



# Big Data and AI for Data Analysis

26/10/2020

Alberto Di Meglio - CERN openlab Head, CERN IT Department

Sofia Vallecorsa - CERN openlab AI and Quantum Research, CERN IT Department

João Fernandes - EOSC ARCHIVER project coordinator, CERN IT Department

Value

Veracity

Variability

Variety

Velocity

Volume

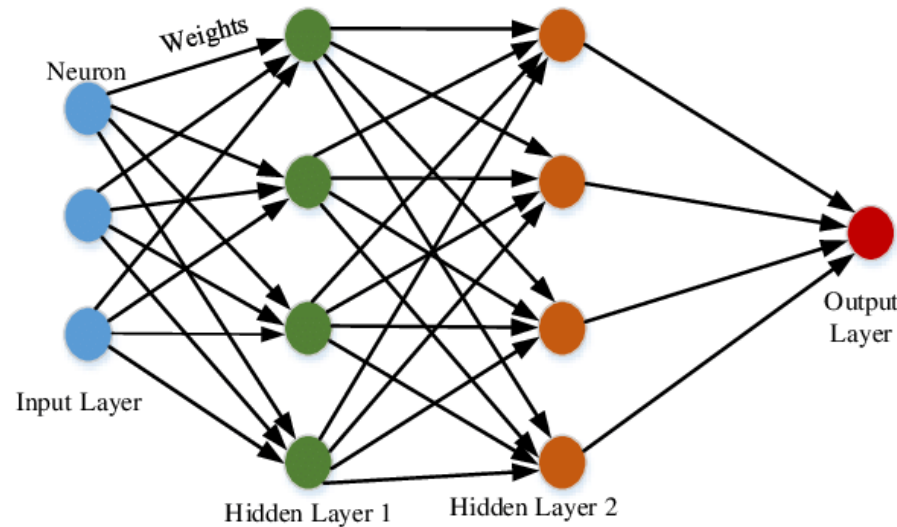
3-V  
Definition

富嶽三十六景 神奈川沖  
浪裏

丁舟の島一景

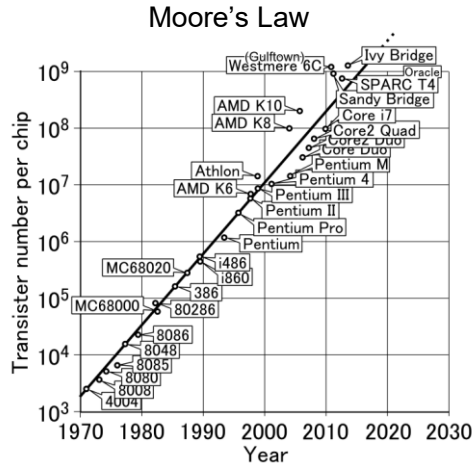
# Artificial Intelligence (ML/DL)

The coincidental combination of the availability of increasing amounts of (noisy!) data, more efficient algorithms, and faster hardware makes Machine/Deep Learning approaches more and more appealing



But... how do we implement it at large scale?

# New Computing Platforms



2000



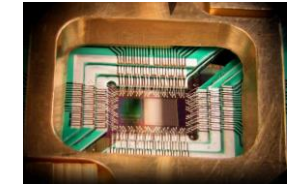
2010



2014



2015



2019

Radically new computing platforms are rapidly moving from pure computer science to realistic devices, e.g. **Neuromorphic Computing and Quantum Computing**

A **Quantum Computing Initiative** has been launched in 2020 as a long-term investigation activity

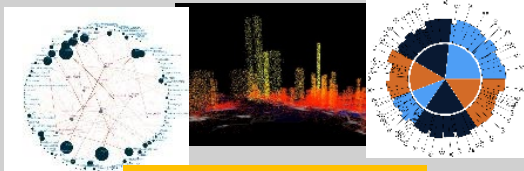


# What does it Really Mean?

Information

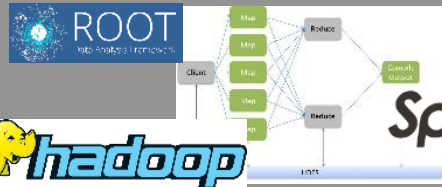


Applications designed for "big data"



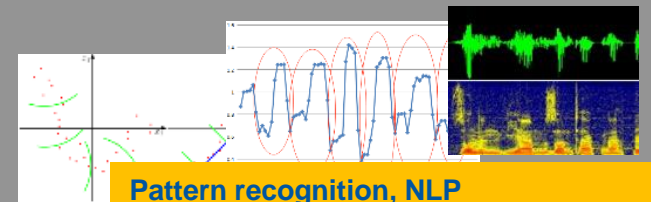
Data visualization, representation

Software



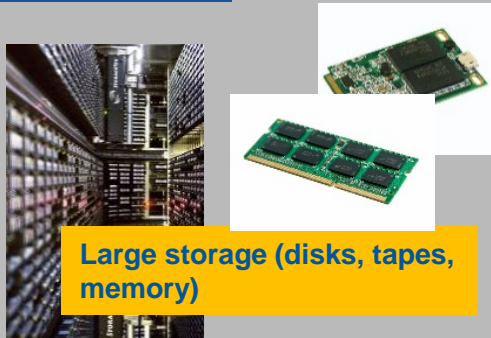
Data analysis and analytics platforms

Value  
Veracity  
Variability  
Variety  
Volume  
Velocity



Pattern recognition, NLP machine learning, predictive analytics

Platform



Large storage (disks, tapes, memory)



