# Tackling Computing Challenges at CERN

Maria Girone
CERN openlab CTO

LHC
- 50-175m underground
- 27 km circumference tunnel
- Four giant experiments
- Particles travelling at 99.9999991% the speed of light
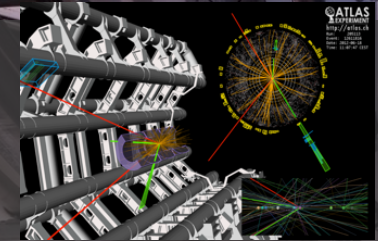- 11245 turns every second

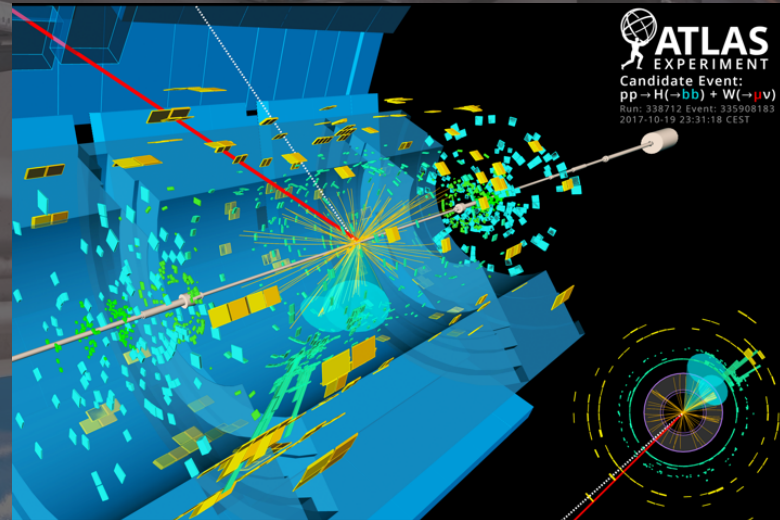The most powerful discovery machine at CERN is the Large Hadron Collider

Maria Girone
CERN openlab CTO

2

- Raw data:
  - Was a detector element hit?
  - How much energy?
  - What time?

- Reconstructed data:
  - Momentum of tracks (4-vectors)
  - Origin
  - Energy in clusters (jets)
  - Particle type
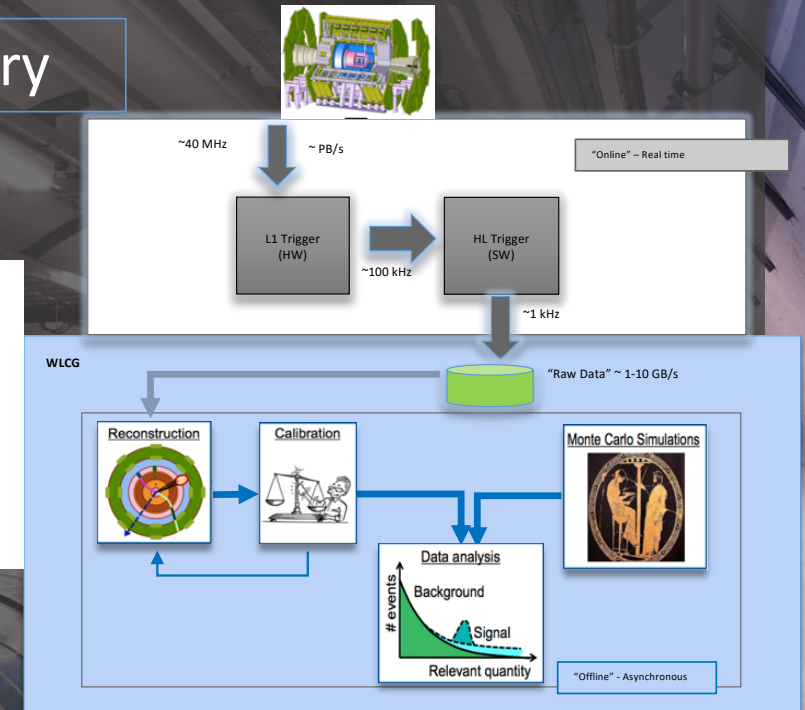  - Calibration information
  - Analysis Objects
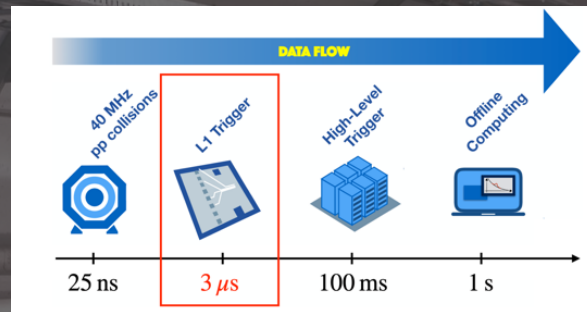  - …

