# The HPC Collaboration:
## CERN, SKAO, GÉANT PRACE

As you learned in the talk of Joseph Flix on Tuesday, the needs for computing in data intensive science are increasing dramatically

- The HL-LHC and SKAO face unprecedent computing challenges over the next 5-10 years

The needs for new resources are driving an extensive R&D effort

- Heterogenous hardware and High-Performance Computing (HPC)
  - There are big national investments in HPC
  - Improvements in heterogenous hardware are driving increases in capacity
- New methods and opportunities in AI/ML

# New Computing Challenges

Upgraded Accelerator
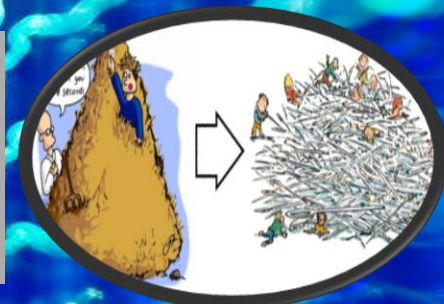- Higher Luminosity

Upgraded Detectors
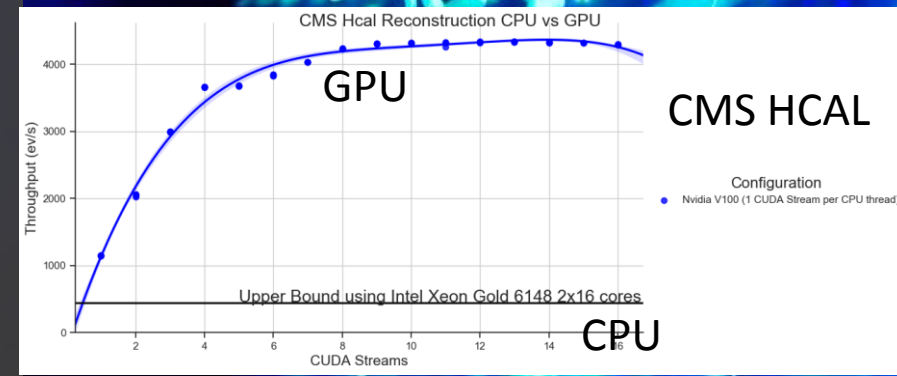- Higher Granularity
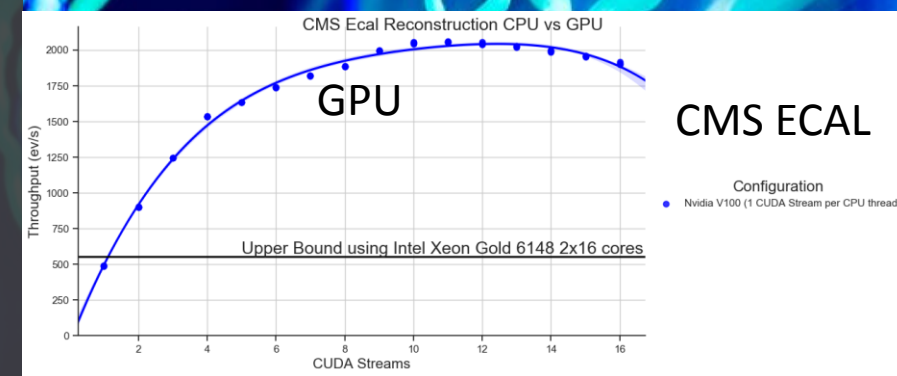- Higher Occupancy

Changing Filtering Paradigms
- Higher Sensitivity
- Higher Data Rates

New Computing Challenges

CERN openlab CTO

- Explorations in heterogeneous hardware are one of the drivers of innovation in computing
  - Large improvements recently in processing performance have come from offloading work to GPUs
    - E.g. CMS HLT (**ECAL, HCAL)** with GPUs gives ~factor 10 in improvement in throughput (on NVIDIA V100)

- Reengineering the software is necessary
  - Results in performance gains, easier adoption of heterogeneous hardware and better maintainability

- Experiments are encouraged to make use of supercomputers (growing to exascale)
  - Common across US, EU, and Asia
  - HPC sites are early adopters of new technologies and sources of expertise

**CMS ECAL**

CMS Ecal Reconstruction CPU vs GPU

GPU

Configuration
Nvidia V100 (1 CUDA Stream per CPU thread)

Upper Bound using Intel Xeon Gold 6148 2x16 cores

Throughput (ev/s)