Fast data acquisition systems



Flexible networks

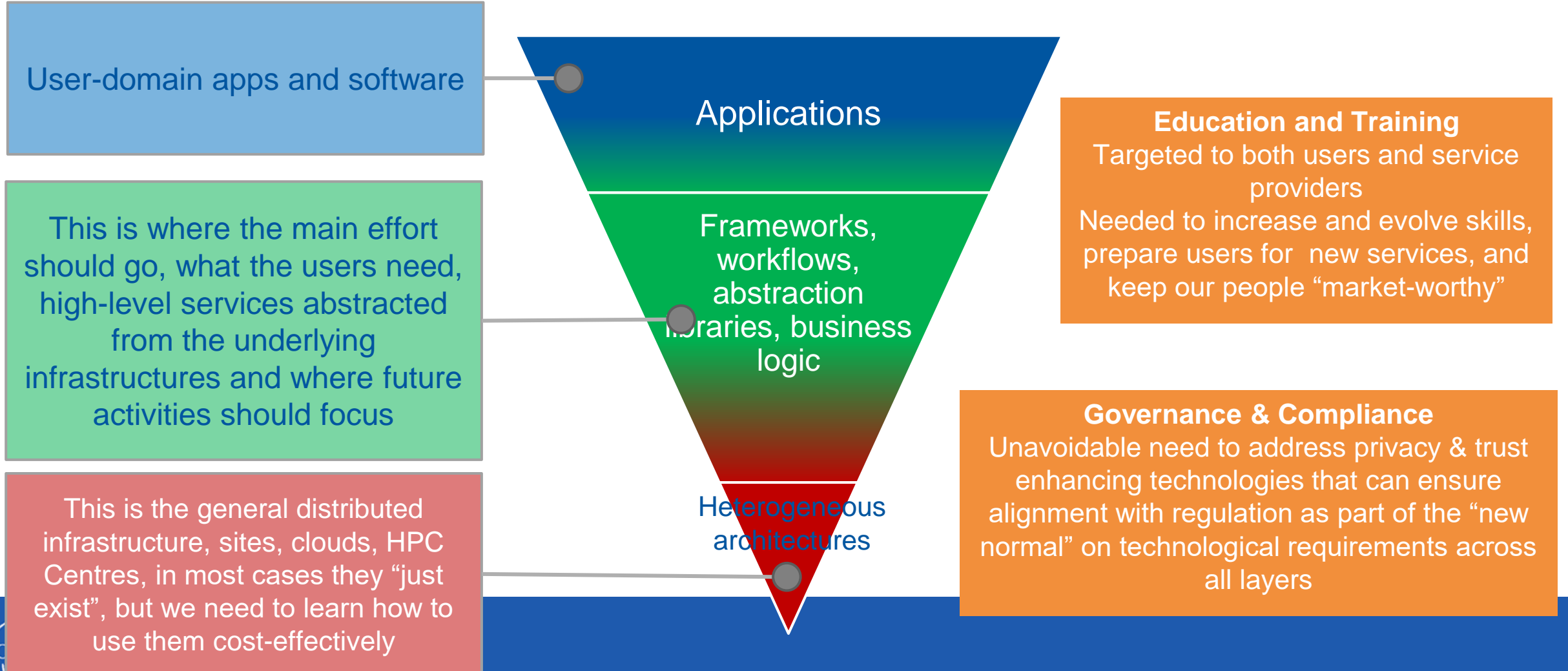


Distributed Computing and Data Grids, Clouds, HPC, Crowd Computing

Infrastructure

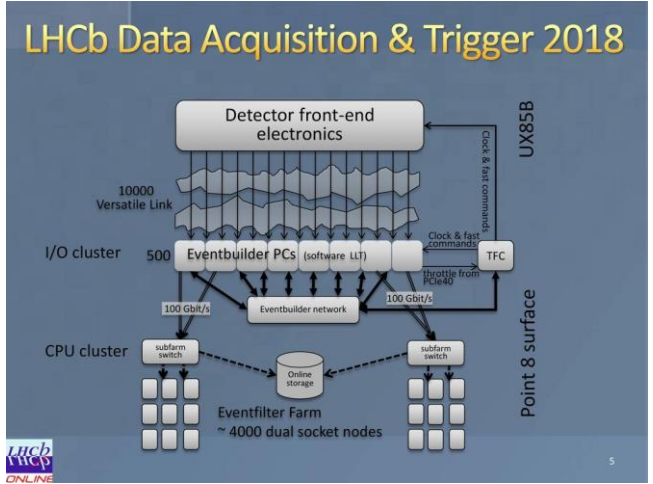
# As Easy as **ABCDE**:

***A**I, **B**ig data, **C**louds, **D**irection, **E**ducation*



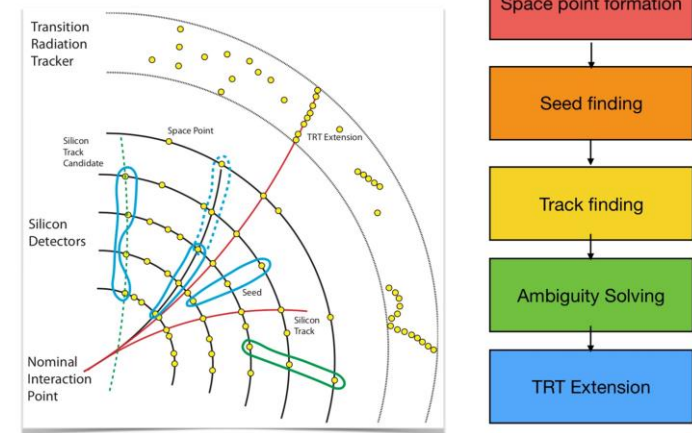
# Typical LHC Experiments Workloads

© Niko Neufeld - LHCb

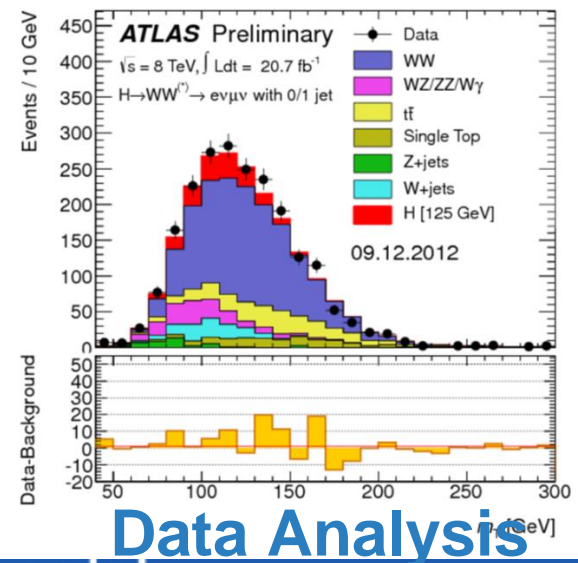


## Data Acquisition

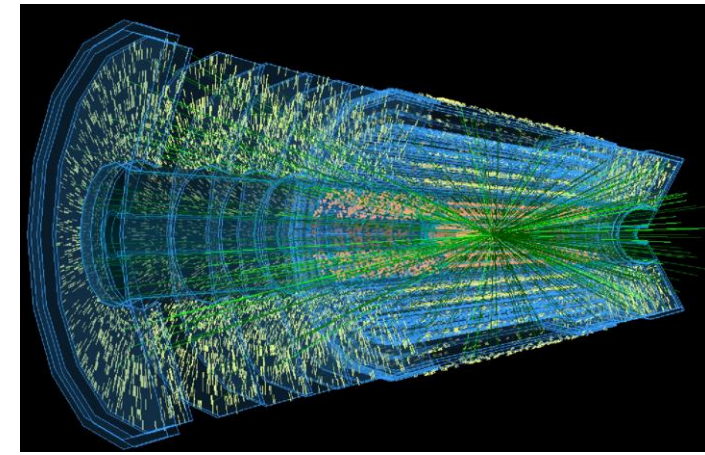
Multi-step iterative Kalman filter approach



## Track Reconstruction



## Data Analysis

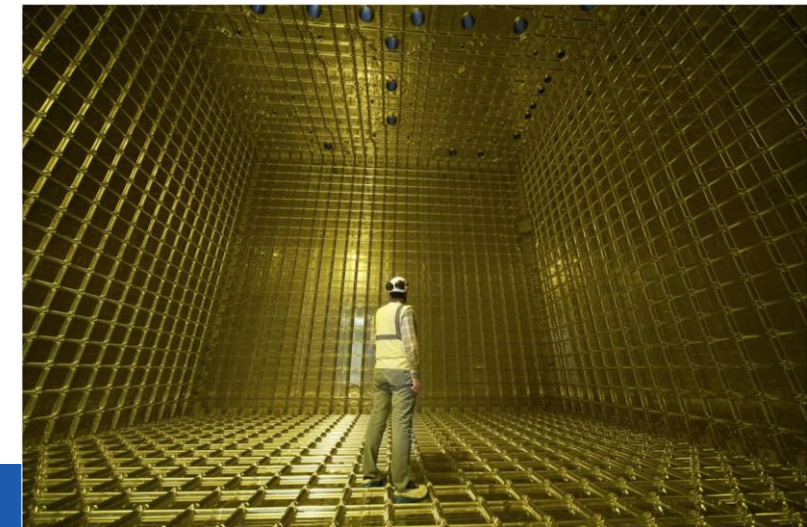
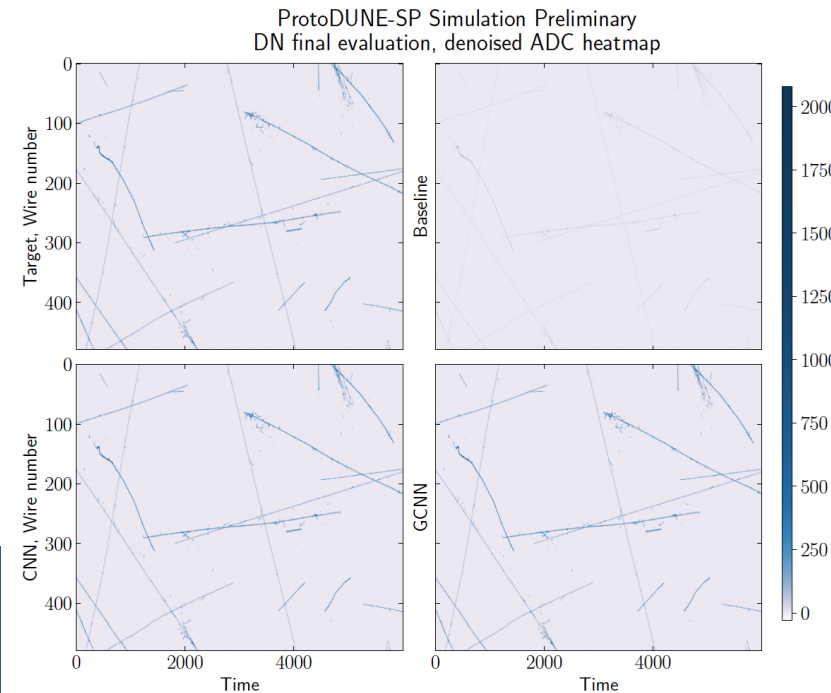
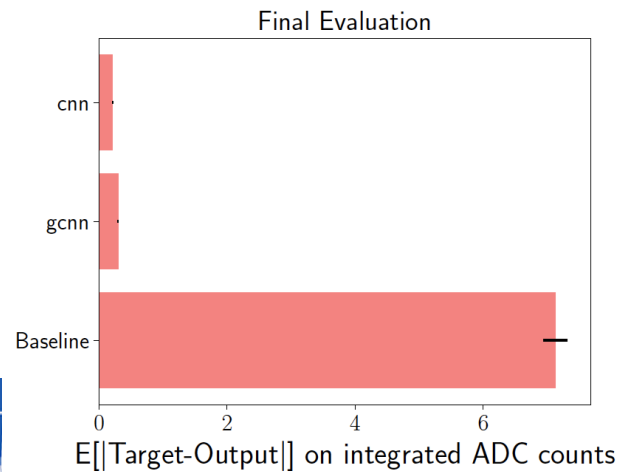
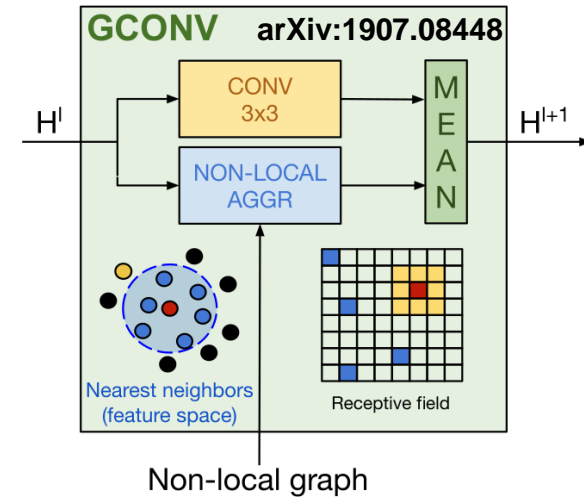