- 150 Million sensors deliver data …
40 Million times per second

- Generates ~ 1 PB per second



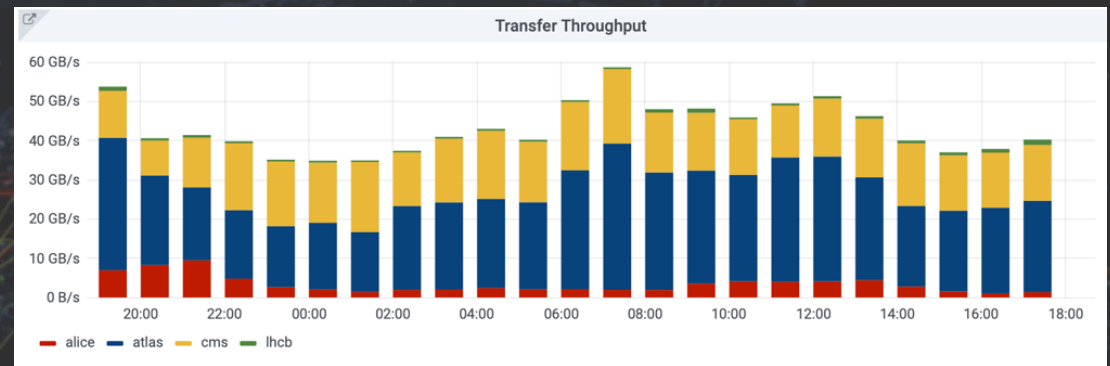# The LHC Data

# Data processing and analysis drive physics discovery

**Software**

- 50M lines dominated by C++ and Python
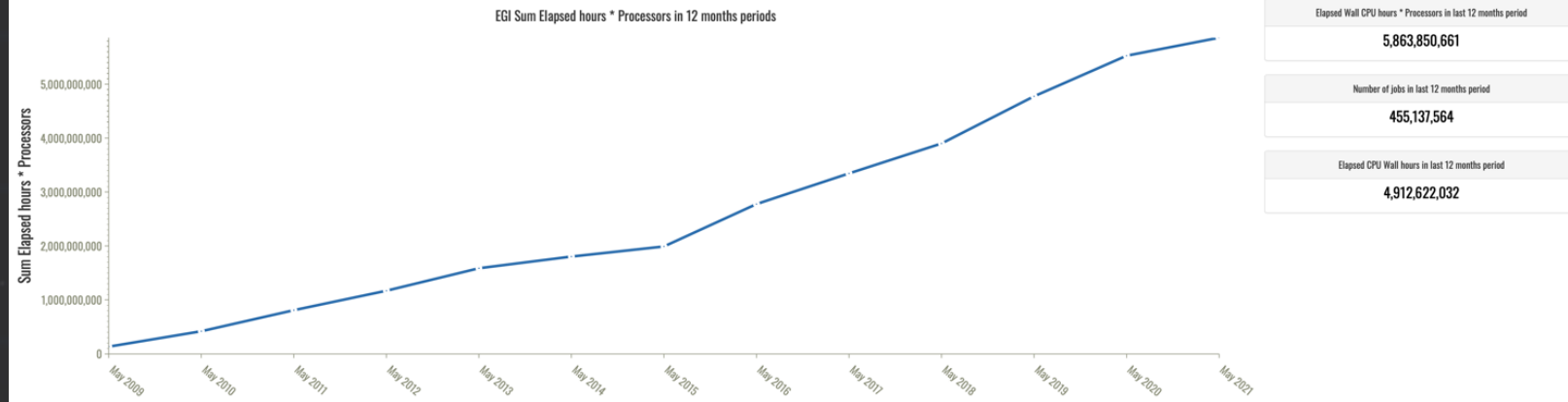- Contributions by hundreds of scientists
- Much 20+ years old



DATA FLOW

40 MHz pp collisions — L1 Trigger — High-Level Trigger — Offline Computing

25 ns — 3 µs — 100 ms — 1 s

~40 MHz — ~ PB/s — "Online" – Real time

L1 Trigger (HW) — ~100 kHz — HL Trigger (SW) — ~1 kHz

WLCG — "Raw Data" ~ 1-10 GB/s

Reconstruction — Calibration — Monte Carlo Simulations

Data analysis
# events
Background
Signal
Relevant quantity

"Offline" - Asynchronous

# Dataflow at LHC

Maria Girone
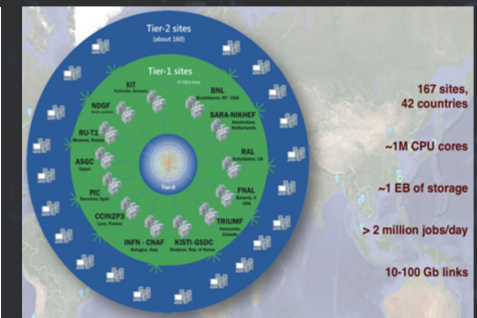CERN openlab CTO

# Hardware

- Primary computing resources is the WLCG, a globally distributed storage and processing infrastructure
- 167 sites over 42 countries
- **~1M CPU cores** and **~1 exabyte** of storage (disk and tape)



