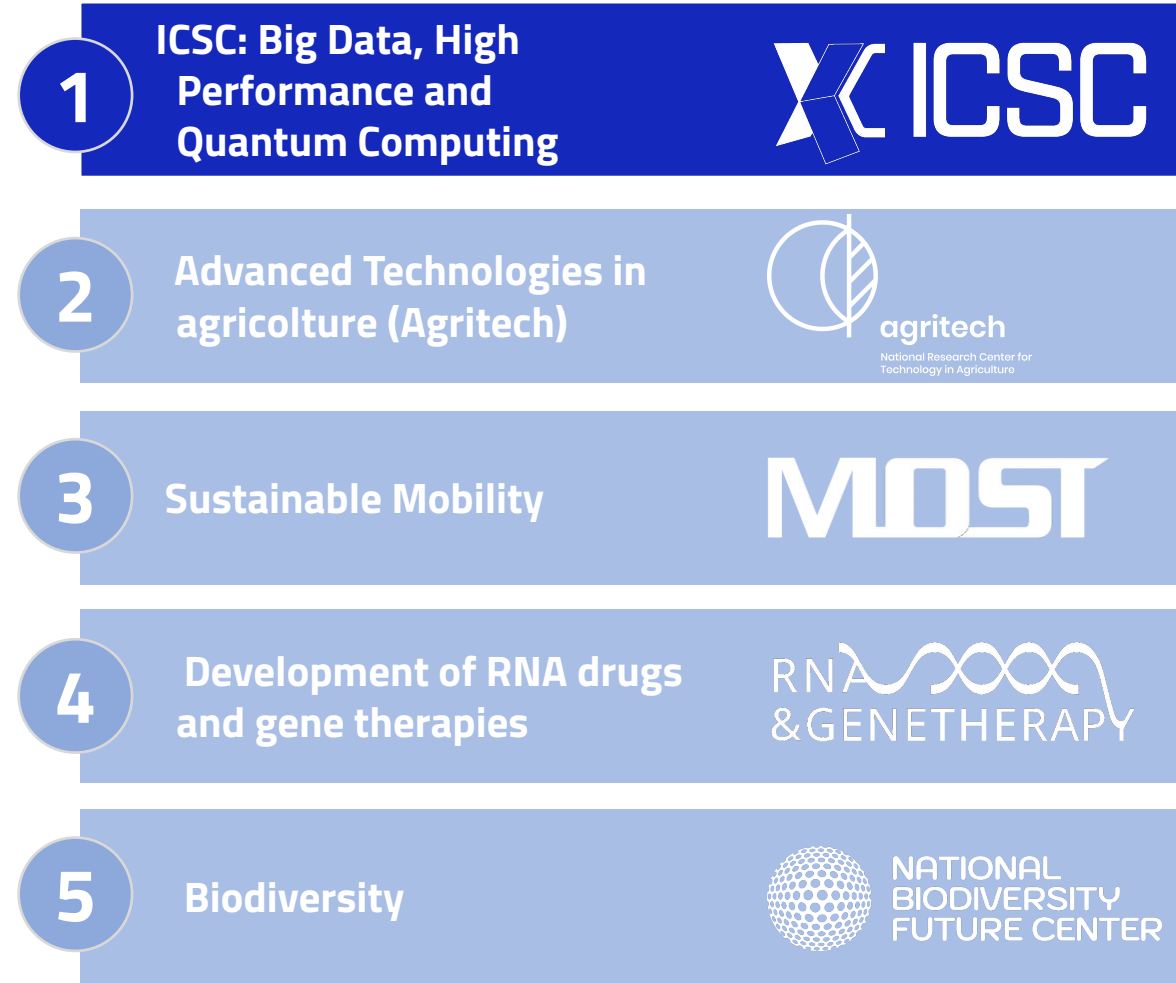# Introduction to High Rate analysis
## Tommaso Diotalevi (UniBO)

**Workshop on "Quasi-Interactive Analysis of Big Data with High Throughput" - 8/10 Jan 2024 - Bologna**

# ICSC: The National Center for HPC, Big Data and Quantum Computing

## what is it?

- Italy has funded, with NRRP (pandemic recovery) funds, **5 large National Centers**, for a total of **1.6B €** over 3 years, on key future technologies.

- One of them, coordinated by **INFN**, focuses on modern IT technologies, with the final goal of deploying a <u>long-term distributed infrastructure</u> (>> 3y) <u>for national research and industrial development</u>.

- The project started on September 2022, lasting until December 2025.

