



Introduzione a INFN DataCloud

Barbara Martelli

(many slides courtesy of
Giacinto Donvito, Claudio
Grandi, Davide Salomoni)

DataCloud: l'infrastruttura di calcolo nazionale

- Nata come WG del Comitato Coordinamento Calcolo INFN, con l'obiettivo di armonizzare e integrare l'infrastruttura tradizionale «a Tier» tipica di WLCG e il modello «cloud-native»
 - Integrazione: di risorse, metodi, persone, soluzioni
- Utilizzata e in fase di evoluzione nell'ambito di ICSC e TeRABIT, progetti PNRR che prevedono la creazione di una infrastruttura cloud nazionale
- Lo scopo è accedere alle risorse in modo trasparente ed efficiente
- Gli attori principali sono: INFN, CINECA, GARR
- Ma anche: CMCC, ENEA, SISSA, IIT, UniTO, Sapienza e altre universita' e centri di ricerca



...un po' di storia

The development of **Scientific Computing within INFN** was driven by the needs of its own theoretical and experimental communities; however, being at the forefront of computing in research seeded **many projects with a much broader scope.**



“preparing the GRID”

“preparing the Cloud”

“expanding beyond HEP”

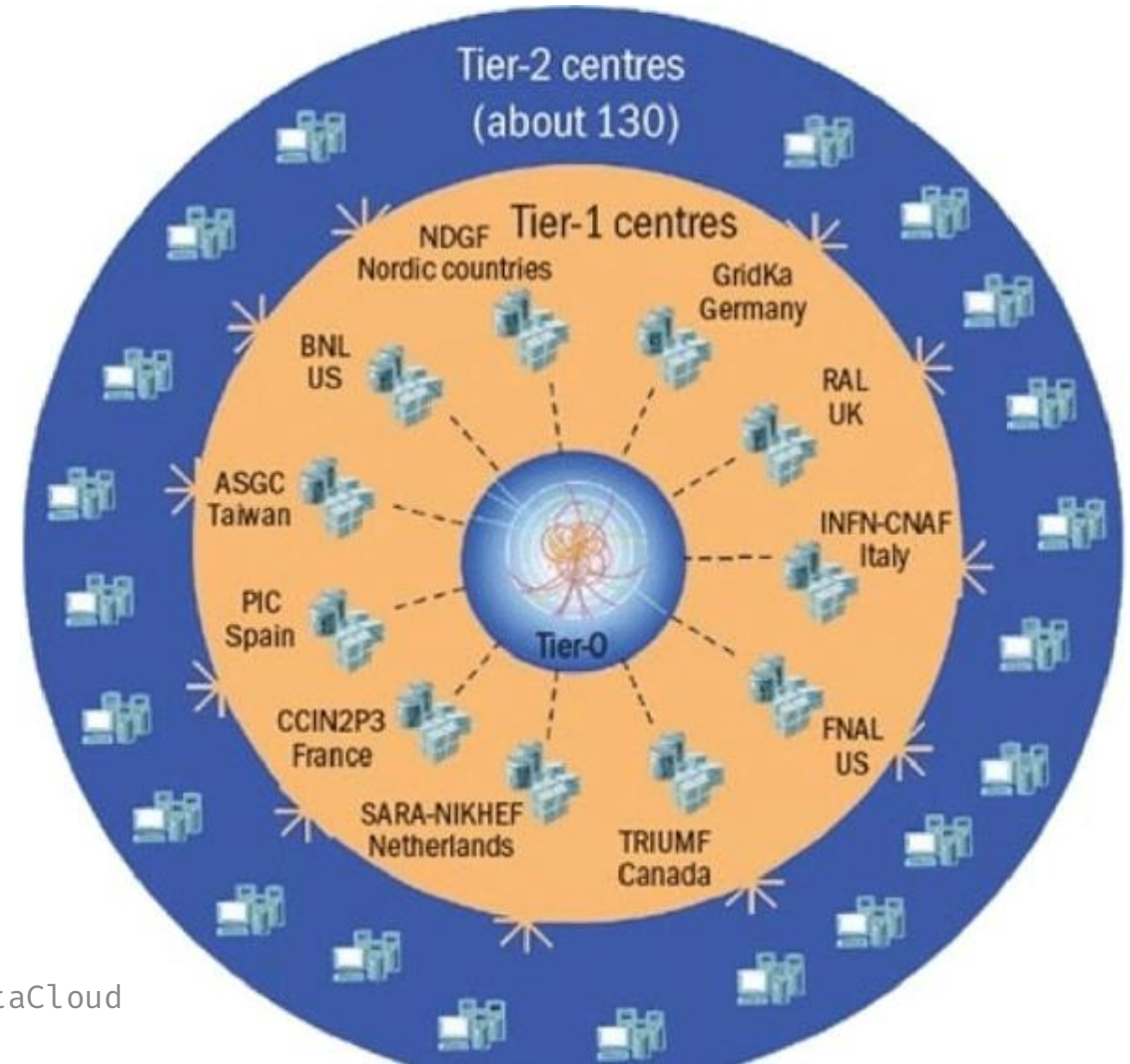
World Wide LHC Computing Grid



The MONARC model

“A model of large-scale distributed computing based on many regional centers, with a focus on LHC experiments at CERN. As part of the MONARC project, a simulation framework was developed that provides a design and optimisation tool.

The MONARC model has been the initial reference for building the WLCG infrastructure and to organise the data transfers around it.”



The HEP Software Foundation., Albrecht, J., Alves, A.A. *et al.* A Roadmap for HEP Software and Computing R&D for the 2020s. *Comput Softw Big Sci* 3, 7 (2019).
<https://doi.org/10.1007/s41781-018-0018-8>

Cloud computing essential characteristics