Transfer Throughput



This graph shows the Sum Elapsed hours * Processors in the whole EGI infrastructure. Only non-local jobs on official EGI VOs are accounted. Each point represents a period of 12 months counting backwards from the last complete period.
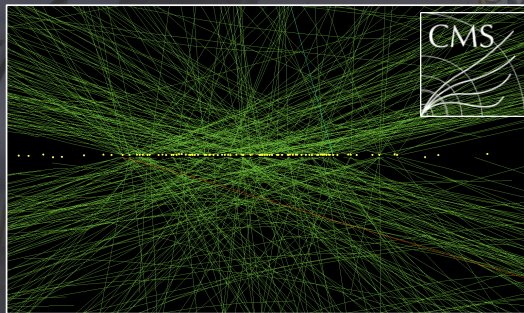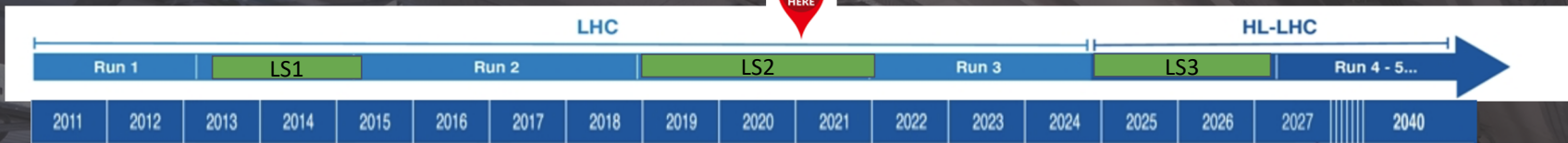
EGI Sum Elapsed hours * Processors in 12 months periods

| Elapsed Wall CPU hours * Processors in last 12 months period |
| --- |
| 5,863,850,661 |

| Number of jobs in last 12 months period |
| --- |
| 455,137,564 |

| Elapsed CPU Wall hours in last 12 months period |
| --- |
| 4,912,622,032 |

167 sites, 42 countries

~1M CPU cores

~1 EB of storage

> 2 million jobs/day

10-100 Gb links
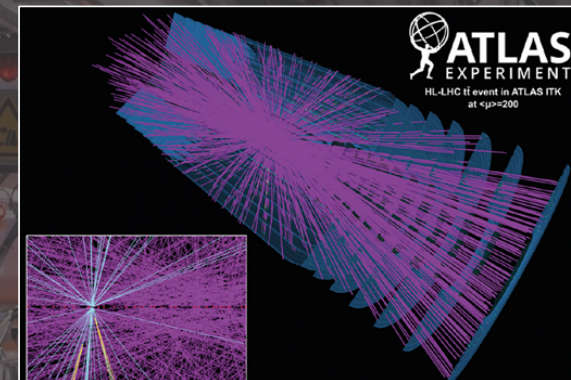
# The Worldwide LHC Computing Grid

Maria Girone
CERN openlab CTO

# The Challenges of HL-LHC

LHCb and ALICE will be upgraded for Run3 and will collect much more data

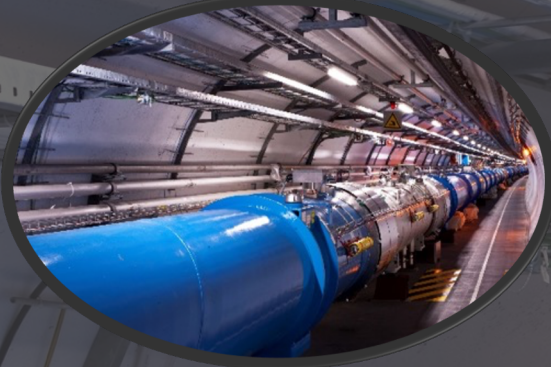The LHC will be upgraded for Run4 to the High Luminosity LHC (HL-LHC). This will deliver:
- x10 increase in luminosity over LHC design
- great increase in event complexity
- more collisions and more complex data will result in a compute challenge at the Exascale level

