# WP2.5

Elvira Rossi

Università Federico II di Napoli
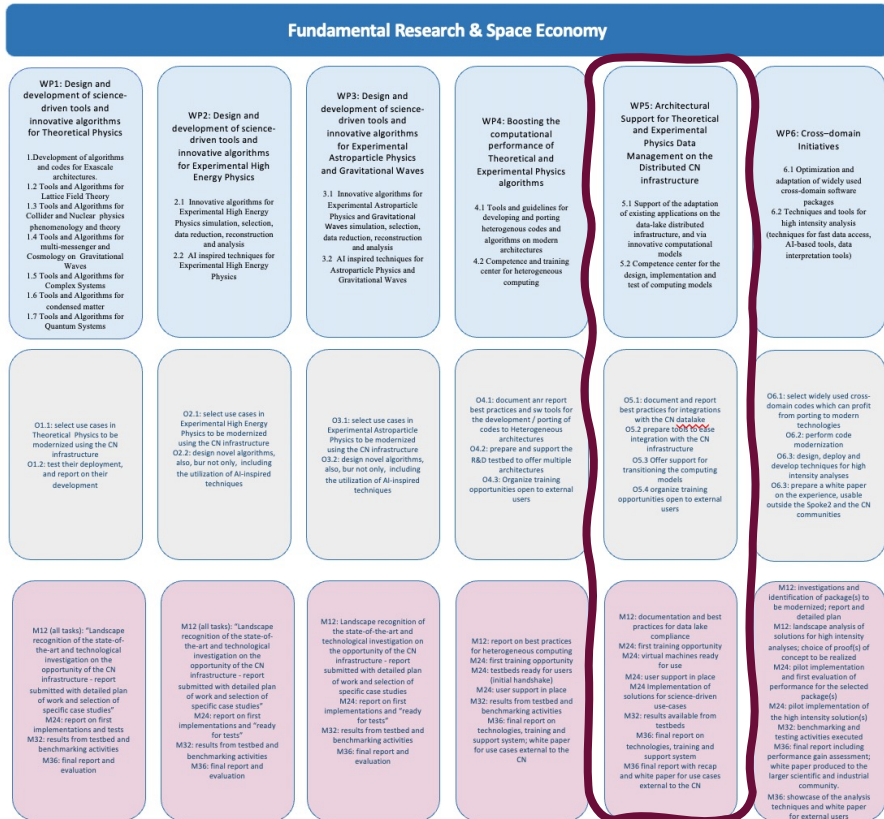
Contributors: INFN, UNIMIB, UNINA, ROMA1, UNITS, UNIBO, UNIPD, UNIFE

- Support of the adaptation of existing applications on the data-lake distributed infrastructure, and via innovative computational models
- Competence center for the design, implementation and test of computing models

# Spoke 2 - Fundamental Research & Space Economy

The activities in Spoke 2 "Fundamental Research and Space Economy" focus on boosting the science capabilities of current and future science initiatives.

The Spoke intends to address the needs of theoretical and experimental physics with accelerators, astroparticle physics with space- and ground-based detectors and gravitational wave investigation designing, developing and testing solutions apt to the current and next-generation experiments, and fitting the opportunities provided by the PNRR and the National Centre (CN) "Big Data, HPC and Quantum Computing".



**WP2.5**

*The crucial aspects: the creation and/or optimization of algorithms and, in general, computing solutions capable of maximizing the potential physics output from experimental data and theoretical and phenomenological simulations, by using the tools made available by the Centre: e.g., heterogeneous and high-performance computing (via standard programming and AI-based solutions) and the ability to process large quantities of data beyond the capabilities of traditional methods.*

*A common denominator will be the utilization of more efficient strategies, reducing the computational costs and their power consumption footprint ➔ the project aims to be a key player in developing and testing solutions which ensure the sustainability of computing for the next generation of scientific experiments.*

*All the activities in Spoke 2 will be executed in strict collaboration with the major players in the respective scientific domains; it is thus important to maintain active and frequent communication channels with them.*

# WP5: Architectural Support for Theoretical and Experimental Physics Data Management on the Distributed CN infrastructure

**Purpose:**

Support for the adaptation of existing applications on the data-lake distributed infrastructure, and via innovative computational models (*for example sharing of gauge configurations in lattice field theories, long-term data preservation, streaming access to data, tiered storage solutions, ...*). The solutions implemented will be tailored to the needs of the scientific fields, easily extendible not only to the nearby scientific domains in the Centre, but also to all academic and industrial realities where needs to access distributed computing and large amounts of data exist. In particular, the industrial partners in the Spoke have expressed interest in using the same technologies for their specific use cases.

**Objectives:**

O5.1: Document and report best practices for integrations with the CN data lake
O5.2: Prepare tools to ease integration with the CN infrastructure
O5.3: Offer support for transitioning the computing models
O5.4: Organize training opportunities open to external users

Similar approach with the respect to WP 2.4 and 2.6:

- ✓ *Planning and identification:* landscape recognition for best solutions for the realization of heterogeneous and portable code (e.g. software frameworks, compilers, programming models, ...), for the integration of services into a data-lake infrastructure; cross domain software and services will be identified if appropriate. Moreover, solutions for handling user support, user fora, and training opportunities will be identified;

- ✓ *Realization phase:* phase in which the services and the support systems are put into place, at least in alpha/beta phase. These include the testbeds to be used for benchmarking of scientific and industrial solutions, the user support system, the training opportunities.

- ✓ *Validation phase:* phase in which experience on the supported services and codes are reported, to be used as a touch base before the end of the project.