**1** ICSC: Big Data, High Performance and Quantum Computing — ICSC

**2** Advanced Technologies in agricolture (Agritech) — agritech National Research Center for Technology in Agriculture

**3** Sustainable Mobility — MOST

**4** Development of RNA drugs and gene therapies — RNA &GENETHERAPY

**5** Biodiversity — NATIONAL BIODIVERSITY FUTURE CENTER

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

# The participants

## The ICSC foundation: public and private members
an initiative spread across Italy

**25**
**Universities**

### National institutes

INFN
Istituto Nazionale di Fisica Nucleare

CINECA

ENEA

Consiglio Nazionale delle Ricerche

INAF
ISTITUTO NAZIONALE DI ASTROFISICA

ISTITUTO NAZIONALE DI GEOFISICA E VULCANOLOGIA

Consortium GARR

**12**
**Research institutes**

### HU

UNIMORE
UNIVERSITÀ DEGLI STUDI DI MODENA E REGGIO EMILIA

UNIVERSITÀ DI PARMA

OGS

### Privates

enel

ENGINEERING
THE DIGITAL TRANSFORMATION COMPANY

eni

FERROVIE DELLO STATO ITALIANE

FINCANTIERI

fondazione innovazione urbana

autostrade per l'italia

HUMANITAS
RESEARCH HOSPITAL

iFAB
INTERNATIONAL FOUNDATION
BIG DATA AND ARTIFICIAL INTELLIGENCE
FOR HUMAN DEVELOPMENT

INTESA SANPAOLO

LEONARDO

sogei

ThalesAlenia Space
a Thales / Leonardo company

Terna
Driving Energy

UnipolSai
ASSICURAZIONI

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing
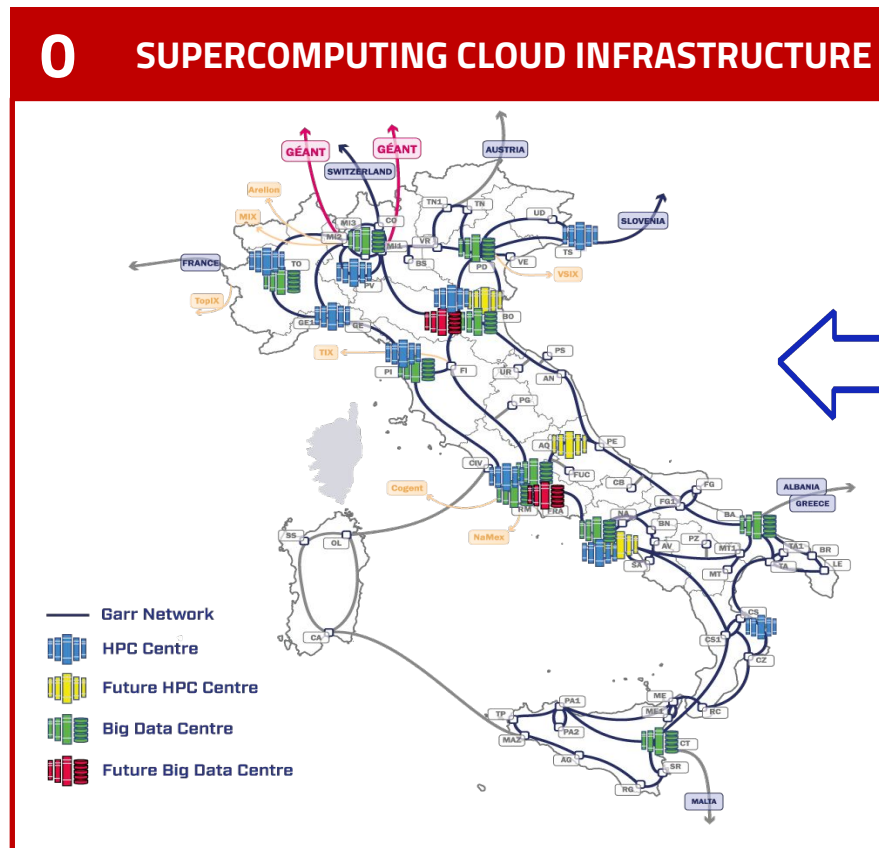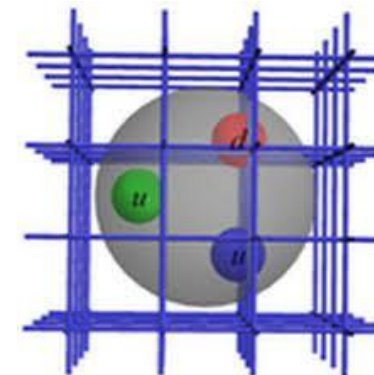
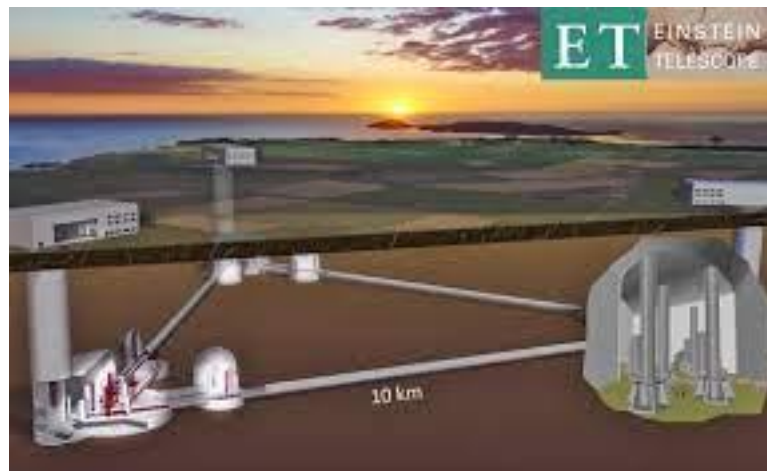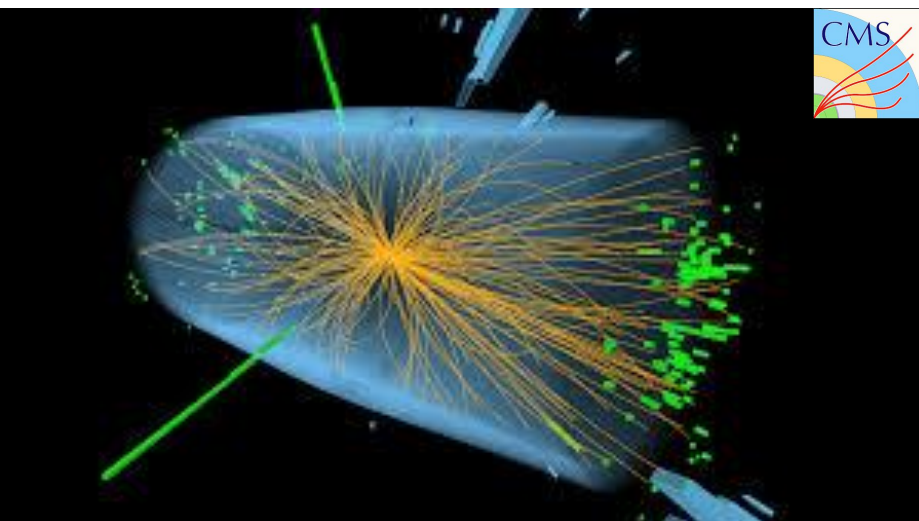# The structure of the ICSC National Center