## Simulation



# Data Acquisition

Deep graph-convolution NN to select and denoise raw detector data (Dune) custom architecture



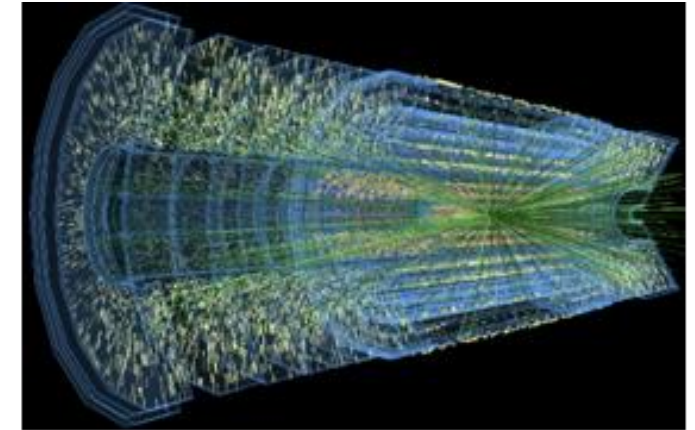
# Track Reconstruction

Exa.TrkX project introduces **Graph Neural Networks** for particle trajectory reconstruction

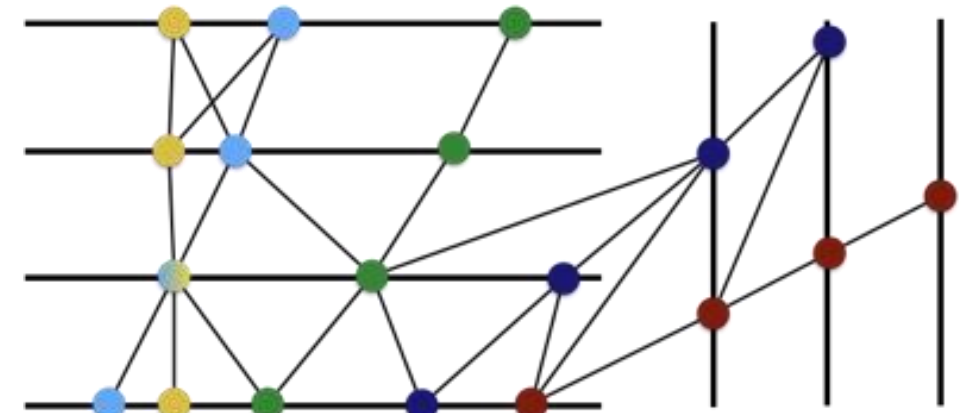
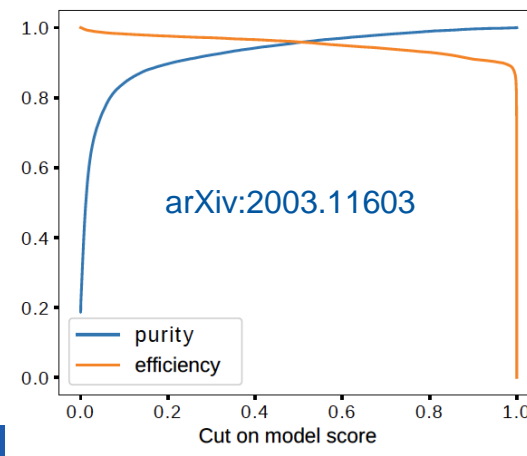
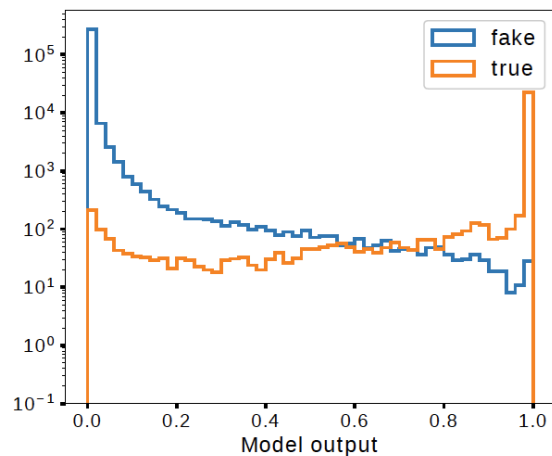
Data as a **graph of connected hits**

Connect plausibly-related hits using geometric constraints

Use a cascade of **Input, Edge and Node GNN Networks**



<https://exatrnx.github.io/>



# Event Classification

Simple, portable and scalable ML platform running on Kubernetes

- Auto Scaling, GPU and Accelerators, Cloud Bursting

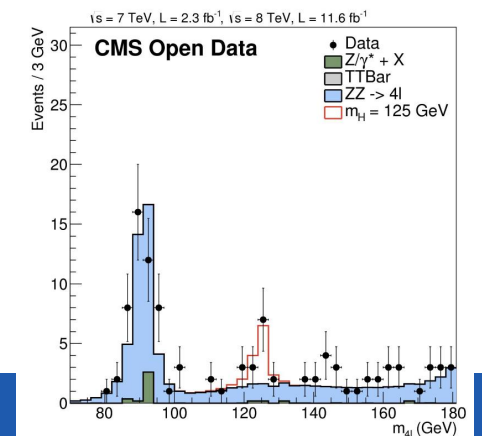
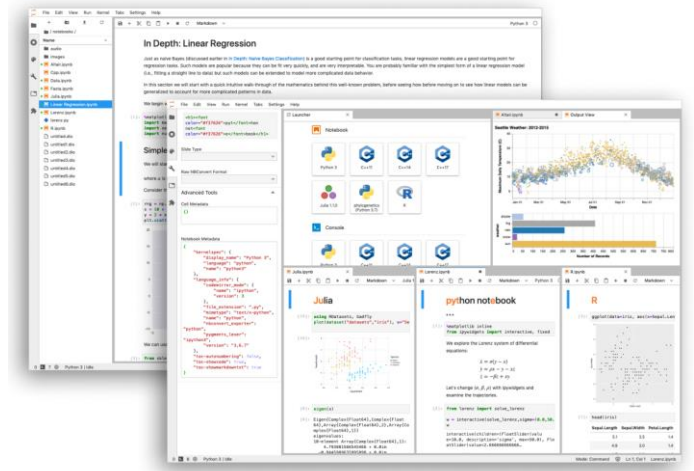
Support entire ML lifecycle

- Development, Training, Inference/Model Serving
- Notebooks, ML Pipelines, Hyper-parameter Optimization
- Tensorflow, PyTorch, scikit-learn, MXNet, MPI, ...