CUDA Streams

**CMS HCAL**

CMS Hcal Reconstruction CPU vs GPU

GPU

Configuration
Nvidia V100 (1 CUDA Stream per CPU thread)

Upper Bound using Intel Xeon Gold 6148 2x16 cores

CPU

Throughput (ev/s)

CUDA Streams

# Motivation for an HPC Collaboration

Maria Girone
CERN openlab CTO

# Changing Computing Landscape

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

General Purpose X86 processing resources

Code ported to Power

Low power highly parallelized

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

Amazon
Graviton2
Google Cloud TPU
Microsoft Azure
Intel DevCloud

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

4

The current landscape is pushing HEP and other sciences to integrate HPC resources

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

Engagement with the HPC Community can be a catalyst for progress

EXPANDING RESOURCES FOR DATA INTENSIVE SCIENCES

HPC

ADOPTING AI/ML TECHNIQUES

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

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

# Challenges

Software and Architectures

Benchmarking and Accounting

 Data Processing and Access

Authorization and Authentication

Runtime Environments and Containers

Provisioning

Wide and Local Area Networking

Supercomputers are early adopters of heterogenous architectures

Performance on diverse architectures needs to be understood

Enormous data volumes to stage, process, and export

Strict cyber security

Resources are shared, environment needs to be brought with the workload

Resources allocated for periods of time through allocations

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

As we adapt

- Our consortium is ideally composed
  - HL-LHC and SKA have a burning physics need and in depth knowledge of the algorithms employed
  - PRACE provide considerable experience in the system adaptation of software environments
  - GEANT provides the infrastructure to take the computing to the many nodes that are needed to tackle the demand

**PRACE | Tier-0 Systems in 2020**

**MareNostrum**: IBM BSC, Barcelona, Spain #38 Top 500

**Piz Daint**: Cray XC50 CSCS, Lugano, Switzerland #10 Top 500

**NEW ENTRY 2018/2019 SuperMUC NG** : Lenovo cluster GAUSS @ LRZ, Garching, Germany #13 Top 500

**NEW ENTRY 2018 JUWELS (Module 1):** Atos/Bull Sequana GAUSS @ FZJ, Jülich, Germany #39 Top 500

©FZ Jülich / R.-U. Limbach

**NEW ENTRY 2018 JOLIOT CURIE** : Atos/Bull Sequana X1000; GENCI @ CEA, Bruyères-le-Châtel, France #34 Top 500

**MARCONI-100: IBM** CINECA, Bologna, Italy #9 Top 500

**NEW ENTRY 2020 HAWK:** HPE Apollo GAUSS @ HLRS, Stuttgart, Germany

**Close to 110 Petaflops total peak performance**

5   The Partnership for Advanced Computing in Europe | PRACE

From the HPC Collaboration Kick-off-Workshop

Signature Cerimony

# CERN, SKAO, GÉANT, PRACE Consortium

- Four areas of work have been identified as foundational. Progress will be evaluated by a series of common challenges and demonstrators
  - **Benchmarking**
  - **Data Access**
  - **Authentication and Authorization**
  - **Building a Common Center of Expertise**

- The next crucial step is to address the challenges through a **common program of work**
  - The **roadmap** is outlined on the next slides
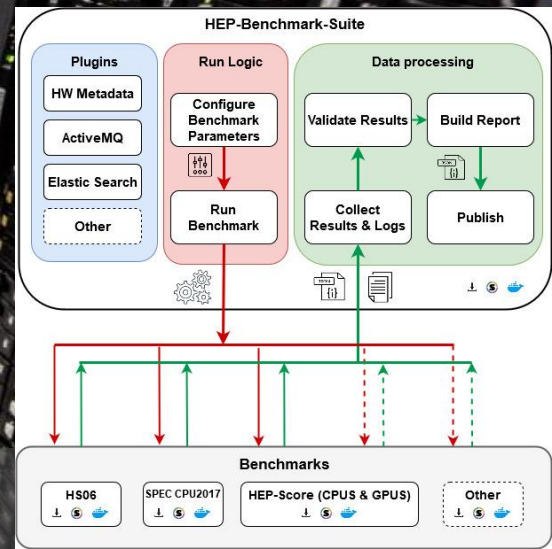
# The Four Pillars of the Collaboration