The ICSC includes:
**10 thematic Spokes** and **1 Infrastructure Spoke**

**0   SUPERCOMPUTING CLOUD INFRASTRUCTURE**



Garr Network
HPC Centre
Future HPC Centre
Big Data Centre
Future Big Data Centre

**EDUCATION AND TRAINING, ENTREPRENEURSHIP, KNOWLEDGE TRANSFER, POLICY MAKING, OUTREACH**

**1**
FUTURE HPC
& BIG DATA

**2**
FUNDAMENTAL
RESEARCH
& SPACE ECONOMY

**3**
ASTROPHYSICS &
COSMOS OBSERVATIONS

**4**
EARTH
& CLIMATE

**5**
ENVIRONMENT
& NATURAL DISASTERS

**6**
MULTISCALE MODELING
& ENGINEERING APPLICATIO

**7**
MATERIALS &
MOLECULAR SCIENCES

**8**
IN-SILICO MEDICINE
& OMICS DATA

**9**
DIGITAL SOCIETY
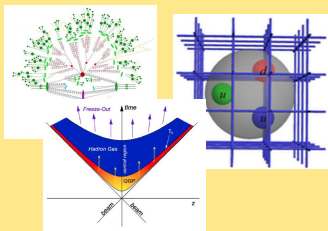& SMART CITIES

**10**
QUANTUM
COMPUTING

- One Spoke dedicated to building the infrastructure

- Ten thematic Spokes, one of which dedicated to the HEP and Astroparticle research domains.
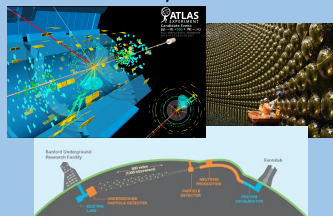
# Spoke 2 – Who are «we»?