- ✓ *Wrap-up phase:* phase in which results are reported for executed activities, and are disseminated via white papers for future and external use cases.

*Realization and Validation phase:*

*High integration with the industrial partners:*

- ➢ industrial partners can provide testbed platforms, by offering their infrastructure in-kind, or procuring in their centres new R&D platforms not available in the CN;

- ➢ industrial partners are expected to test their typical-use cases on the same platforms, under the paradigm that "data are data", and once one can abstract from the specific domain, the technology to treat them is similar. This second part is expected to siphon experience and technologies between the academic and productive sectors, and vice versa.

*The testbed and benchmarking ("validation") phase will be executed partially via the "innovation grants" available via the project, in which the solutions developed in alpha/beta level will be tested together with industrial partners, on hardware either provided by the CN, or acquired via the same grants.*

# WP2.5: Approach and workflow

*Realization and Validation phase:*

▶ **Shared testbeds and proof of concept**, using the infrastructure of the CB and other R&D platforms procured as needed, in order to test the processing of large, dispersed and heterogeneous data sources. Solution to be tested tentatively include:
  - ▶ Test on processing advantage using smart/ dumb caches for remote/local data;
  - ▶ test using tiered storage systems (from tape to rotating disks to solid state disks);
  - ▶ test on remote streming vs lazy download vs caches w or w/o prefetch.

  ▶ **Use cases:**
  - ▶ processing of O(10-100) TB of data from collider experiments, for example using workflows such as data processing in HEP;
  - ▶ typical data-intensive-use cases from companies, such ad the processing of agricultural data, as suggested by Intesa Sanpaolo.

▶ Metrics to be considered for the testbeds are (at least) the processing efficiency / total time and power consumption using the various storage solutions, and the cost and scalability of the storage systems.

▶ All the activities will be executed in strict collaboration with the major players in the respective scientific domains; it is thus important to maintain active and frequent communication channels with them.

▶ **Open Calls:** number of support and ancillary services, like a web portal, a ticketing / support system, and help in organizing activities like benchmarking and training. We intend to use a part of the "Open Calls" to this purpose, selecting professionalities from companies or academic institutions with experience on the subjects from previous projects.

# WP2.5: Architectural Support for Theoretical and Experimental Physics Data Management on the Distributed CN infrastructure

## Milestones

**M12** — Documentation and best practices for data lake compliance

**M24**
- M24: first training opportunity
- M24: virtual machines ready for use
- M24: user support in place
- M24 Implementation of solutions for science-driven use-cases

**M32** — M32: results available from testbeds

**M36**
- M36: final report on technologies, training and support system
- M36 final report with recap and white paper for use cases external to the CN

# WP2.5: People

Very Very preliminary!!!
Only first declaration of interest not yet completed!

| INFN | Tommaso | Boccali | Ricercatore II Livello | WP2, WP4, WP5 |
|------|---------|---------|------------------------|---------------|
| INFN | Lucia | Silvestris | Ricercatore II Livello | WP2, WP4, WP5 |
| INFN | Stefano | Bagnasco | Tecnologo II Livello | WP3, WP4, WP5 |
| INFN | Daniele | Spiga | Tecnologo III Livello | WP2, WP4, WP5 |
| INFN | Alessandro | De Salvo | Tecnologo II Livello | WP2, WP5 |
| INFN | Domenico | Elia | Ricercatore II Livello | WP3, WP4, WP5 |
| UNINA | Guido | Russo | PO | WP5 |
| UNINA | Elvira | Rossi | RTDB | WP5, WP2 |
| UNIMIB | Mattia | Bruno | RTD-B | WP1,WP4,WP5 |
| UNIFE | Eleonora | Luppi | PO | WP2,WP5 |
| UNITS | Andrea | Bressan | PA | WP2,WP4,WP5 |
| UNITS | Giuseppe | Della Ricca | PO | WP2,WP5 |

Please fill the google doc with your interest and/or ask your affiliate to fill it

Don't be shy...
Anyone is very welcome to contribute and to join WP2.5

Contact me: elvira.rossi@unina.it; elvira.rossi@cern.ch;
Skype user: ellyufetto

# Next steps

▶ *As soon as possibile (before summer break):* Create a e-mail list of all the people interested in contributing to WP2.5

▶ *Estabilish a first meeting in September:*

  ▶ to get to know each other better

  ▶ to make brain storming and share past expertice

  ▶ to start giving ideas on how to proceed, define a preliminary workflow and start to distribute tasks

▶ *Set regular meetings to update and be updated on the ongoing work*

▶ *Set regular interactions with other WPs*

▶ *Have a very strict and regular contact with Industries*

# Backup

## Past experiences and responsibilities relevant to the spoke activities:

➢ physics and computing coordination for the most computing-intensive experiments;
➢ experience in obtaining large computing grants on HPC systems from international bodies;
➢ project coordination of multi-million research projects, including design, construction, operations and analysis of Exabyte-scale datasets;
➢ design and operation of global scale e-Infrastructures serving science and industry, also via the coordination of multi-million Europe funded projects.

### Rucio:

Rucio is a project that provides services and associated libraries for allowing scientific collaborations to manage large volumes of data spread across facilities at multiple institutions and organisations. Rucio was originally developed to meet the requirements of the high-energy physics experiment ATLAS, and now is continuously extended to support the LHC experiments and other diverse scientific communities.

Rucio offers advanced features, is highly scalable, and modular. It is a data management solution that covers the needs of different communities in the scientific domain (e.g., HEP, astronomy, biology). Below are some resources to help you get you started on your journey.