Use Cases: 3DGAN/Fast Simulation, CNNs DUNE, CMS HLT

H $\rightarrow$ 4l re-discovery from CMS Open Data **in 4 minutes using 25k core on Google cloud**: cloud native approach to ML-based analysis for HEP

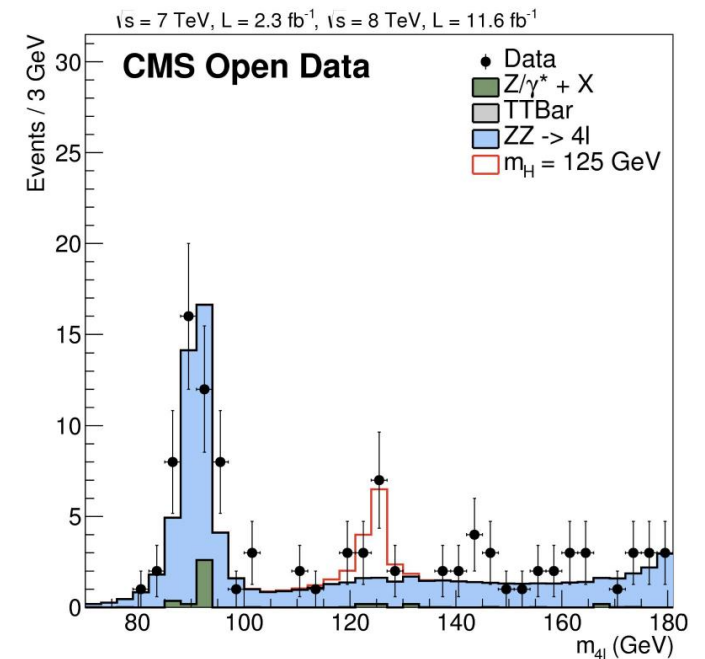
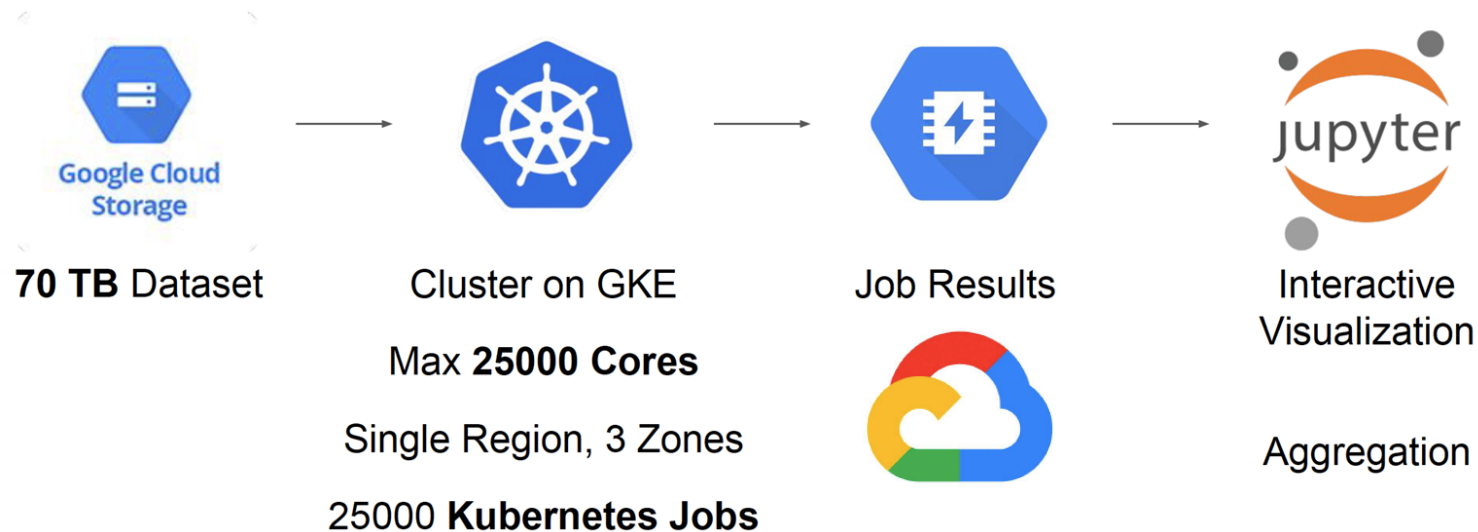
L.Heinrich, R. Rocha, [Keynote at KubeCon EU 2019](#)



# Event Classification

H→4l re-discovery from CMS Open Data **in 4 minutes using 25k core on Google cloud:** cloud native approach to ML-based analysis for HEP

L.Heinrich, R. Rocha, [Keynote at KubeCon EU 2019](#)