# The structure of Spoke 2

**Scientific**
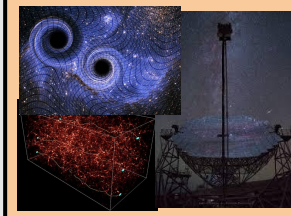
**WP1**: tools and algorithms for Theoretical Physics

**WP2**: tools and algorithms for Experimental High Energy Physics

**WP3**: tools and algorithms for Experimental Astroparticle Physics and Gravitational waves

ICSC-SPOKE2
Centro Nazionale di Ricerca in HPC, Big Data and Quantum Computing

**WP6**: cross domain initiatives + space economy

**WP5**: Boosting computational performance on the distributed CN infrastructure

**WP4**: tools for porting/optimization on new architectures (low power, GPU, FPGA, ...)

**Technologic**

## FUNDAMENTAL RESEARCH & SPACE ECONOMY

Istitution leader: INFN Istituto Nazionale di Fisica Nucleare

Istitution co-leader: INAF ISTITUTO NAZIONALE DI ASTROFISICA

Istitutions and Universities: Politecnico di Bari, SAPIENZA Università di Roma, UNIVERSITÀ DEGLI STUDI DI BARI ALDO MORO, ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA, UNIVERSITÀ DELLA CALABRIA, Università di Catania, Università degli Studi di Ferrara, UNIVERSITÀ DEGLI STUDI FIRENZE, UNIVERSITÀ DEGLI STUDI BICOCCA, UNIVERSITÀ DEGLI STUDI DI NAPOLI FEDERICO II, UNIVERSITÀ DEGLI STUDI DI PADOVA, UNIVERSITÀ DEL SALENTO, UNIVERSITÀ DEGLI STUDI DI TRIESTE

Companies: eni, INTESA SANPAOLO, LEONARDO, ThalesAlenia Space, UnipolSai Assicurazioni, iFAB, sogei

| | |
|---|---|
| Staff Researchers | 195 |
| (kEur) | 6333 |
| Rectruited researchers | 28 |
| (kEur) | 5067 |
| Phd positions | 25 |
| (kEur) | 1992 |
| Budget Innovation Grants (kEur) | 1800 |
| Bugdet Cascade Calls (kEur) | 3200 |
| **Total Budget (kEur)** | **18391** |

**Scientific**
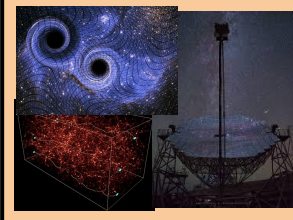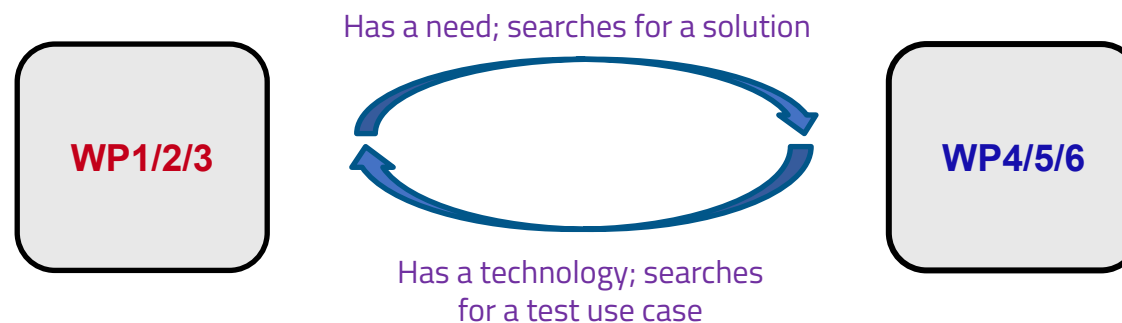
**WP1**: tools and algorithms for Theoretical Physics