Maria Girone
CERN openlab CTO

Benchmarking Activities

**Open Question:** How do we show we can effectively use HPC systems and account for contributions?

- 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

# HEP Workflows on HPC

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

Courtesy of D. Southwick

**Open Questions:** How do we bring data to process on supercomputers?. Can we use HPC sites to produce simulation and reconstruction datasets at exascale?
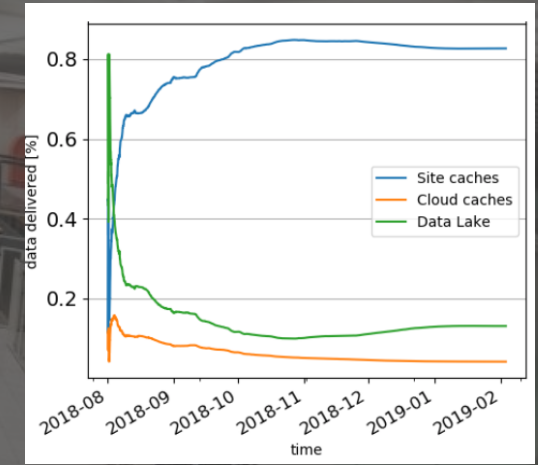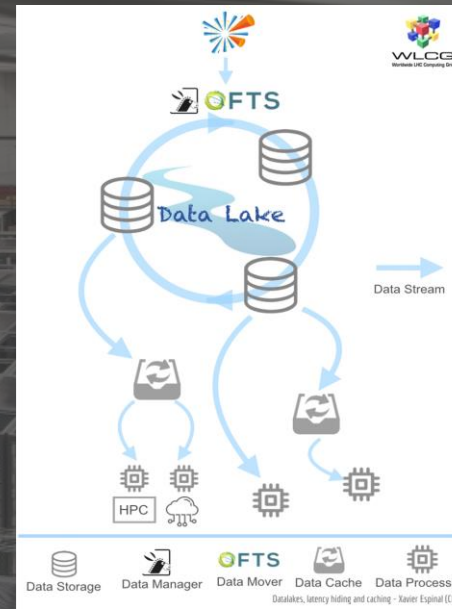
WLCG *Data Lake* model separates storage and processing functionality. HPC will be a part of the *Data Lake* model

- Relies on caching and networking

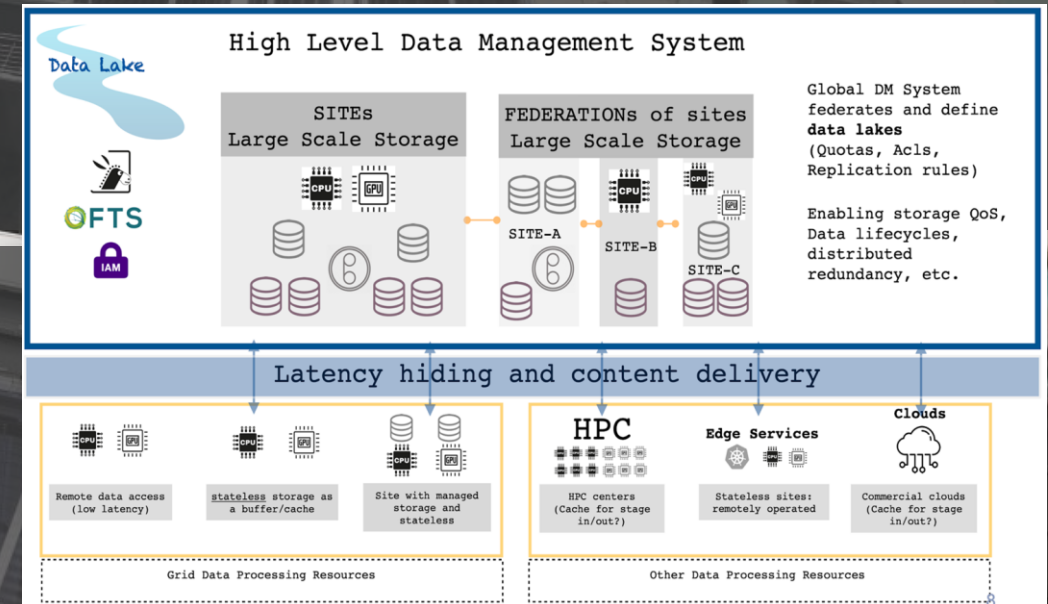- EuroHPC will have significant WAN connectivity and disk space





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*



Maria Girone
CERN openlab CTO

- **With the help of Summer of HPC Students we started an IO Benchmarking project on HPC**
  - Evaluate and rank system performance.
  - Find limitations (more computing power isn't necessary more efficiency).

