capabilities might specialise in specific

² <http://wlcg-docs.web.cern.ch/wlcg-docs/MoU/MoU-blank-example-28APR2015.pdf>

³ http://wlcg-docs.web.cern.ch/wlcg-docs/MoU/Annexes/Annex3_min_membership.pdf

workflows, for example organised analysis, data filtering, reconstruction, processing of special datasets.

Many funding agencies have anticipated the possibility of a consolidation into fewer, larger facilities, particularly for storage services. Those facilities will serve as a backbone for data curation and data provisioning. The storage (data provisioning) services will be characterised by different Quality Of Service classes, so that the experiments will be able to optimise the usage of storage for different use cases. The processing capacity might be co-located with the storage service or not, allowing a local cost and efficiency optimisation between storage and compute.

A content delivery network consisting of services for caching and buffering will therefore be needed to serve the data to such compute capacity. The binding between compute and storage services at the same site is then less critical. The model offers sites the possibility to deploy the set of services that they can best procure, support and operate. The aim is reducing the overall total cost of ownership, or, in other terms, providing a better resourced service at the same cost.

The model heavily relies on efficient network services with enough capacity, at the level of both LAN and WAN. The network proved to be a very reliable service during LHC Run-1 and Run-2 and therefore the strategy for Run-3 and HL-LHC is to further rely on the network. A recent analysis of hardware trends in the hardware outlook [section](#) of this document notes that network bandwidth has increased at a faster pace than CPU and storage over several years, with small fluctuations. An analysis of the network needs for HL-LHC can be found one of the following sections.

Implementing this model, often referred to as a "datalake", is one of the core activities of the DOMA project and is described in detail in a separate document. The plan is to demonstrate with the datalake model an overall higher efficiency and processing capacity at the same cost, both in terms of operations and hardware.

Most WLCG facilities also support other HEP experiments and more generally other sciences. WLCG feels it is in the interest of scientific computing to organise ourselves to collaborate on the same infrastructure, sharing services and expertise, rather than compartmentalising. This we understand is also the preferred evolution as seen by our funding agencies. As input to the Open Symposium for the European Strategy of Particle Physics in May 2019, we proposed a possible evolution of the WLCG collaboration⁴. The idea is to separate the purely LHC computing project part from the infrastructure related aspects and expanding the latter to be inclusive initially of other HEP experiments. For the implementation, we decided to take a pragmatic approach and start creating synergies with the Belle-2 and DUNE collaborations, which are now invited to participate and input to the WLCG/HSF workshops and the Grid Deployment Board and and to attend the WLCG Management Board as observers. The co-location of WLCG and DUNE on network resources already offers an opportunity to understand

⁴ http://wlcg-docs.web.cern.ch/wlcg-docs/technical_documents/HEP-Computing-Evolution.pdf

how to practically share the infrastructure with mutual benefits (see Network section of this document). The actual plan presented in May 2019 will be reviewed in light of the first experiences and following several recent discussions.

Finally, WLCG plans to evolve its services and policies according to FAIR (Findable, Accessible, Interoperable, Reusable) data principles. HL-LHC is one of the science projects in the ESCAPE European funded initiative and drives the evolution of a data infrastructure for open science in HENP and astronomy. The work in ESCAPE is done in synergy with the DOMA project and focuses on FAIR data. Services such as Zenodo, today in use by many scientists of very different domains, and the CERN Analysis Preservation portal, can be evolved and further integrated with the WLCG services and offer an open data frontend to some of the HEP data. Policies on data preservation and open data are discussed in a dedicated working group organised by the CERN management with the participation of the experiments. WLCG will follow and implement the recommendations.

Opportunistic resources

The WLCG experiments benefit from a relatively large amount of opportunistic resources, on average 20% beyond the CPU pledge level, with large fluctuations between experiments. There is today no such concept as opportunistic storage in WLCG. However the R&D work in DOMA about caching creates the possibility to use opportunistic storage as a buffer for data processing. While there has been always a worry about the long term implications in relying on opportunistic resources, they have proven to be available over the last decade at a relatively constant level and therefore some experiments consider them for their future plans. For more details we refer to the ATLAS and CMS documents.