**WP2**: tools and algorithms for Experimental High Energy Physics

**WP3**: tools and algorithms for Experimental Astroparticle Physics and Gravitational waves

**WP4**: tools for porting/optimization on new architectures (low power, GPU, FPGA, …)

**WP6**: cross domain initiatives + space economy

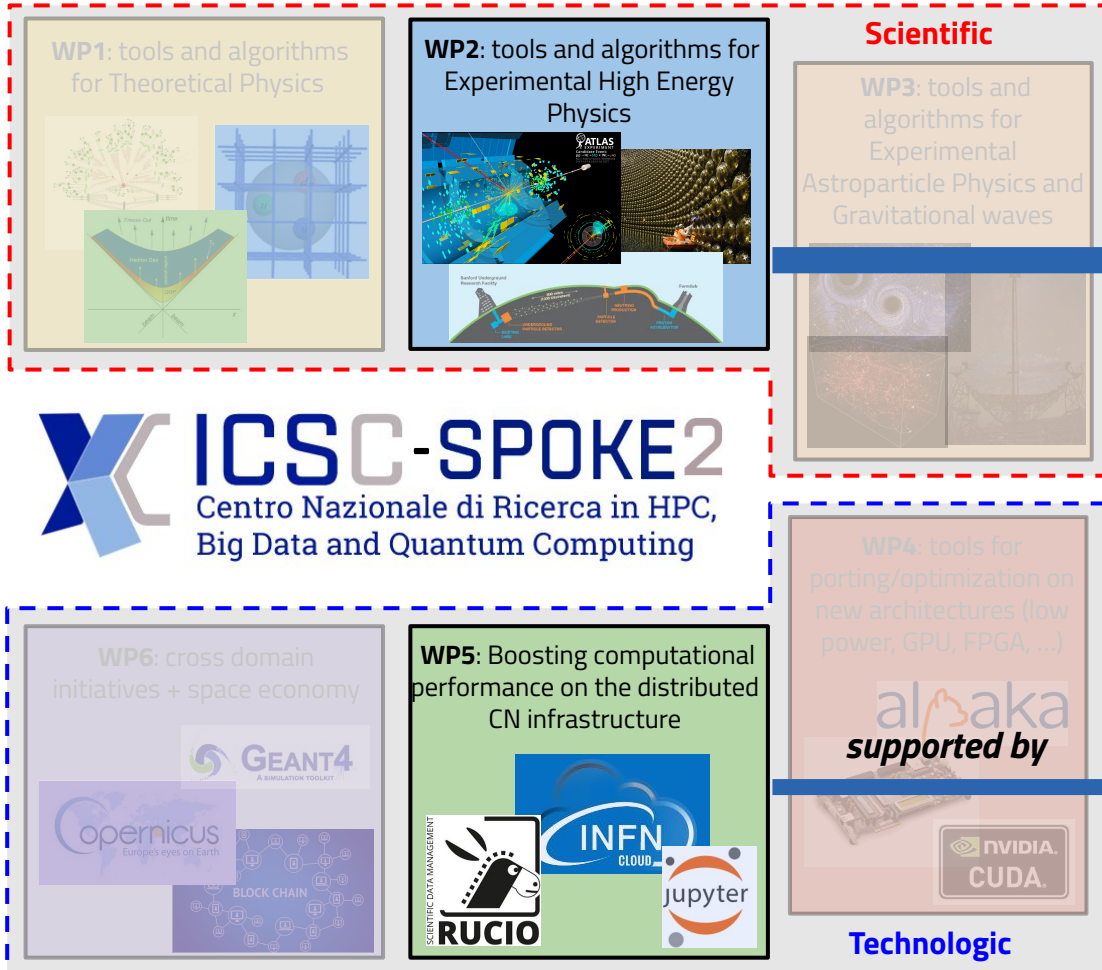**WP5**: Boosting computational performance on the distributed CN infrastructure

**Technologic**

- We defined 2 types of Work Packages (WP):
  - **"Scientific"** WPs: they analyze the needs of the (sub-)domain, and pose open problems for which advanced computing solutions are needed;
  - **"Technological"** WPs: they harvest/investigate technical solutions in computing, on the infrastructure of the ICSC and beyond, and provide support / training for these; at the same time propose these to a larger audience, including industries.

Has a need; searches for a solution

**WP1/2/3**

**WP4/5/6**

Has a technology; searches for a test use case

# Quasi interactive analysis of big data with high throughput



WP1: tools and algorithms for Theoretical Physics

WP2: tools and algorithms for Experimental High Energy Physics

**Scientific**

WP3: tools and algorithms for Experimental Astroparticle Physics and Gravitational waves

WP6: cross domain initiatives + space economy

WP5: Boosting computational performance on the distributed CN infrastructure

WP4: tools for porting/optimization on new architectures (low power, GPU, FPGA, ...)