- **The goal of this project**:
  - Scale out I/O mock-ups using HPC benchmarks .
  - Evaluate and visualize performance metrics (e.g., bandwidth utilization).
  - Report performance under heavy dataflow load.

metric

Courtesy of V. Khristenko

# Studying Data Access

Maria Girone
CERN openlab CTO

# ❑CSCS Grand Tavé[1]

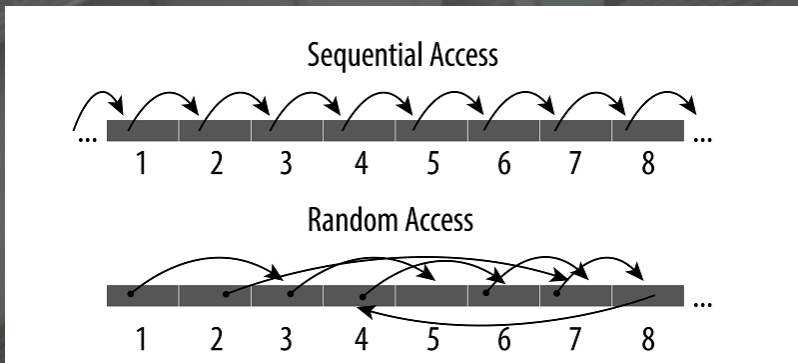| Specifications | |
|---|---|
| Model | Cray XC40 KNL |
| Compute Nodes | 64 cores Intel(R) Xeon Phi @ 1.30GHz |
| Memory Capacity per Compute Node | 96 GB, 16 GB HBM |
| Login Nodes | 8 cores Intel(R) Xeon(R) @ 2.60GHz |
| Memory Capacity per Login Node | 256 GB |
| Theoretical Peak Performance | 436.63 TFlops |
| Max number nodes | 164 |
| Scratch capacity | /scratch/snx2000 904 TB |

The $SCRATCH space (/scratch/snx2000/$USER) is connected via Infiniband interconnect to the system.
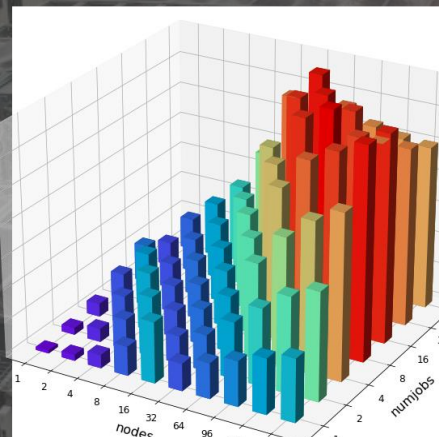
Courtesy of V. Khristenko

# HPC Data Access Testbed

FIO (Flexible I/O tester) is an open-source synthetic benchmark tool
Generates various I/O type workloads (sequential/random reads and writes, etc.)
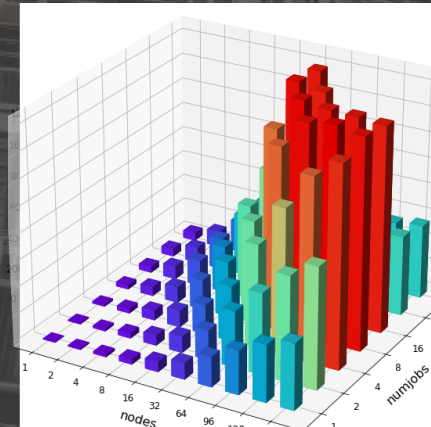
BS=16K
Peak: 1.37 GB/s

BS = 16K
Peak: 8.3 GB/s



RANDOM

SEQUENTIAL

Courtesy of V. Khristenko

# Total Bandwidth vs Nodes

In addition to studying the file systems and storage performance, HEP is investigating how to improve the IO performance of the stored data files

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

- **Scaling:** multi-threaded (>200 cores), distributed backends (dask / spark /...)