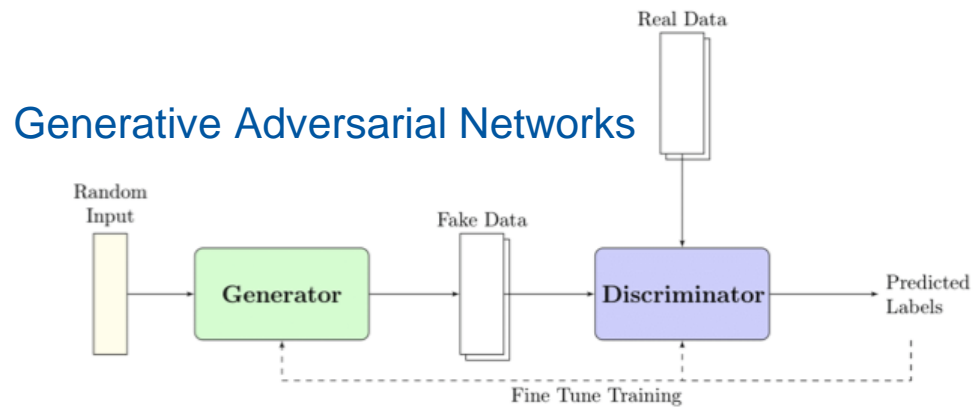


# Simulation

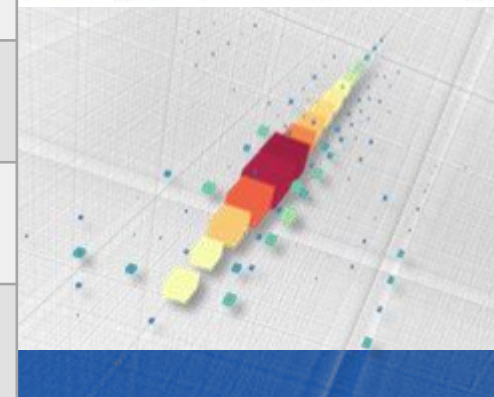
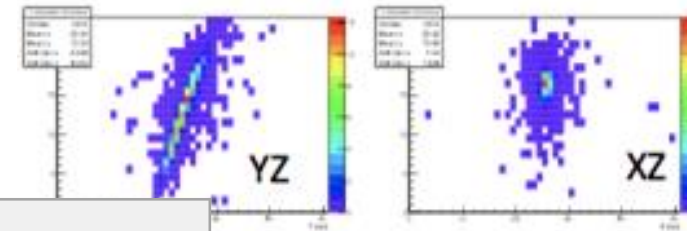
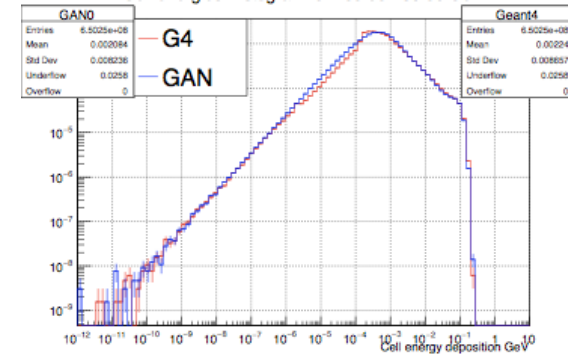
Simulation is a **major workload** in terms of computing resources.

With High Luminosity LHC we expect a x100 increase in simulation need

Investigate **Deep Generative Models** to replace Monte Carlo



G.Khattak, ICMLA2019



Time to create an electron shower		
Method	Machine	Time/particle (msec)
<b>MC Simulation</b>	Intel Xeon Platinum 8180	17000
<b>3DGAN (batch size 128)</b>	Intel Xeon Platinum 8160 (TF 1.8)	1

# Quantum Machine Learning

Quantum linear algebra is **generally faster** than classical counterpart

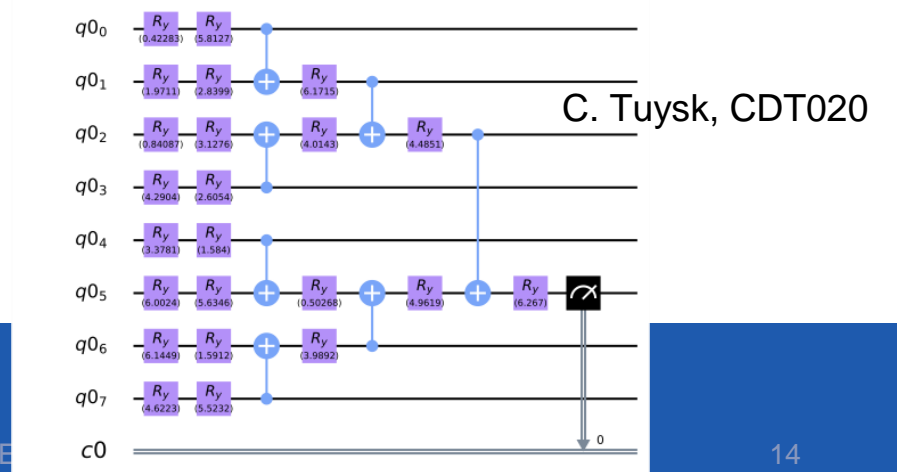
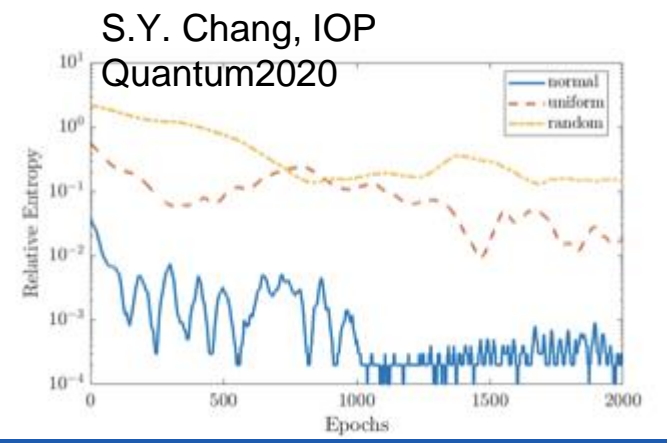
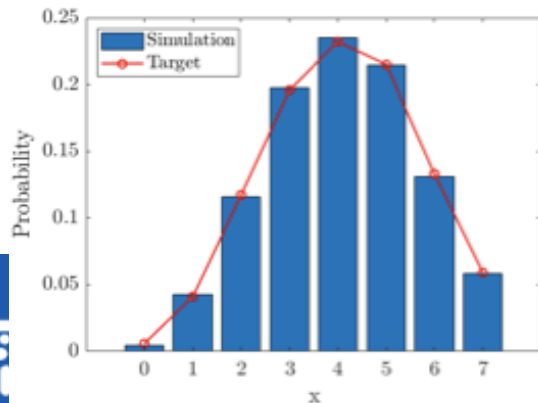
Some standard ML techniques estimate the **ground state of Hamiltonians**