- **On-demand self-service.** A consumer can **unilaterally provision** computing capabilities, such as server time and network storage, as needed **automatically without requiring human interaction** with each service provider.
- **Broad network access.** Capabilities are available over the network and **accessed through standard mechanisms** that promote use by heterogeneous thin or thick client platforms (e.g., mobile phones, tablets, laptops, and workstations).
- **Resource pooling.** The provider's computing **resources are pooled to serve multiple consumers using a multi-tenant model**, with different physical and virtual resources **dynamically assigned and reassigned according to consumer demand**. There is a sense of location independence in that the customer generally has no control or knowledge over the exact location of the provided resources but may be able to specify location at a higher level of abstraction (e.g., country, state, or datacenter). Examples of resources include storage, processing, memory, and network bandwidth.
- **Rapid elasticity.** Capabilities can be **elastically provisioned and released**, in some cases automatically, **to scale rapidly outward and inward** commensurate with demand. To the consumer, the capabilities available for provisioning often appear to be unlimited and can be appropriated in any quantity at any time.
- **Measured service.** Cloud systems automatically control and optimize resource use by leveraging a metering capability at some level of abstraction appropriate to the type of service (e.g., storage, processing, bandwidth, and active user accounts). Resource usage can be monitored, controlled, and reported, providing transparency for both the provider and consumer of the utilized service.

“The NIST definition of Cloud Computing”, 2011

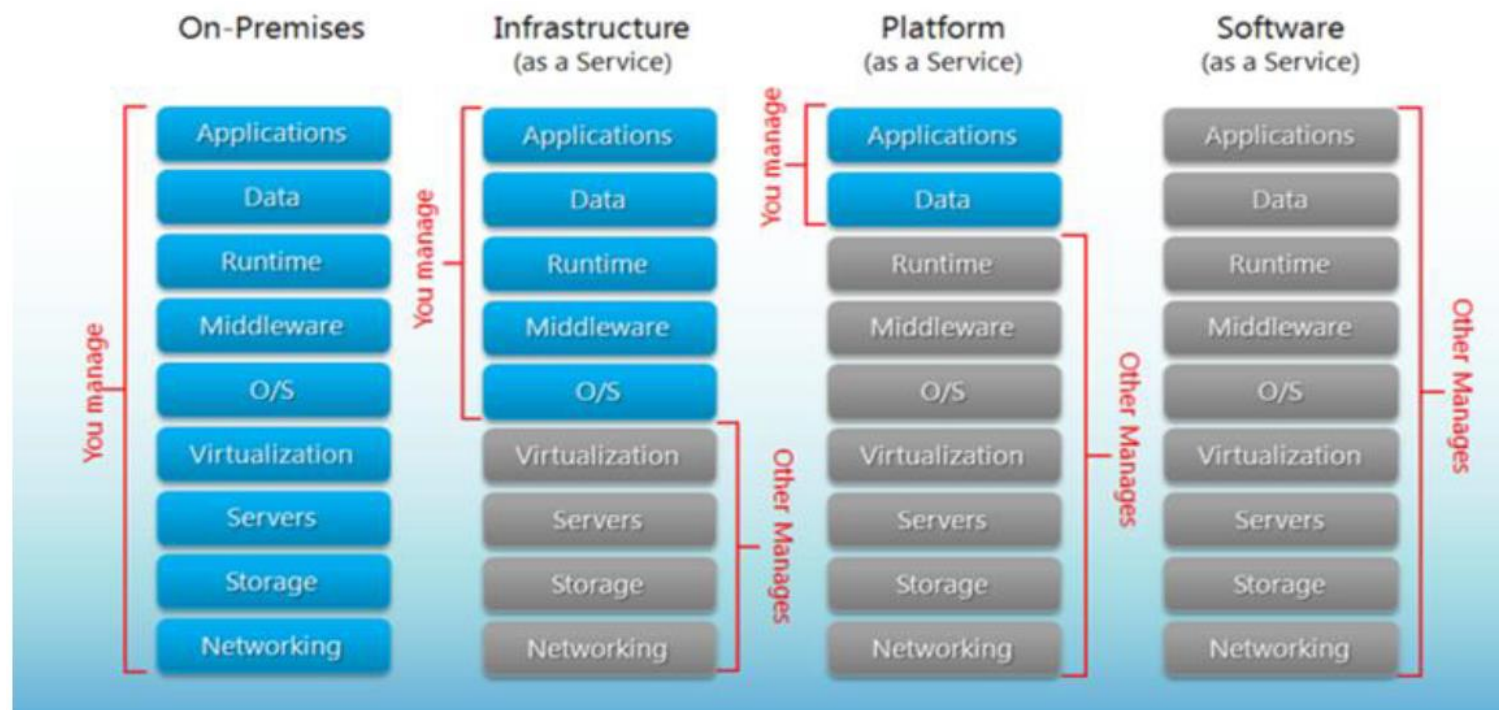
<https://nvlpubs.nist.gov/nistpubs/legacy/sp/nistspecialpublication800-145.pdf>