Run2 – Average 40 collisions per crossing

Run4 – Average 200 collisions per crossing

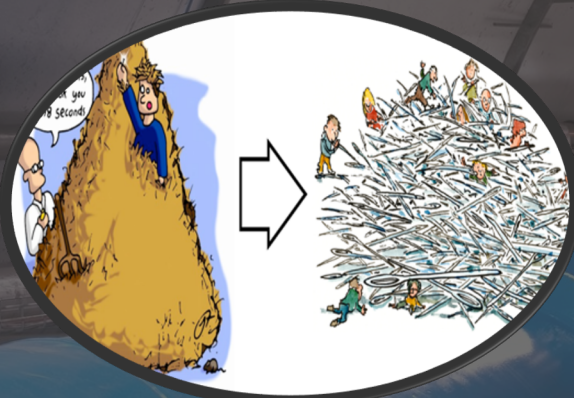# Scheduled Upgrades of LHC Program

Maria Girone
CERN openlab CTO

Upgraded Accelerator
- Higher Luminosity

Upgraded Detectors
- Higher Granularity
- Higher Occupancy

Changing Filtering Paradigms
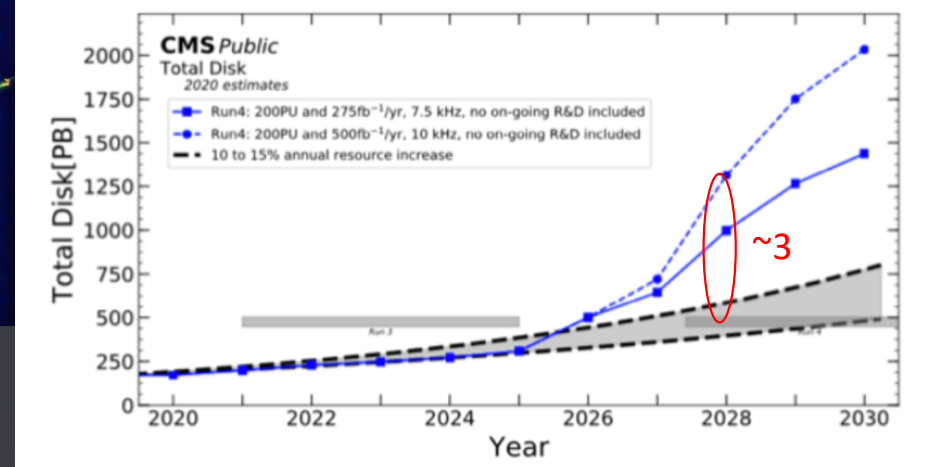- Higher Sensitivity
- Higher Data Rates
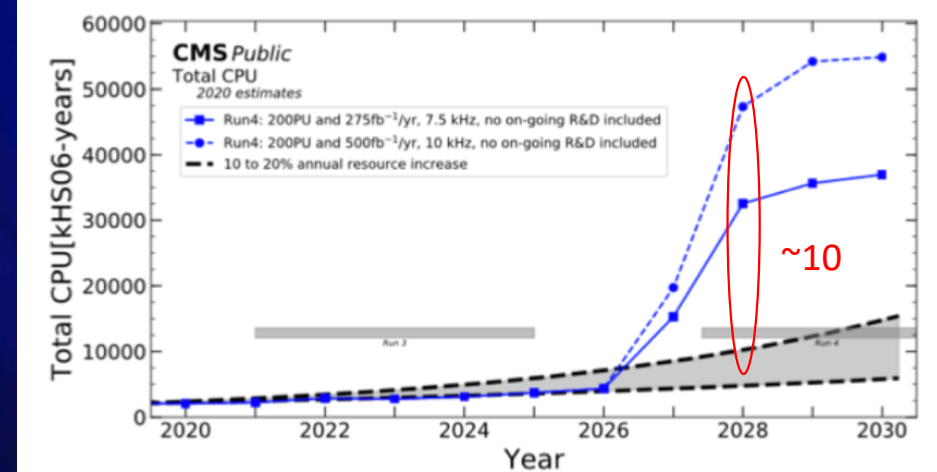
New Computing Challenges

# Upgraded Program = New Challenges

Gap between needed/available computing resources

- CMS estimate computing power needs to be ~6-10x higher

- CMS estimate disk needs to be ~3-5x times larger

- ATLAS estimates are similar

Investments for R&D in

- Code modernization and optimization

- Adapting code to hardware accelerators and HPC

- Reducing storage needs

- New techniques like AI and ML

# Computing Challenges

General purpose CPU performance increases have slowed

Optimized heterogenous architectures have evolved faster, **HEP Is investing heavily in development to use new hardware resources**

- **GPUs** are the most common
- **FPGAs** currently used mostly in low latency applications
- **TPUs** and specialized ASICs are available

| | | Accelerated Reconstruction And AI/ML | Accelerator | | Low latency Online applications | |
| --- | --- | --- | --- | --- | --- | --- |
| | | Intel | NVidia | AMD | FPGA | Other |
| CPU | Intel | Aurora | Cori Piz Daint Tsukuba Mare Nostrum | | Tsukuba | |
| | AMD | | Perlmutter JUWELS Booster | Frontier El Capitan LUMI | | Amazon Graviton2 Google Cloud TPU Microsoft Azure Intel DevCloud |
| | IBM | | Summit Sierra Mare Nostrum | | | |
| | ARM | | Wombat | | | Astra Fugaku |

General Purpose X86 processing resources

Code ported to Power

Low power highly parallelized

# HEP and the Computing Landscape

https://cacm.acm.org/magazines/2019/2/234352-a-new-golden-age-for-computer-architecture/fulltext

End of the Line ⇒ 2X/20 years (3%/yr)
Amdahl's Law ⇒ 2X/6 years (12%/year)
End of Dennard Scaling ⇒ Multicore 2X/3.5 years (23%/year)
CISC 2X/2.5 years (22%/year)
RISC 2X/1.5 years (52%/year)

Performance vs. VAX11-780

100,000
10,000
1,000
100
10
1

1980  1985  1990  1995  2000  2005  2010  2015

# Four Pillars of Activity

**XT**
**eXascale Technologies**

A comprehensive investigation of HPC and Cloud infrastructures, frameworks, tools to support key scientific workloads and applications

**AI-S**
**Artificial Intelligence for Science**

Analysis and development of algorithms, optimisation for new architectures, interpretability, synergies between Physics and other sciences

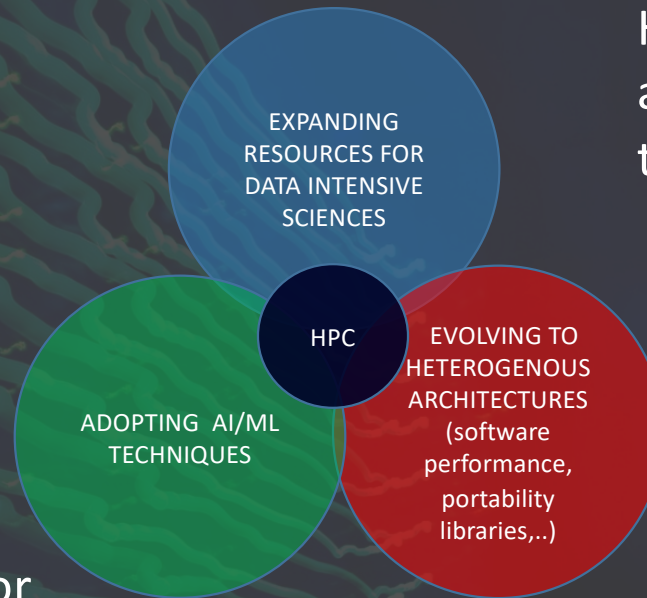**QTI-C**
**Quantum Technology Initiative - Computing**