ML algorithms have some **tolerance to errors**

Specific **quantum techniques** can be exploited to bring further improvement

## Partnership with ESA

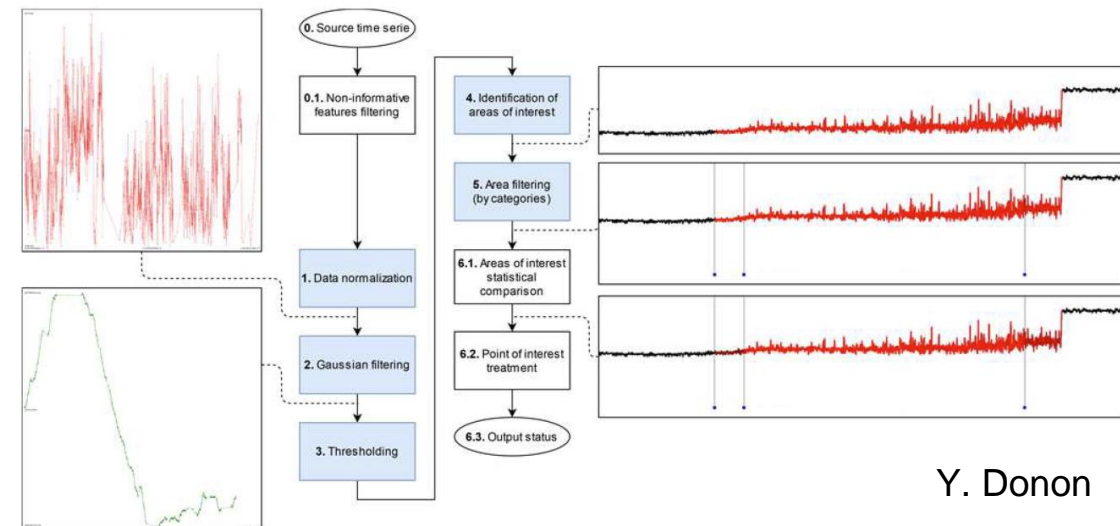
Investigate quantum Generative Models Applications to applications to Earth Observation



# Anomaly Detection, Preventive Maintenance, Control Systems

## Many use cases and ongoing projects

- Anomaly identification in noisy data using SNIFF (Series with Noise Featuring), tests with LINAC4 at CERN, collaborations on medical and industrial LINACs
- CNN-based image analysis for anomaly detection in industrial components
- Edge-computing, IoT, block chains for distributed control systems
- Beam stability and quality control for particle injectors



# AI for Earth Observation

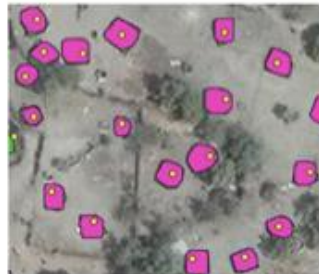
Automatic scan of high-resolution satellite images for **disaster relief**

High precision is required



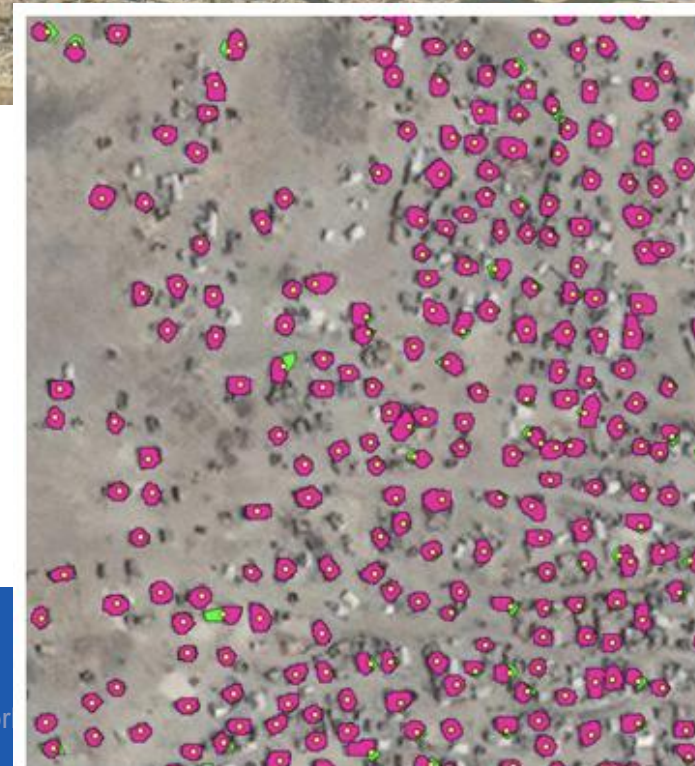
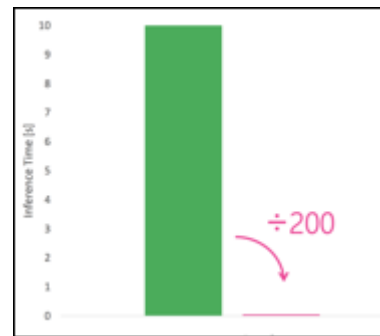
Detectron Framework (FacebookAI)