Cloud computing service and shared responsibility models



- **Software as a Service (SaaS):** The capability provided to the consumer is to use the provider's **applications** running on a cloud infrastructure
- **Platform as a Service (PaaS):** The capability provided to the consumer is to **deploy** onto the cloud infrastructure consumer-created or acquired **applications created using programming languages, libraries, services, and tools supported by the provider**. The consumer [...] has control over the deployed applications and possibly configuration settings for the application-hosting environment
- **Infrastructure as a Service (IaaS):** The capability provided to the consumer is to **provision processing, storage, networks, and other fundamental computing resources** where the consumer is able to deploy and run arbitrary software, which can include operating systems and applications

“The NIST definition of Cloud Computing”, 2011
<https://nvlpubs.nist.gov/nistpubs/legacy/sp/nistspecialpublication800-145.pdf>



Cullum, Paul. (2020). **A SURVEY OF THE HOST HYPERVISOR SECURITY ISSUES PRESENTED IN PUBLIC IAAS ENVIRONMENTS AND THEIR SOLUTIONS.** International Journal of Engineering Applied Sciences and Technology. 5. 10.33564/IJEAST.2020.v05i08.006.

Cloud computing deployment models



- **Private cloud:** The cloud infrastructure is provisioned for **exclusive use by a single organization** comprising multiple consumers (e.g., business units). It may be owned, managed, and operated by the organization, a third party, or some combination of them, and it may exist on or off premises.
- **Community cloud:** The cloud infrastructure is provisioned for exclusive use by a specific community of consumers from organizations that have shared concerns (e.g., mission, security requirements, policy, and compliance considerations). It may be owned, managed, and operated by one or more of the organizations in the community, a third party, or some combination of them, and it may exist on or off premises.
- **Public cloud:** The cloud infrastructure is provisioned for open use by the **general public**. It may be owned, managed, and operated by a business, academic, or government organization, or some combination of them. It exists on the premises of the cloud provider.
- **Hybrid cloud:** The cloud infrastructure is a **composition of two or more distinct cloud infrastructures (private, community, or public)** that remain unique entities, but are bound together by standardized or proprietary technology that enables data and application portability (e.g., cloud bursting for load balancing between clouds).

“The NIST definition of Cloud Computing”, 2011

<https://nvlpubs.nist.gov/nistpubs/legacy/sp/nistspecialpublication800-145.pdf>

15 marzo 2021: INFN Cloud entra in produzione <https://www.cloud.infn.it/>



- The **starting point** for a National Datalake for research and beyond, building on (existing|renewed|new) e-Infrastructures.
- The **base of the evolution** of the INFN Distributed Computing vision.
- Built on a **thin middleware layer** running on top of federated clouds, decoupling physical and logical views via a **service composition** mechanism.