Opportunistic resources today are provided for the vast majority in the form of Grid resources at WLCG and non-WLCG (Tier 3) sites. This leverages the flexibility of the Grid model in sharing resources between different communities and integrating local resources into the system. We do not have indications about this model changing in the future and therefore we expect some level of such resources continuing to be available. The LHC experiments leverage opportunistic resources accessible through non-Grid interfaces such as Cloud and HPC facilities. We have to foresee the possibility in the future for those facilities to provide also pledged computing capacity, as it has been indicated by some funding agencies. Technical challenges will be highlighted in the next section. As an opportunity, large facilities such as commercial cloud providers and Leadership Class HPCs can provide a large amount of computing resources for peak usage and the experiments have demonstrated the ability to use them.

Clouds and HPCs

Commercial Cloud facilities offer elastic capabilities and can absorb bursts of activity without prior notice. It is not obvious however if this flexibility is needed for the expected needs of HL-LHC and probably this will be experiment specific. At the current cost, commercial cloud resources are not economically a good alternative to on-premise, particularly for heavy I/O applications given the ingress and egress costs of the cloud providers today. The funding and

the procurement model is also not obvious in some cases. So far, cloud resources have been integrated in some cases as an extension of a WLCG facility as in the examples⁵ of the CERN batch farm in the HelixNebula project, the Fermilab farm through HEPcloud and the CNAF extension to Azure and Aruba. In other cases they were integrated directly with the Workload management Systems of the experiments. In terms of future plans, WLCG needs to continue ensuring its capability to integrate transparently cloud resources so that those can be provisioned when economically viable. Part of this work will focus on ensuring that resources are properly benchmarked and accounted for in the standard WLCG tools.

High Performance Computers (HPC) are generally designed for massively parallel processing (up to hundred-thousand-core MPI jobs). Node interconnection and similar capabilities are not needed in HEP, where processing happens at the much smaller granularity of independent physics events. HPCs however are seen as an opportunity to complement the capacity offered by High Throughput Facilities in WLCG. The challenges integrating HPC centres to run experiments' workflows is summarised in a WLCG document⁶. The challenges can vary depending on the HPC facility and span very different areas. The R&D activities in the DOMA project, and particularly the ones on caching and high throughput network challenges, contribute to addressing data ingress and egress to and from those facilities. One of the main challenges for the offline software is ensuring a capability to benefit from accelerators or coprocessors. HPCs are very heterogeneous in this respect, they do not always deploy x86_64 processors and most of the capacity comes in the form of accelerators such as GPUs. The challenges and R&D activities in this area are explained in the software and experiments' documents, input to this review. We note that such work in software portability and use of heterogeneous architectures has benefits expanding much beyond the use of HPC facilities. It will improve the overall quality of the software and improve its long term sustainability. Additionally, more data centres deploy accelerators as processing capacity and it is likely this trend will also affect WLCG sites, in the need to support heterogeneous applications from different sciences.

Software Portability & Efficiency

Maintaining the software code base and user community across unavoidable technology changes is one of the key challenges for long term projects such as the LHC. Software may need to migrate several times in order to stay efficient with the evolving hardware. This will require continual adaptation of the skills of the software developers, across several generations of developers, while at the same time retaining the domain knowledge that is represented in the (current) software base. The core codebase for HEP must in future exist in a very different environment than has been the case for the past 25 years. In the future all of the code such as applications, common software, libraries and tools must be capable of running efficiently in a variety of different CPU architectures and on various accelerators and coprocessors. This will not be a one-off port, but the portability of the software environment must be an inherent and ongoing part of the software lifecycle. In order for such an activity to be successful, the strategy must be defined in collaboration between all parties involved in LHC software development. It is important in