*supported by*

**Technologic**

**Quasi interactive analysis of big data with high throughput**

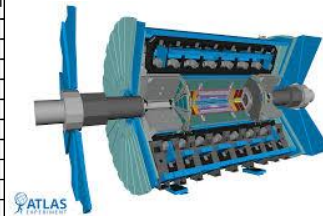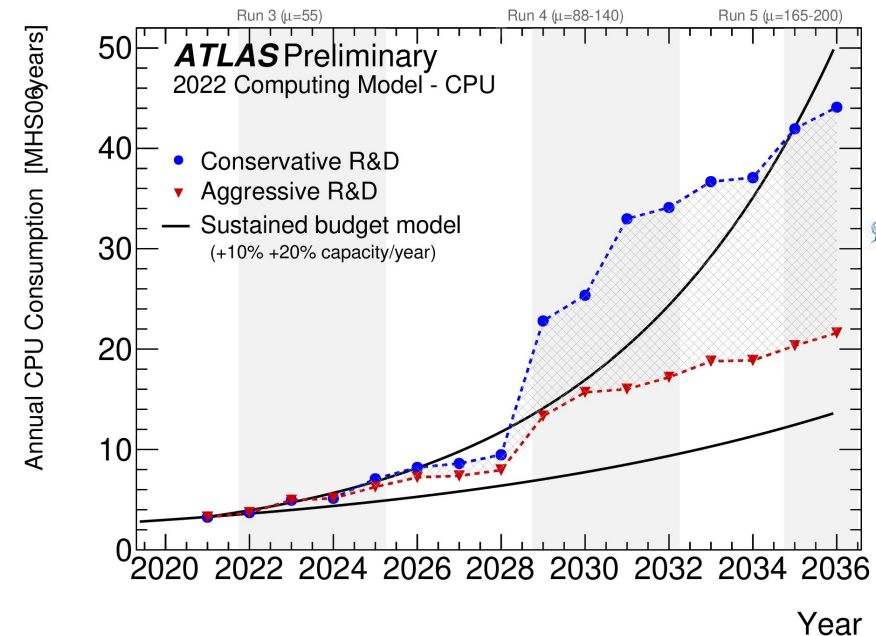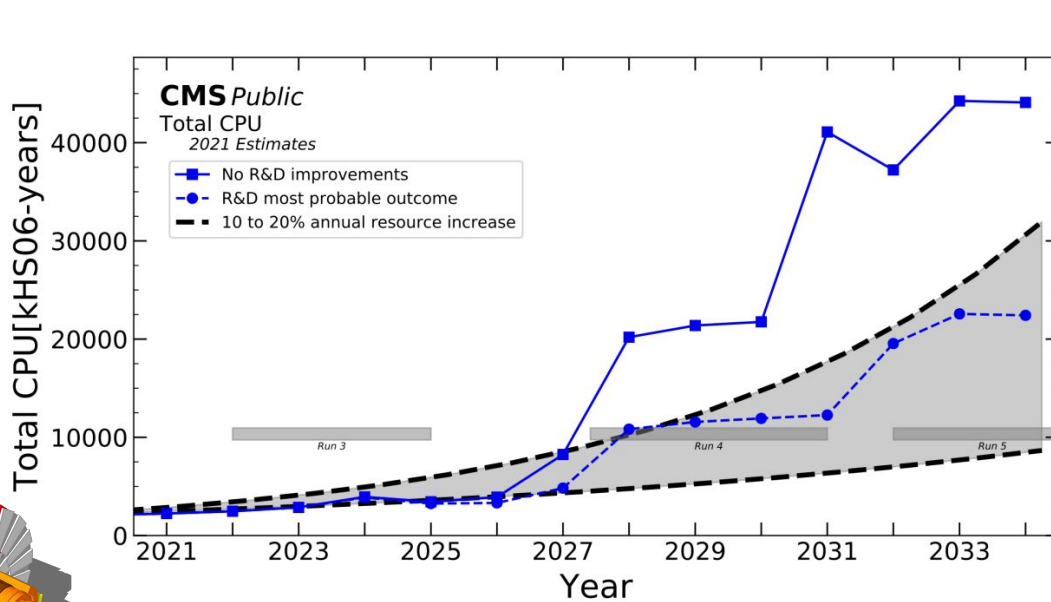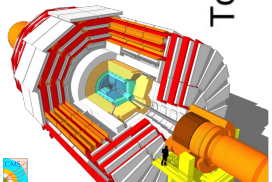| Spoke | 2 |
|---|---|
| WP | 2, 5 |
| Use case short name | Quasi interactive analysis of big data with high throughput |
| Use case ID | UC2.2.2 |
| Expected Completion | 31/8/2025 |

*Approval workflow*

| Status | Version | Date | Submitter | Note | Signature |
|---|---|---|---|---|---|
| Draft | 1.0 | 03/07/23 | WP Leaders | First version | |
| Final Version | 1.1 | 1/9/2023 | WP Leaders | | |
| Approved by Spoke Leaders | 1.1 | 11/9/2023 | Spoke Leaders | | |

*More than 40 people involved in this activity!*

# Why we need this activity?

- Analysing large amounts of data efficiently, exploiting the available resources as much as possible, is a <u>common challenge</u> both for research and industry.

- From the beginning, the High Energy Physics (HEP) experiments at CERN, gave much attention to the computing and data management aspects. Nevertheless, the **next phases of the Large Hadron Collider** (HL-LHC) will require <u>an even greater effort</u>.