Assess the potential impact of quantum computing in HEP and other sciences, investigate quantum machine learning algorithms and areas of potential quantum advantage, set up a collaborative quantum computing (simulation) platform

**MSC**
**Multi-Science Collaborations**

Share the expertise and knowledge generated across all activities with other sciences, work with CERN KT to explore novel applications of CERN computing systems and ideas, create collaborations and contribute to common solutions

CERN openlab R&D's: HPC, AI, and QC

HPC falls at the intersection of several important R&D areas

EXPANDING RESOURCES FOR DATA INTENSIVE SCIENCES

HPC

ADOPTING AI/ML TECHNIQUES

EVOLVING TO HETEROGENOUS ARCHITECTURES (software performance, portability libraries,..)

Engagement with the HPC Community can be a catalyst for progress

HPC Supercomputers will grow by a factor of 10 on the time scale of the HL-LHC

A thorough R&D program has been established

Unified programming models facilitate HPC adoption

# High Performance Computing

- **An HPC Collaboration agreement was signed by CERN, SKAO, GÉANT and PRACE, CERN and SKA on 22.07.2020**
  - Engages at the community level
    - Bringing together data intensive sciences, high-performance computing infrastructures and networking

- Collaboration built around 4 pillars
  - Building a common centre with expertise to support heterogenous hardware
  - Benchmarking Demonstrator
  - Data Access Demonstrator
  - Authentication and Authorization Demonstrator

An Exascale project for an Exascale problem



Eckhard Elsen (top left), Director for Research and Computing at CERN; Philip Diamond (top right), SKA Director-General; Erik Huizer (bottom left), Chief Executive Officer of GÉANT; and Philippe Lavocat (bottom right), PRACE Council Vice-Chair, signed the agreement for the new collaboration.

https://home.cern/news/news/computing/cern-skao-geant-and-prace-collaborate-high-performance-computing

# Working Together: The HPC Collaboration



CERN, SKAO, GÉANT and PRACE to collaborate on high-performance computing

The next generation of high-performance computers holds significant promise for both particle physics and astronomy but key challenges remain to be addressed

22 JULY, 2020 | By Andrew Purcell

Participation on the path to exascale

Maria Girone
CERN openlab CTO

14

**Proven CERN capability** ✓

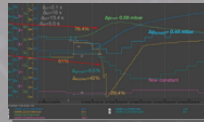Use case specific

**Fast ML**

Ultra-fast on-edge inference under strict latency constraints

**Anomaly detection**

Object identification, classification, anomaly detection in big and noisy data sets
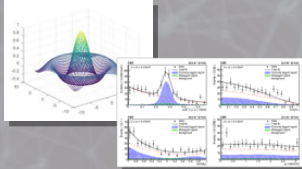
**Industrial controls**

Machine efficiency and predictive maintenance with industrial control systems

**Distributed computing**

Optimization of distributed computing, storage, and networks; fast I/O for large files

**Large scale, science grade data analytics and visualization**

Cross use case

- Optimization and evaluation for science-grade precision of large data sets using advanced data analytics
- Data visualization, interactive plotting (e.g., statistical visualizations, uncertainties, distributions), model visualization
- Large-scale, quality-controlled CERN data as testbed/benchmark (e.g., single data set with 100m examples, >1TB)

**In development, opportunity for joint R&D**

**Simulation**

Simulation and reconstruction with generative DL for efficient computation

**Graphs**

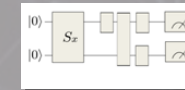Exploring Graph NNs for high-multiplicity problems with non-linear distances

Determining optimal machine design and component configuration
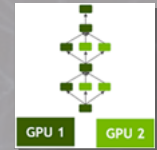
**ML in Robotics**