Retrain & encode point data cleverly



Unosat Adapted model

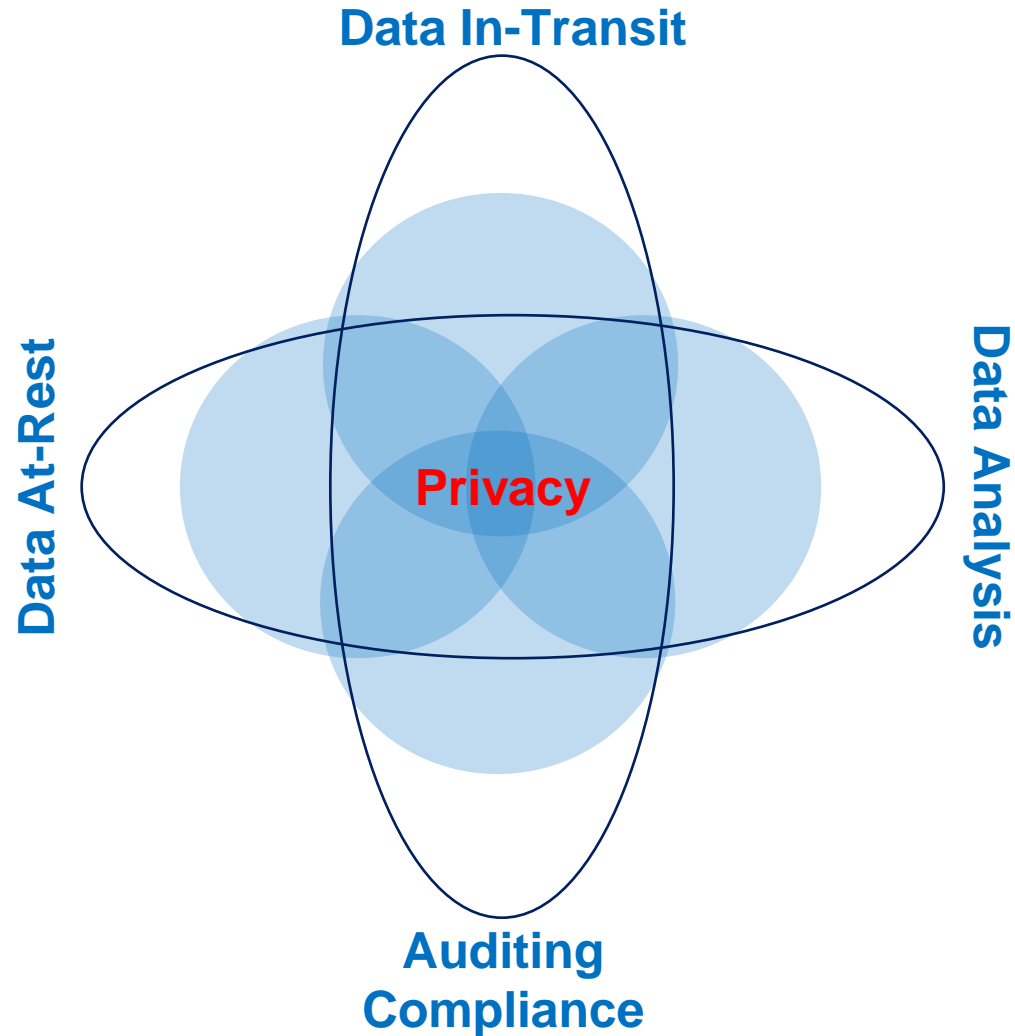
■ Human ■ Neural Net \*



Transfer learning from Region-based CNN  
Average precision is 82%, 200x speedup



# Privacy-Preserving Analysis for Medical Research



**The possible risks of large-scale use of data must be managed**

Use cases at CERN openlab based on methods developed for HEP

- Parkinson's detection from wearable devices using different statistical/ML/DL methods (event classification)
- Explicable AI for phenotype/genotype ("deconvolutional" Neural Networks)
- Image classification and segmentation for neurological diseases research using homomorphic encryption (noise-resilient networks for DAQ, radiation-hardening)
- Block chains and other transaction ledgers (security, but also analysis reproducibility)
- Applications of classic and quantum key distribution methods

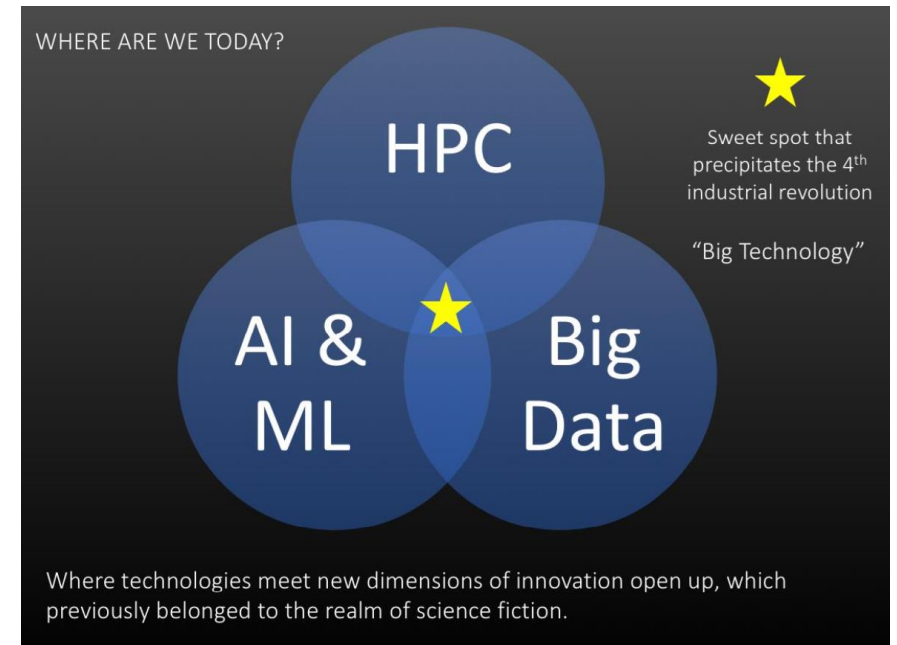
# Changing Business Models

Transparent integration of commercial cloud services, with default data governance: in DMPs, considering procurement cycles & research grant periods

Promote choice: ecosystem for innovation

Stay mainstream by adopting internationally recognized standards

Create a sustainable level playing field:  
same requirements for commercial and not-for-profit providers



# European Infrastructures for Science: EOSC

Framework to make access to scientific data widely available for public/private sectors across domains

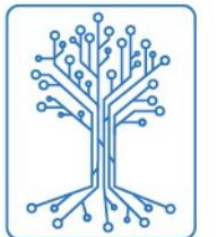


Key point on data preservation & reuse

- ensure a continuum after a mission/experiment finishes



Create the conditions to develop a sovereign European data sharing environment (e.g. explore links with GAIA-X)



**GAIA-X**

# Collaboration and joint initiatives

Joint investigation across distributed heterogeneous platforms are critical for the future of research

- CERN and ESA announced in September a joint partnership on Quantum and AI research for Earth Observation and Physics
- CERN and ESRF are discussing about collaborations on ML/DL for accelerators operations and maintenance
- CERN, SKA, PRACE, GEANT signed an agreement to work on common access models for HPC
- CERN and EMBL have a long history of collaboration on Grid and Cloud infrastructures

# Parting messages



**Tipping point for large-scale research**



**ABC computing offers unprecedented opportunities**



**Distributed platforms and services must become foundational**



**Collaborations across Europe and internationally are critical**



**All domains of science, technology, industry, society can benefit**



**Good governance, trust, reliability, explainability must be built-in**



# Thanks!

*alberto.di.meglio@cern.ch*  
*@AlbertoDiMeglio*