⁵ <https://doi.org/10.1051/epjconf/201921408002>

⁶ <https://zenodo.org/record/3647548#.Xob5FS2B1N0>

fact that this is a community-wide effort, as we must ensure the engagement of the key software projects e.g. Geant4, ROOT, and the experiment core software teams. The WLCG facilities must play a role in this effort: it is in their interest for the community to build software expertise focusing on portability and efficiency, and understanding the adaptability of the software for which these facilities are being built is essential to make them efficient. Such expertise will allow the facilities to work with the applications to improve software performance on the deployed systems. At the same time, it will contribute to the decision making process with respect to future acquisitions. This has to be therefore seen as a common investment and strategy around portability and sustainability of code as we do not want different solutions for each experiment and we need the common solutions to be properly supported by the WLCG sites. This has been a common strategy at HPC centres. The certainty of x86 is gone: portability will not only be required due to differences between sites, but also to support future generations of systems.

We believe therefore that an R&D activity on long-term software portability and efficiency should be recognised as a key area in the strategy and formally organised and coordinated. This will help in the work being properly acknowledged and funded. As part of the acknowledgement, the profile of experts in this area should be given the proper relevance in HEP in the form of career opportunities, which appears not to be the case at the moment.

The “Common Tools and Community Software” document elaborates more on the aspects of software portability. From the strategy point of view there should be a regular review of the key architectures for which the effort should be prioritised, based on the future market trends (driven outside of the science community) and availability of such resources. This has been the case for GPUs for example because of their rapidly increasing availability at many large centres such as HPCs that WLCG could leverage in the future. The activity on software portability should also be supported in terms of infrastructure, by offering solutions for rapid software development and testing cycles on those new architectures. Such infrastructure could be in some cases provided in house at WLCG sites, including CERN. It could also leverage collaboration with HPC centres offering access to testbeds, with the advantage of providing a system close to the hardware that might be available in future. Finally, resources could be purchased from cloud providers if this is the most effective solution. A non person-intensive software validation process (including physics validation) with a fast turnaround time will be very important. Any form of long-term portability requires such a facility as an essential ingredient. We should foresee adequate investments in developing further the existing frameworks for software validation and improving them.

Networking

The WLCG network infrastructure was identified as a limiting factor before the start of LHC data taking. Instead, it proved to be a very reliable service in Run 1 and Run 2, with an adequate level of capacity and considerable headroom. The annual growth of such capacity has been consistent over the last 20 years and its cost is less exposed to market trends than other types of hardware. The LHCOPN and later LHCONE initiatives contributed to organising the network infrastructure and operations for WLCG, in collaboration with the R&E national

network providers and consortia such as GEANT. For those reasons, the experiment computing models and analysis models became more and more network centric and this trend continues toward HL-LHC. Leveraging network connectivity to reduce the needs of storage and reduce the operational cost is part of the strategy.

The vast majority of WAN connectivity is today used for scheduled traffic, replicating data between storage at the sites. The WAN traffic serving data from storage to CPUs at a different sites today accounts for around 10% of the total. This number is likely to increase in HL-LHC due to a more distributed storage model, as described in the DOMA document, while we foresee most of the WAN traffic still being in the form of an organised activity. The use of caches would be part of the reason. Scheduled traffic is managed using the FTS service (except for ALICE) and triggered by the experiments' data management systems. The outlook for HL-LHC, where ATLAS and CMS will be the main network consumers, foresees FTS and Rucio as the natural services to interact with the network for monitoring and shaping purposes.

The DOMA project discussed the possible network needs for HL-LHC in terms of capacity, collecting input from the experiments. More details can be found in the DOMA document. As stated above, we foresee a smaller number of larger sites providing storage. We expect those sites to require high bandwidth connectivity (1 Tbps) between themselves. Some storage endpoints will be regionally distributed and also in this case sufficient internal connectivity would be needed. Computing sites accessing data from the storage backbone would leverage caching techniques and therefore would require a lower level of connectivity, by at least a factor 10 or so. In particular, user analysis use-cases should have less impact on the network if the trimming of the data format for user analysis (as explained in the ATLAS and CMS documents) enables the data set to be permanently cached at a site.