Remote maintenance and safety with autonomous robots and computer vision

**Quantum ML**

Research quantum algorithms to solve pattern recognition, classification and generation problems
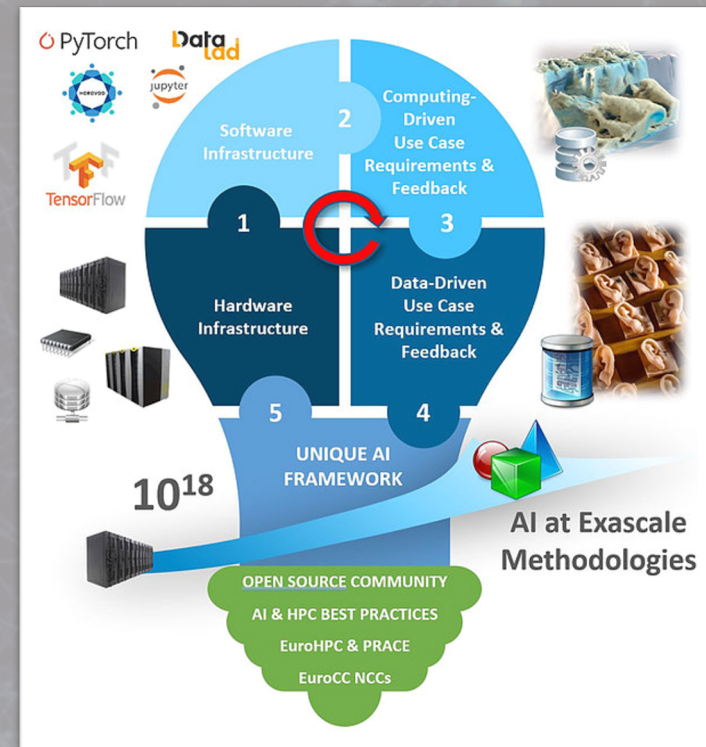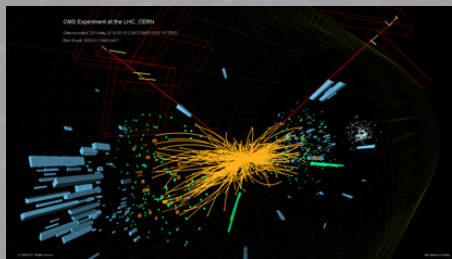
**Computing parallelization**

Training and optimization of complex NNs on parallelized GPU infrastructure

# Progress on AI/ML Capabilities

15

Launched in January the RAISE Center of Excellence enabled researchers from science and industry to develop novel, scalable Artificial Intelligence technologies towards Exascale along representative use-cases from Engineering and Natural Sciences

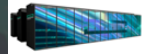CERN is leading the leading the data driven use-cases



# AI/ML Projects

HEP, HPC, and Commercial Clouds

**Challenges**

Software and Architectures      Supercomputers are early adopters of heterogenous architectures

Benchmarking and Accounting      Performance on diverse architectures needs to be understood

Data Processing and Access      Enormous data volumes to stage, process, and export

Authorization and Authentication      Strict cyber security

Runtime Environments and Containers      Resources are shared, environment needs to be brought with the workload

Provisioning      Resources allocated for periods of time through allocations

Wide and Local Area Networking      Processing and storage resources are separate

# Challenges in HPC Integration

The common challenges for HPC integration into LHC Computer were described in an engagement document

https://zenodo.org/record/3647548#.YBnA1y2cbVs

Develop an energy-efficient system architecture that fits HPC and HPDA workloads





Build a fully working Modular Supercomputer Architecture prototype to Exascale



Large variety of hardware available supporting the different requirements of HPC, Big Data Analytics and Machine Learning with highest efficiency and scalability

# Optimising HEP applications towards Exascale

Within the DEEP-EST project, re-engineered CMS ECAL and HCAL local reconstruction workloads to use GPUs using CUDA



- Achieving between 3x(ECAL) and 8x(HCAL) using Nvidia V100 vs filling in 2-socket Intel Xeon Gold 6148

- Now integrated with CMS framework and will be used in the CMS HLT reconstruction for Run3

Performance studies on the HEP MC generator code on GPUs using CUDA

- The idea is to enable utilization of heterogenous architectures for MC generation as well

We are investigating **unified programming models** to create sustainable code that can be supported on multiple architectures

Results on Open Data:
http://opendata.cern.ch/record/12303

# Progress Using Heterogenous Architectures

# How do we bring large datasets to supercomputers

WLCG *Data Lake* model separates storage and processing functionality.  HPC will be a part of the *Data Lake* model
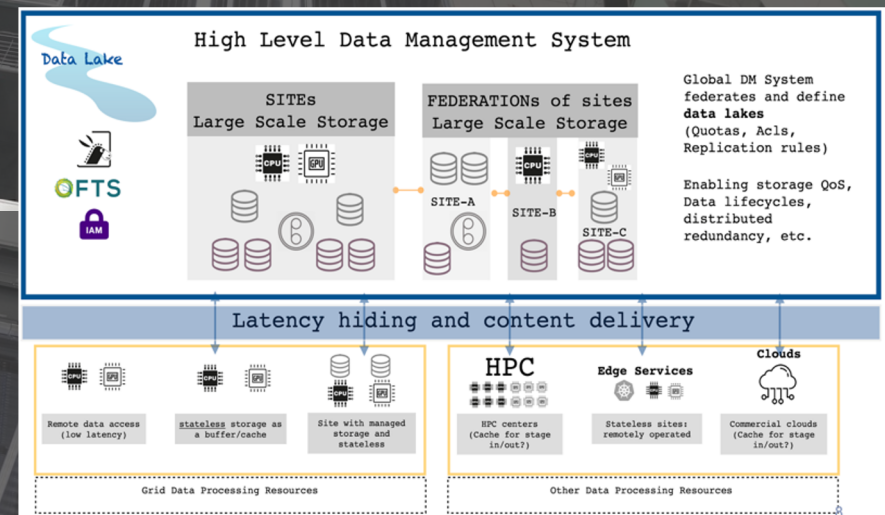
- Relies on caching and networking
- EuroHPC will have significant WAN connectivity and disk space





Emulation of cache delivery vs. time including regional caches

## Technical Activities

Execute a series of data challenges to demonstrate the feasibility of the *Data Lake* model on a path to Exascale