# Why we need this activity?

**Some estimate for the next 5-10 years of CMS operation:**

- ~30 Billion collision events + 30 Billion simulation events;
- Each event: 2-4 kB;
- The last update of the CMS Computing model foresees this throughput:

| Name | Length | % of the dataset | Data to process | Event, data rate |
|---|---|---|---|---|
| "A coffee" | < 5 min | 1% (~0.6B evts) | ~2 TB | ~1.7MHz, ~7GB/s |
| "A lunch break" | 1 hour | 10% (~6B evts) | ~20 TB | ~1.5MHz, ~6GB/s |
| "A night" | 12 hours | 100% (60B evts) | ~200 TB | ~1.2MHz, ~5GB/s |

- Difficult to get more than 100 Hz/CPU core → needs efficient distribution on a few tens of machines;

**New analysis paradigm based on:**

**Not only concerning the HEP domain ("Data is data"):**

- More and more scientific / industrial / societal domains have or will have soon needs similar to those from LHC:

ProtoDune: 2-3GB/s (like CMS); Real Dune: 80x

SKA: up to 2 PB/day;

CTA projects: up to 10PB/y
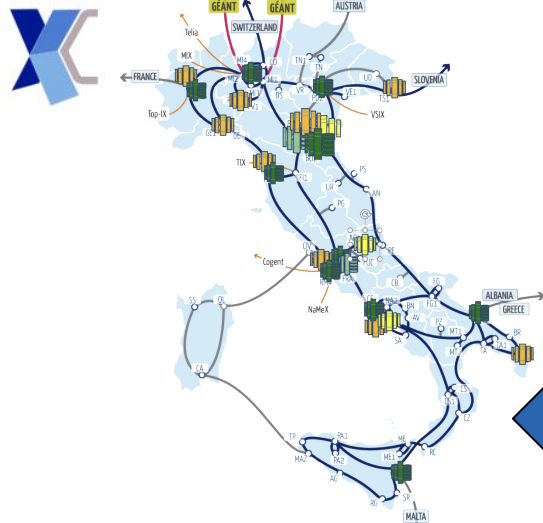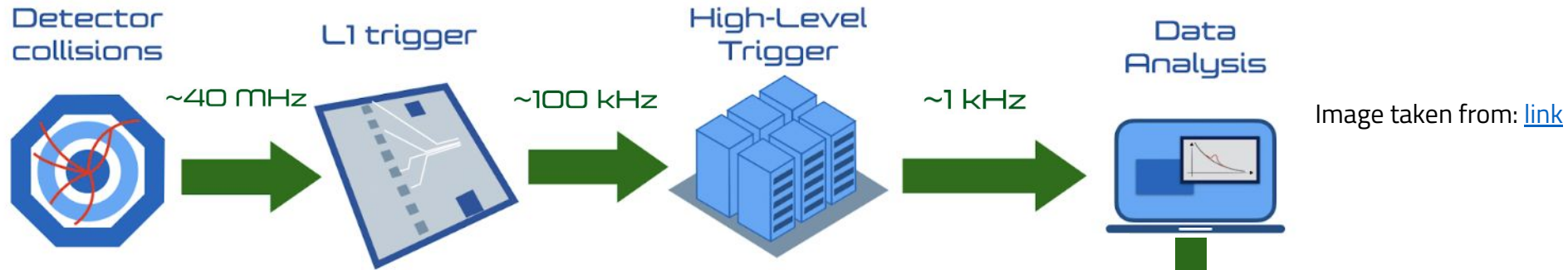
A single genome: ~100GB, a 1M survey=100PB

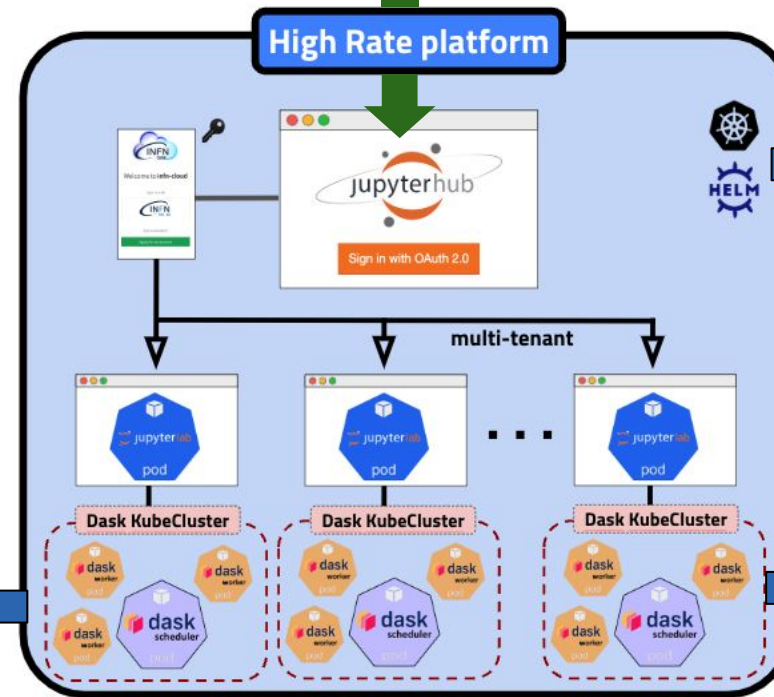O(50 TB/y) per sensor; ~10-100 sensors: O(5 PB/y)

- Declarative programming and interactive workflows;
- Distributed computing on geographically separated resources.