Courtesy of A. Naumann

# HEP Data in HPC

ROOT
Data Analysis Framework
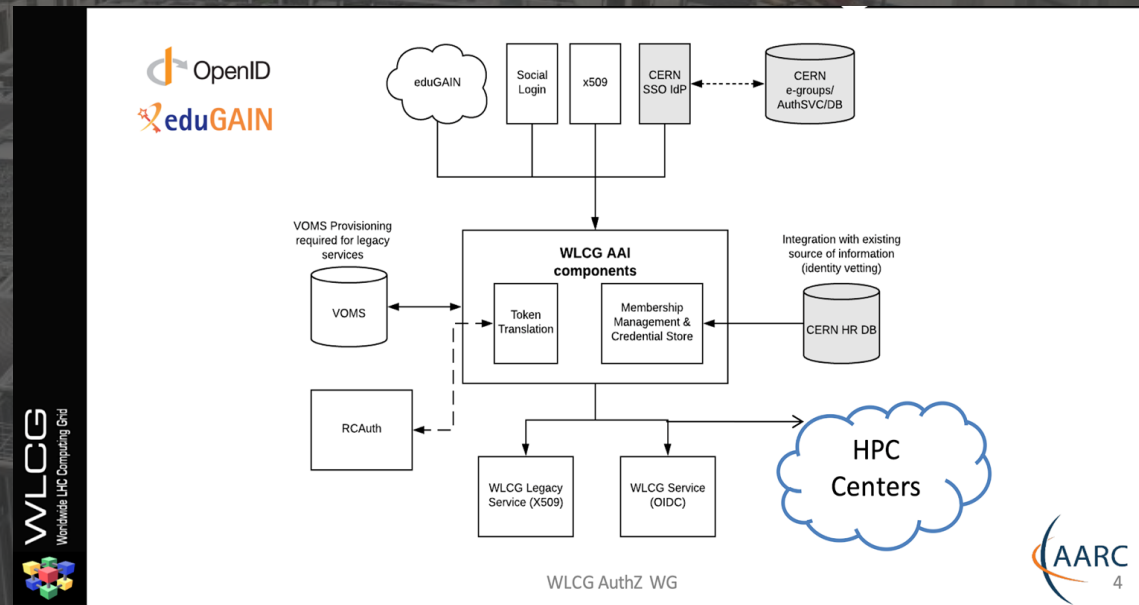
Maria Girone
CERN openlab CTO

**Open Questions**

- Can we accommodate the strict cybersecurity requirement of HPC while enabling access by large international scientific collaborations?
- Can we capitalize on the efforts to modernize AAI infrastructure on the grid infrastructure?

**Technical Program**

- **Take advantage of growing experience in federated identity management (e.g. from AARC and GN4-3 EU Projects)**
- **Test OAuth2 token-based finer grained authorization on HPC, trusting token issuers that belong to scientific collaborations**
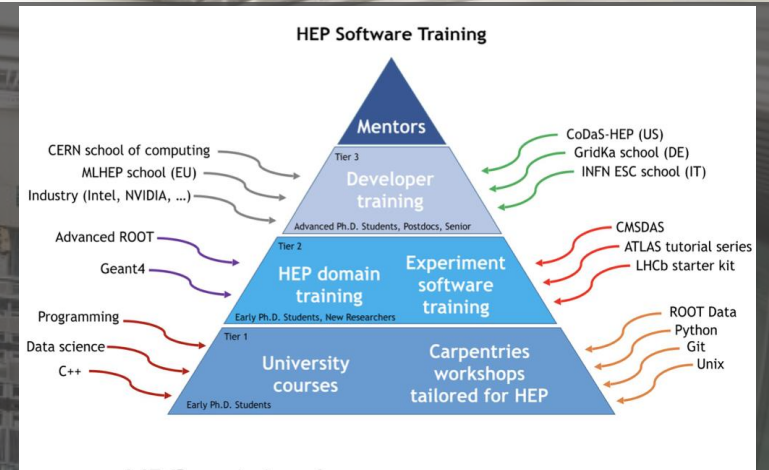  - Token issuer for integration tests available

# Authentication and Authorization Infrastructure

Maria Girone
CERN openlab CTO

**Open questions**

- How to build a common Center of Expertise
- How to make training relevant, scalable, and sustainable in HEP
  - Community investment in software is large, turnover is high, huge user and developer base

**Training Program**

- Dedicated events throughout 2021 on accelerator programming and performance tuning
  - Large expertise in PRACE

- Joint summer internship programs between PRACE Summer of HPC/CERN openlab in Q2 2021, focused on demonstrators
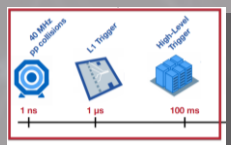