# Data Access – *Data Lakes*

HEP data is primarily stored as files, optimized for highly parallel HTC
- **ROOT** is the HEP analysis framework
- **ROOT** defines columnar data layout tailored for HEP: extreme throughput compared to alternatives
- https://root.cern

**ROOT** Challenges
- Maximize throughput I/O and optimize for HPC
- Optimize **persistent** data layout to facilitate conversion for CPU, GPU, SIMD (LLAMA), read patterns, and storage backend

Ongoing R&D, bringing >4GB/s from off-the-shelf desktop to HPC

ROOT team bringing heterogenous computing and environments to physicists

- **Declarative programming:** physicists define data + analysis flow; "kernel graph" built behind the scene with runtime-detected input data types
- **Transparent acceleration:** algorithms (modelling / minimization) in multi-arch libraries selected at runtime, covering architectures' SIMD to GPU
  - Optimal abstraction? Autovec + CUDA / std::simd / alpaka SYCL
- **Enabling feature:** C++ just-in-time compilation (cling); also supports runtime-CUDA. Use of C++ automatic differentiation (clad) for minimization
- **Scaling:** multi-threaded (>200 cores), distributed backends (dask / spark /...)
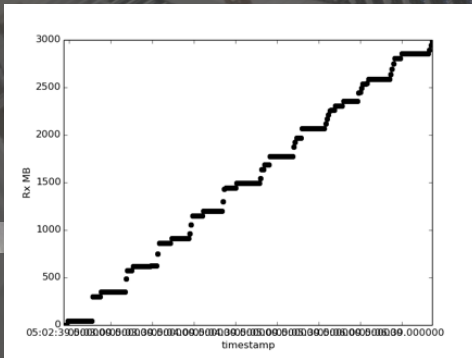
Courtesy of A. Naumann

# HEP Data in HPC

**ROOT**
Data Analysis Framework

Maria Girone
CERN openlab CTO

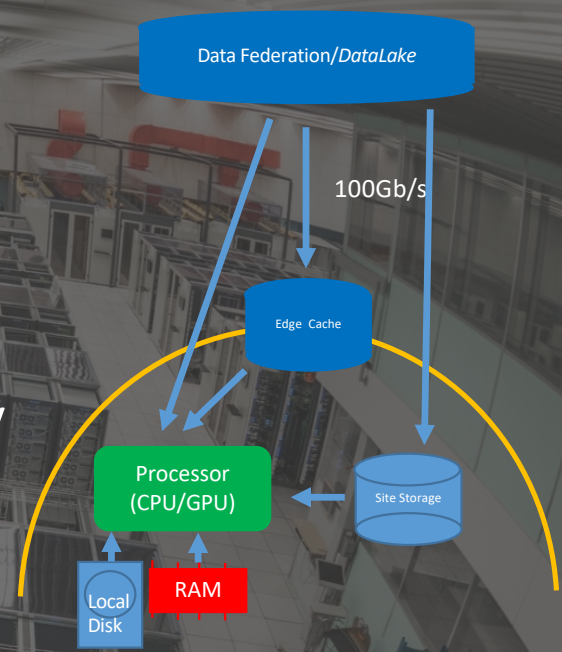Establishing a data access framework
- Multiple containerized applications with different IO profiles
- An aggregator to analyse, study, and optimize facility and workflow deployments

Studying workflow performance progressively farther from the processing resources

- Goal is a series of data processing challenges that will eventually reach 10PB a day

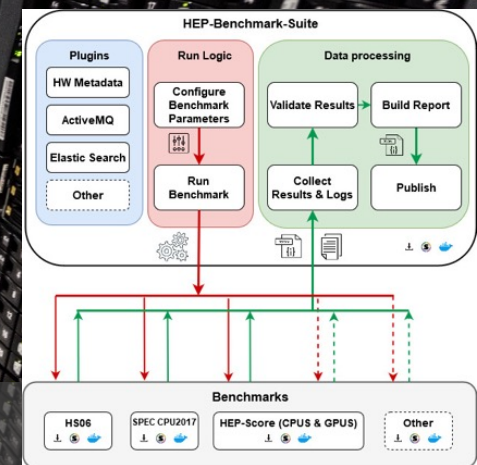For Exascale HPC O(1M) cores
-> O(10K) nodes -> O(150) GB/s

Data Federation/*DataLake*

100Gb/s

Edge Cache

Processor (CPU/GPU)

Site Storage

Local Disk

RAM

Courtesy of V. Khristenko

# Studying Data Access

Benchmarking Activities

- PRACE-CERN-GÉANT-SKAO collaboration brings opportunity to expand capabilities using tools already developed for HPC sites by each community:

  - Unified European Applications Benchmark Suite (UEABS)- 13 workloads for HPC

- CERN is evolving the approach to benchmarking in HEP to embrace HPC:

  - Builds on experience from WLCG computing environment tools

  - Developed with secure, self-contained workload images (Singularity)

  - Assumes no privileges, no docker, limited/restricted node connectivity

# Benchmarking Demonstrator



PRACE
Unified European Applications Benchmarking Suite



HEP-Benchmark-Suite

Benchmarking Heterogenous architectures
- Multi-architecture as workflows become available (ARM, IBM Power)
- GPU accelerators (NVIDIA, AMD)