Hiding the effects of network latency will be an important aspect as some computing capacity might need to process data from distant sources, and this is taken into account in the caching and I/O technology studies summarised in the DOMA document. Analyzing network and cache use and effectiveness for HEP workflows with respect of cache sizes is also a major aspect of the future DOMA work. It will quantify the cost-optimization by estimating both disk and network costs.

The risk of possible bottlenecks has been also identified and those should be studied and addressed in preparation for HL-LHC. The transatlantic link is a concrete example: the LHC experiments will produce around 1 EB of RAW data per year at HL-LHC. About 30% of such data will be stored in the US. There is a concrete need to do this as quickly as possible, in quasi-real time with data taking, to secure the data and reduce the pressure on the Tier 0 resources. That might require peaks of bandwidth above the aforementioned 1 Tbps. In the same way, a large computing allocation, at an HPC site or cloud provider could become available for a short period of time and enough connectivity will be required to feed such resources with data from the storage backbone. Validating those capabilities is part of the plans of the DOMA project in the next few years, as described in its document.

We plan to work on improving network monitoring as issues in network services remain difficult to identify. In particular, we are progressing with implementing a service to collect

monitoring metrics from various sources such as network elements, high level service and network probes and organising them in a data analytic infrastructure. For the future it will be important to progress on making our network use visible, by marking traffic by experiment or activity. This will allow in a second stage to shape WAN data flows and finally orchestrate the network to enable multi-site infrastructures. The R&D activities launched for this purpose are highlighted in the DOMA document.

Other HEP experiments and other sciences, in particular astronomy and astroparticle physics, will compete for R&E network resources in the next decade. Particularly for HEP experiments, there is a question about how to best provide networking following the LHCONE model. The LHCONE overlay network has been used to accommodate part of the traffic of DUNE and Belle-2 and experts are evaluating if this is the best option for the future. In particular, the effort of marking traffic mentioned above will allow us to understand if a further segregation of the traffic is needed, for example creating separate overlay networks for other experiments and sciences. In this process we need to ensure the experiments, NRENs and sites collaborate on a simple solution.

Finally, the relatively small LHCONE community has grown considerably since its conception, with the gradual inclusion of more sites and finally experiments, while the original AUP limited its use to sites providing WLCG resources. This raises security concerns loosening the level of trust in the community, which is one of the main benefits. There is an ongoing initiative to strengthen such trust and the security monitoring, while understanding a model for the future.

Authentication Authorization and Identity (AAI)

Our strategy in terms of Authorisation, Authentication and management consists in shifting towards federated identities and adopting new authorization standards from industry. Such a strategy allows for the modernisation of our infrastructure for long term sustainability and an increased level of security, as it will be based on more modern technologies. It also favours the integration of modern services, based on the same standards and it improves user experience. It is also necessary to continue to connect with users globally as well as peer organisation, infrastructures and cloud services.

The needs of the WLCG sites and experiments have been collected, analysed and documented⁷ and a review of the main available authorization architectures was conducted⁸, in order to prepare a transition of WLCG towards a sustainable and highly interoperable authorization infrastructure. It is clear that WLCG has to evolve away from X509, at least for end users, and implementing this strategy has been ongoing for the last few years. Multiple solutions are being developed by communities within the Research and Education sector and a number of translation services will be required to allow interoperable services.