HEP Software Training



PRACE Training Overview in 6 IP

# Building a Center of Expertise

Maria Girone
CERN openlab CTO

**Proven CERN capability** ✓

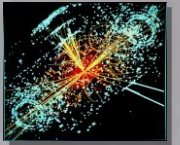**In development, opportunity for joint R&D**

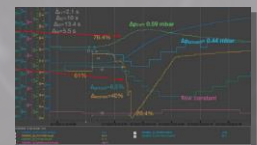**Use case specific**

**Fast ML**
Ultra-fast on-edge inference under strict latency constraints

**Anomaly detection**
Object identification, classification, anomaly detection in big and noisy data sets
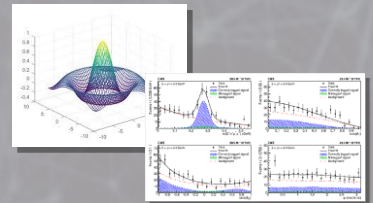
**Industrial controls**
Machine efficiency and predictive maintenance with industrial control systems

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

**Large scale, science grade data analytics and visualization**
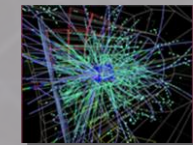
**Cross use case**

- 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)
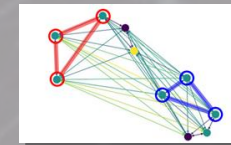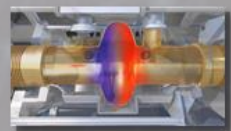
**Simulation**
Simulation and reconstruction with generative DL for efficient computation

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

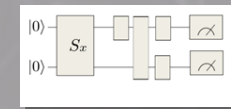Determining optimal machine design and component configuration
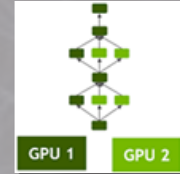
**ML in Robotics**
Remote maintenance and safety with autonomous robots and computer vision

**Quantum ML**
Research quantum algorithms to solve pattern recognition, classification and generation problems

**Computing parallelization**
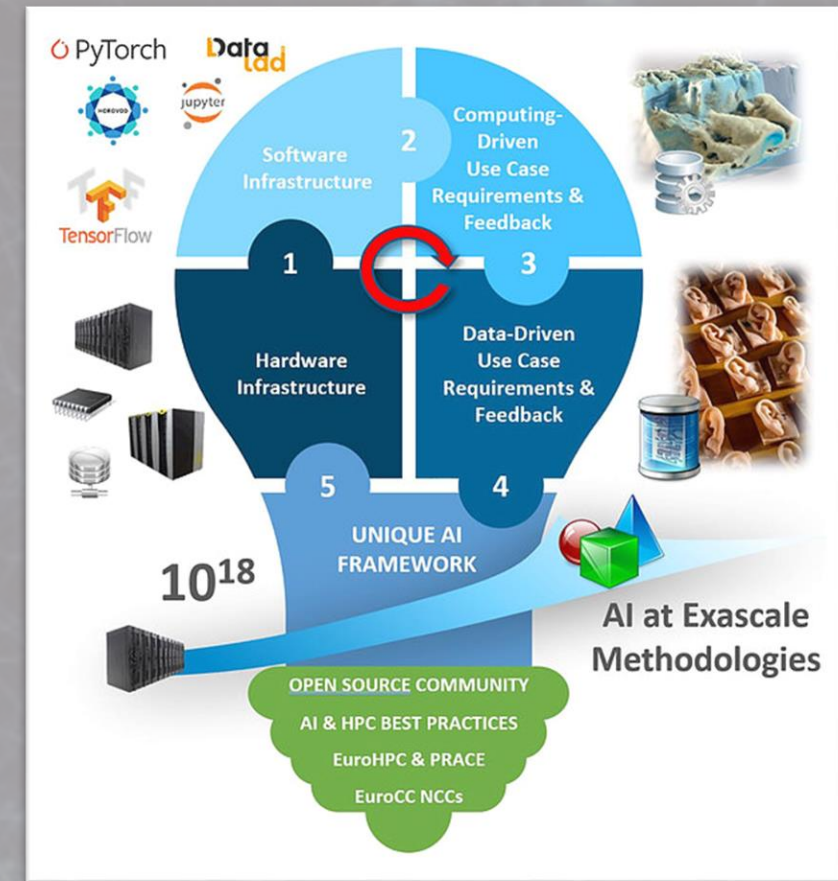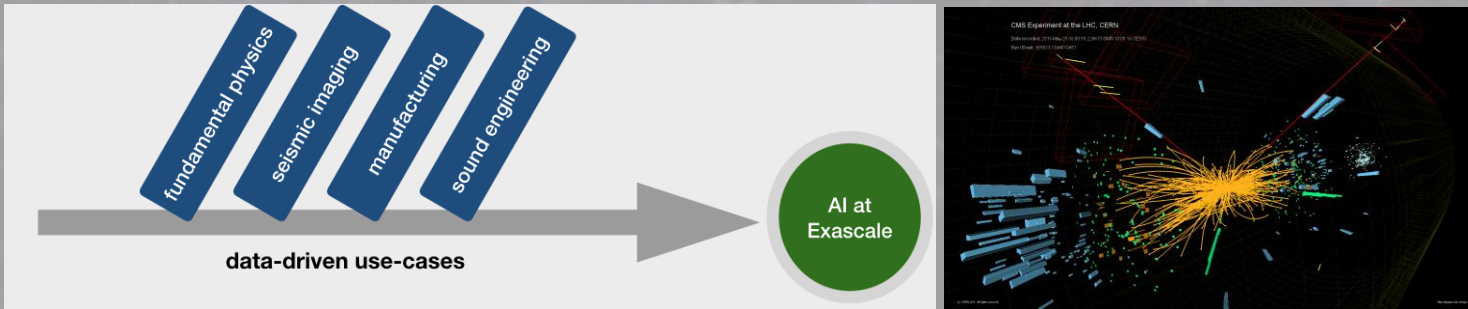Training and optimization of complex NNs on parallelized GPU infrastructure