Automated collection and aggregation



https://gitlab.cern.ch/hep-benchmarks

# HEP Workflows on HPC

Courtesy of D. Southwick

IceCube with OSG and UCSD have been producing simulation with GPUs using the major commercial cloud providers

- Over the 9h test reached 1 exaflop/hour delivered
  - 150PF/h



# Simulation on Commercial Clouds

Demonstration of cloud analysis access was shown on CMS open data during CHEP, Adelaide 2019

- Analyzing 70TB of data and generates the Higgs discovery plot in about 20 minutes



**70 TB** of Physics Data    **~25000** Files



70 TB Dataset → Cluster on GKE → Job Results → Interactive Visualization

Max **25000 Cores**

Single Region, 3 Zones          Aggregation

25000 **Kubernetes Jobs**

# Cloud Use for Analysis

https://www.youtube.com/watch?v=CTfp2woVEkA

# Experiment Activities in HPC and Clouds

Maria Girone
CERN openlab CTO

**ATLAS** Computing Model (CM) designed to use distributed computing centers. CM based on three main pillars : Data Management (Rucio), Workload Management (PanDA) and monitoring

- ATLAS decided to add HPCs  and to integrate High Throughput Computing (Grid) with HPC.  HPCs integrated into the production, analysis and data management systems (also to monitoring and accounting) in 2016

- Over the past 7 years, the ATLAS experiment collaborated with many large HPC sites for full integration into ATLAS distributed computing

**CMS** continuously invests effort to build up expertise on HPC resources integration
- Given the unicity of the HPC Machines multiple approaches have been successfully commissioned:
  - HEPCloud: US-CMS gateway to provide access for CMS to US HPC

  - Site extension:  mechanism exploited at CINECA Marconi A2 and ForHLR2(KIT) and CLAIX(RWTH Aachen)

- Working on the exploitation of CPU resources at HPC centres where compute nodes do not have external network connectivity.

# Experiment Use of HPC

# ATLAS CPU Resource Mix

- HPCs delivered 9% of ATLAS normalized wallclock usage in the past year
- A large number of HPCs contributed worldwide. Clear demonstration that we can integrate diverse mix of HPC systems to enable LHC physics
  - *Incomplete list of integrated centers : CSCS, MareNostrum, OLCF, ALCF, NERSC, TACC, LRZ, Nordic HPCs, RU NRC KI, iT4I, ....*
- ATLAS focused on enabling all HPC centers available to ATLAS into the distributed computing system
- New opportunities with EuroHPC project :
  - New ways how physics analysis and data processing will be done



HPC 9%
Cloud 18%
Grid 70%

ATLAS jobs normalized wallclock consumption. All resources
Jan 2020 - Jan 2021

*ATLAS weekly CPU consumption in 2018 (LHC Run2)*

- *HPC delivered 700 million CPU wallclock hours for Monte-Carlo simulations*

ATLAS jobs normalized wallclock consumption at Tier-1 in Barcelona
Jan 2020 - Jan 2021



Grid 36%
Mare Nostrum 61%



CPU usage 2018
All sites
Grid
Cloud
Special cloud
HPC
Special HPC
500k
300k
100k
Feb  Apr  Jun  Aug  Oct  Dec

Courtesy of D. Benjamin A. Filipic, A.Klimentov

- ATLAS CH Tier 2 is integrated with Piz Daint
  - First use of HPC as a WLCG Tier 2 Center
  - Using ARC-CE and ARC cache

- MareNostrum
  - Served by ARC-CEs located in Spanish Tier 1/2 WLCG centers)
  - Using singularity container (no internet from WN) for ATLAS SW and databases (O(100GB))



ATLAS WLCG T2 center integration with CSCS

Software developed to address HPC challenges :

ATLAS software in containers, granular data processing, seamless integration with grid, preemption, backfill mode, ...

Courtesy of D. Benjamin A. Filipic, A.Klimentov

# EU-HPC Integration with ATLAS Computing

CMS has demonstrated co-located HPC access as well as the use of HEPCloud for accessing HPC and Commercial Clouds

**HPC resources in HEPCloud in 2020**

Top utilization has been in excess of 100k cores

300k Cores

80k Cores

Via Fermilab HEPCloud:
CMS Amazon Web Services (AWS) Usage

Fermilab Tier-1

Courtesy of D. Spiga, C. Wissing

**KIT Site extension**

CMS Resource Usage

HEP is facing an unprecedented computing challenge from the Exabytes of data expected from the HL-LHC

- We have successfully operated our distributed computing environment, the WLCG, for more than a decade and exploited globally distributed computing resources to realize the scientific potential of our data

Looking forward, we will need more resources and opportunities with HPC/commercial clouds may play an important role

- We are involved in projects to exploit exascale capabilities for data-intensive science
  - Continuing explorations of heterogenous architectures
  - Building the expertise in AI/ML
  - Increasing scale of processing and data access solutions
  - Liaising to technology providers through CERN openlab



HORIZON 2020

RAISE
Center of Excellence

HPC Collaboration

CERN GÉANT
PRACE SKA
SQUARE KILOMETRE ARRAY

egi-ace

CERN openlab

We are working with to establish enablers for data intensive science using HPC and commercial cloud resources

# Outlook

Maria Girone
CERN openlab CTO