⁷ https://twiki.cern.ch/twiki/pub/LCG/WLCGAuthorizationWG/WLCG_Authorisation_Requirements.pdf

⁸ <https://zenodo.org/record/3460258>

The WLCG AAI Working Group, in collaboration with the AARC⁹, EOSC Pilot¹⁰ and EOSC-Hub¹¹ EU projects, concluded the following part of the work plan:

1. Collect and agree on a well-defined set of requirements from LHC experiments and WLCG sites regarding VO Membership Management and WLCG Service Authorization. These requirements must support and be consistent with existing security policies, operational security requirements, IGTF Levels of Assurance and the EU General Data Protection Regulation.
2. Review the current AAI (Authentication and Authorisation Infrastructure) and the tools being considered for the future by WLCG partners. Evaluate existing or proposed AAI in the HEP community (e.g. in EGI, INDIGO-Datacloud, OSG) for their suitability for WLCG. Review VO management tools (group management) and evaluate how the current VO registration and user management workflow can be expanded to accommodate federated identities. Analyse the aspects related to user authentication, service authentication and authorization, membership management tools, token translation services and suspensions mechanisms.
3. Propose a design for WLCG that ensures both a suitable production service and maximum interoperability in the long term. The costs of preparing and maintaining the authorization infrastructure, the security model, the compliance with existing data protection and conformance to the WLCG requirements defined previously are all key aspects. The scope of the proposal should include the additional services required, such as token translation services or blocking services.
4. Coordinate the definition of a JSON Web Token schema, the building block for token-based authorization solutions such as OpenID Connect and OAuth2, for a common or compatible authorization token profile to be used by collaborating infrastructures.

We foresee many years of coexistence between the old and new AAI methods. The future plan in preparation for HL-LHC consists therefore of a gradual integration of the new AAI tools and services in the current production infrastructure. Its implementation is tracked by the WLCG Authorization Working Group¹², with time frames defined in accordance with input from the wider community. The plan will require development at the level of services and tools, to implement token based authentication and authorisation following the agreed JWT profiles. The interoperability of different implementations will need testing and attention. The most pressing aspect of the migration to the new AAI components is for services relying on Globus GSI libraries as the support for those currently is unclear beyond the end of 2021. The main component in this respect is gridFTP, the retirement of which is a cornerstone of DOMA future planning. See the DOMA document for more information. We foresee to deploy a testbed, which might include production endpoints, where experts, and in future, end users will be able to perform a more and more inclusive set of testing and commissioning activities for the new model.

⁹ <https://aarc-project.eu/>

¹⁰ <https://eoscpilot.eu/>

¹¹ <https://www.eosc-hub.eu/>

¹² <https://twiki.cern.ch/twiki/bin/view/LCG/WLCGAuthorizationWG>

Security & Security Coordination

The threat landscape for WLCG has significantly evolved in recent years, as WLCG adversaries became increasingly sophisticated and well-funded. This evolution affects not only WLCG, but the whole community. In 2019, for example, Verizon found that “23% of bad actors are identified as nation-state or state affiliated”¹³. The strategy for WLCG to manage these sophisticated and persistent threats is to leverage the community itself. This involves two aspects firstly, collecting threat intelligence information and identifying precisely how the malicious actors are currently attacking the victims; secondly, enabling WLCG sites to act on this information within minutes, collecting local system and network logs, and correlating the resulting data with the threat intelligence information obtained above

Collecting threat intelligence is typically done by mature organisations or security trust groups. CERN is, for instance, sharing and collecting threat intelligence from hundreds of partner organisations and making indicators of compromise available to all interested WLCG sites via the industry-standard platform MISP¹⁴. This said, threat intelligence collection must be done both at the global level (to defend against international groups) and at the local or national level (to manage TLD-based or localised threats). It is essential to maintain and establish further strategic collaborations with law enforcement, industry partners and academic/research peers to share and collect relevant, targeted, quality threat intelligence. That is the foundation of the WLCG security strategy. These collaborations are also often leveraged during incident response phases, and our allies provide invaluable information, tools, and expertise, which would all otherwise be impossible to source or afford, regardless of the financial aspects.

A challenge for the coming years is to make better use of MISP and integrate it within the normal incident response workflow in WLCG. This involves reaching a sufficient adoption rate at the sites, which is a key goal of the WLCG Security Operations Centre Working group (SOC WG). MISP alone is not sufficient: in fact, a number of MISP users in WLCG currently do not leverage the indicators of compromise available to them.