# Why we need this activity?

**Some estimate for the next 5-10 years of CMS operation:**

- ~30 Billion collision events + 30 Billion simulation events;
- Each event: 2-4 kB;
- The last update of the CMS Computing model foresees this throughput:

| Name | Length | % of the dataset | Data to process | Event, data rate |
|---|---|---|---|---|
| "A coffee" | < 5 min | 1% (~0.6B evts) | ~2 TB | ~1.7MHz, ~7GB/s |
| "A lunch break" | 1 hour | 10% (~6B evts) | ~20 TB | ~1.5MHz, ~6GB/s |
| "A night" | 12 hours | 100% (60B evts) | ~200 TB | ~1.2MHz, ~5GB/s |

- Difficult to get more than 100 Hz/CPU core → needs efficient distribution on a few tens of machines;

**New analysis paradigm based on:**

**High Throughput Platform**

- Declarative programming and interactive workflows;
- Distributed computing on geographically separated resources.

**Not only concerning the HEP domain ("Data is data"):**

- More and more scientific / industrial / societal domains have or will have soon needs similar to those from LHC:

DUNE — DEEP UNDERGROUND NEUTRINO EXPERIMENT: ProtoDune: 2-3GB/s (like CMS); Real Dune: 80x

SKA — SQUARE KILOMETRE ARRAY: SKA: up to 2 PB/day;

CTAO: CTA projects: up to 10PB/y

A single genome: ~100GB, a 1M survey=100PB

O(50 TB/y) per sensor; ~10-100 sensors: O(5 PB/y)

# Re-thinking the analysis pipeline



Detector collisions → ~40 MHz → L1 trigger → ~100 kHz → High-Level Trigger → ~1 kHz → Data Analysis

Image taken from: link

**More on T.Tedeschi talk**

Deployment of the Kubernetes resources handled via HELM Charts.

(Scalable on available resources)

offloaded to...

- The execution happens in the Dask Cluster;
- Users choose on how many cores parallelly distribute the analysis.

# Activities (so far) orbiting around the flagship



Vector Boson Scattering ssWW analysis in hadronic tau and light lepton

Neural Network hyperparameter optimisation applied to future colliders (FCC-ee)

Heavy Neutral Lepton search on heavy neutrinos in the $D_s$ decays

Declarative paradigms for analysis description and implementation

Muon detector performance analysis

Search of rare events in tau to 3 muons decay

Continuous Integration pipeline, triggering analysis execution on HTP

top quark+MET analysis

di-Higgs decaying to two b quarks and two muons

Benchmark interactive analysis for future colliders (FCC-ee)

Differential cross section measurement for ttbar inclusive production

Search for new phenomena in events with two opposite-charge leptons, jets and missing transverse momentum

**Specific analyses demonstrators will be shown on the closed CMS session!**

**With the infrastructural support of WP5**

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

# Scientific production in conferences

- Poster at the "International Workshop on Advanced Computing and Analysis Techniques in Physics Research (ACAT 2024)":
  - *Declarative paradigms for analysis description and implementation.*
  - *Quasi interactive analysis of High Energy Physics big data with high throughput.*

- Talk at the "Incontri di Fisica delle Alte Energie (IFAE 2024)":
  *Analisi quasi-interattiva per big data con alto throughput per la Fisica delle Alte Energie.*

- Talk at the "International Conference on High Energy Physics (ICHEP 2024)":
  *Enhancing CMS data analyses using a distributed high throughput platform.*

- Talk at the "2nd European Committee for Future Accelerator (ECFA) Workshop on Higgs/EW/Top Factories":
  *Benchmark interactive analysis for future colliders.*

- Talk at the "Conference on Computing in High Energy and Nuclear Physics (CHEP 2024)":
  *Leveraging distributed resources through high throughput analysis platforms for enhancing HEP data analyses.*

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

# Why this workshop?

- The challenge presented by the next LHC phases requires a strong development effort of new tools, for making data analysis as efficient and as modern as possible;

- Several analysis from the HEP world are <u>already testing</u> such infrastructure, for performance measurements;

- Thanks to the ICSC national center, we have the unique opportunity to <u>build a modern infrastructure</u> for research, aligned with the needs from High Energy Physics and beyond;

- What we need, at this point, is to <u>form scientists (young and senior)</u>, to exploit and perform their research activities in the most efficient and modern way possible!

  - This is our motivation for this event (**hoping for more to come**)!

*For now...*

*Enjoy the workshop!*