# Progress on AI/ML Capabilities

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

Testbeds are critical resources for development progress

- Representative hardware architectures, network structures and scale

- Some of the PRACE Tier-0 testbeds will be accessible in Q1 2021 through the collaboration for executing the working groups demonstrators
  - Benchmarking
  - Data access
  - AAI



**PRACE | Tier-0 Systems in 2020**

**NEW ENTRY 2018**
**JUWELS (Module 1):** Atos/Bull Sequana GAUSS @ FZJ, Jülich, Germany #39 Top 500

©FZ Jülich / R.-U. Limbach

**MareNostrum:** IBM BSC, Barcelona, Spain #38 Top 500

**Piz Daint:** Cray XC50 CSCS, Lugano, Switzerland #10 Top 500

**NEW ENTRY 2018/2019**
**SuperMUC NG :** Lenovo cluster GAUSS @ LRZ, Garching, Germany #13 Top 500

**NEW ENTRY 2020**
**HAWK:** HPE Apollo GAUSS @ HLRS, Stuttgart, Germany

© by Ben Derzian for HLRS

© CEA

**NEW ENTRY 2018**
**JOLIOT CURIE :** Atos/Bull Sequana X1000; GENCI @ CEA, Bruyères-le-Châtel, France #34 Top 500

**MARCONI-100: IBM** CINECA, Bologna, Italy #9 Top 500

**Close to 110 Petaflops total peak performance**

5        The Partnership for Advanced Computing in Europe | PRACE

# Testbeds

Participation on the path to exascale

- All stakeholders have engaged and are expected to contribute to the fundamental tasks of **effectively** HPC systems, securely accessing them and intensively accessing data
  - Additional contributions are very welcome
- The next months are crucial to demonstrate collaborative use of HPC systems
  - Access to testbeds proving technical benefits
- Leveraging experience from PRACE and GÉANT as centers of expertise is valuable to our community
  - We will provide application knowledge for a joint program on software optimisation on HPC systems

**AAI:** Token Issuer

**Access for WGs to** testbed resources through the CA

**Benchmark:** Include PRACE

**Data Access:** Rucio and FTS to local storage

**Training:** Summer intern program

**Data Access:** Reach multi 100Gb/s

**Collaboration Workshop**

**Benchmark:** Common Framework

**Benchmark:** Include SKA

**Collaboration Workshop**

**Challenges** wrap-up and future planning

| Q4/2020 | Q1/2021 | Q2/2021 | Q3/2021 | Q4/2021 |

# Towards a Successful Collaboration

Maria Girone
CERN openlab CTO