One significant barrier to overcome in the coming years continues to be the gap between the “grid computing” and “campus security teams”. The resulting lack of cooperation and coordination at this level is a major contributing factor to poor or failed response to security incidents. Using the “grid” as a “Worldwide collaboration” to source quality threat intelligence, is an important asset to encourage campuses to recognise the value of and increase their trust in their “grid computing” teams.

Another priority in the strategy is to provide the WLCG sites with a reference SOC implementation and assist them collecting local system and network logging information, and correlate it with the indicators of compromise provided to them. This is a difficult task; in the private sector, a growing number of organisations outsource their SOC to security companies.

¹³ <https://enterprise.verizon.com/resources/reports/dbir/>

¹⁴ <https://www.misp-project.org>

One goal is also to achieve a degree of maturity sufficient to have a significant number of WLCG participants to produce and share their own threat intelligence, based on their observations. This would enable a community-based response to sophisticated threats, and be the best sustainable line of defence for all WLCG participants. In terms of collaborations, it is crucial to reinforce the collaboration between US and EU WLCG sites.

Operational Model and Service Provisioning

The effort needed for computing operations, both from the sites and the experiments is a reason of concern as we move toward HL-LHC. Our strategy in reducing operational complexity and cost focuses primarily in reducing the needs for managed storage and exploiting commonalities across experiments to use the same services as much as possible. Those aspects of the strategy have been elaborated already in the first part of this document and in the DOMA document. In addition to that, we identified other aspects that would simplify the operational model.

The hardware capacity increases by 5-10% every year with constant budget investment. Some services become therefore more complex to operate. The experiments also need to manage more data and more jobs per day and manual interventions become very time consuming assuming a constant failure rate. The available operational effort does not increase at either the sites or in the experiments. On the contrary, there is a risk of losing expertise if the focus on repetitive tasks increases further. Modernising the way operations are run in WLCG is the way to proceed. At several times already in the past the number of hardware resources was scaled up by several factors not by increasing the staffing, but rather by leveraging modern technologies and organisation of the services. One example was the introduction of a virtualisation layer between the fabric and the compute platform. There are more opportunities coming from open source projects supported by large communities and it is in the interest of WLCG to understand how we can benefit from them. Our strategy is to organise the explorative activities in projects or working groups, so that the know-how is shared across sites and experiments and all parties can contribute with their experience and benefit.

The use of virtualized systems to isolate the application from the hosting environment is being explored in more areas in WLCG. It proves to be very useful at the level of the processing farms in providing user traceability and favoring the migration to recent versions of the operating system. We are actively looking at the possibility to leverage container orchestration systems such as Kubernetes¹⁵ to change the way we deploy services at the site, taking into account the different constraints. A more agile model of service deployment based on container orchestration would allow a rapid turnaround of the deployment cycle for testing and integration activities.

¹⁵ <https://kubernetes.io>

The Operational Intelligence Project (OpInt)¹⁶, a joint effort of HEP collaborations, targets the reduction of operational costs of complex computing infrastructures by increasing the level of automation in operations. This will be achieved through the development of a stack of intelligent tools and technologies aiming to detect, analyse, and predict anomalies of the computing environment, to suggest possible actions, and ultimately automate operation procedures. Examples of such tools are smart alerting systems, recommendation systems for operators, or predictive maintenance through log analysis. The landscape of applicable technologies and methods includes but is not limited to data mining, machine learning, predictive analytics, interactive data visualization, visual analytics, and natural language processing techniques.

In general, modernising the operational tools and techniques has the advantage of attracting younger generations of computing professionals, offering them a better portfolio for future employment perspectives. At the same time, retaining some of this expertise in our field is a very important aspect. The preparation for HL-LHC is a unique opportunity to initiate new generations of specialists through challenges of increasing complexity and building a future core of expertise. Such an opportunity needs the engagement of the funding bodies with adequate resources and long term perspectives.

¹⁶ <https://operational-intelligence.web.cern.ch/>