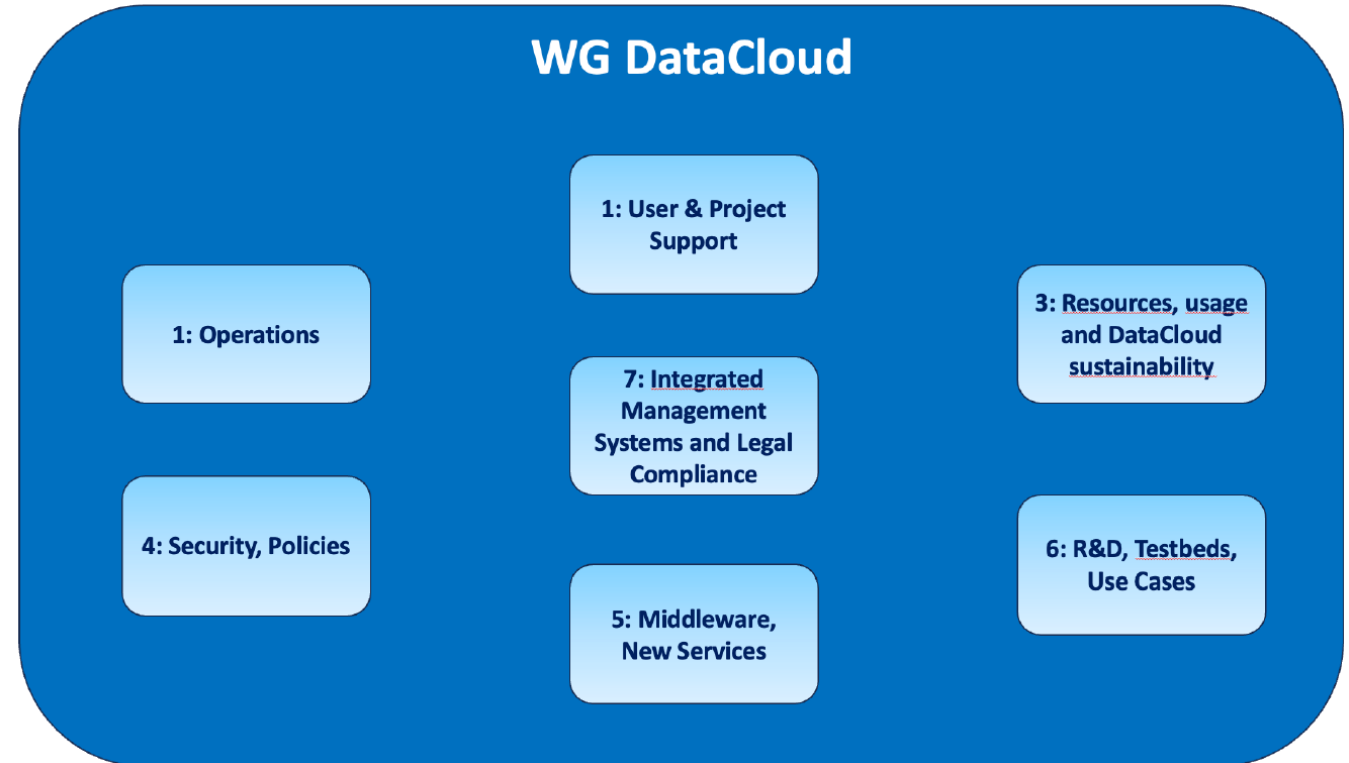


TDR DataCloud (nov 2023)



Governance:

- Coordinatore DataCloud
- Deputy
- Chief Technology Officer (CTO)
- Program Officer (PO)
- Project Management Board (PMB) include i coordinatori di tutti i WP



INFN and Computing / Big Data



- Agli inizi degli anni 2000 una decina di datacenter di INFN furono selezionati per entrare a far parte della «**Worldwide LHC Computing Grid**» (WLCG)
- In Italia è stato costruito un **nuovo grande datacenter a Bologna**: il CNAF (in **rosso** nella figura)
- A questo centro sono stati aggiunti altri **9 centri più piccoli** in varie città italiane (in **giallo** nella figura)
- Da allora INFN è stato pioniere in ambito **Grid** e **Cloud Computing**, creando l'infrastruttura nazionale **HPC-BD-AI infrastructure** (High-Performance Computing / Big Data / Artificial Intelligence)
- Questi centri INFN sono ancora operativi e hanno aumentato le loro dimensioni di circa 100 volte. La loro connettività di rete è oggi di multipli di 100 Gigabit/s
- Complessivamente l'infrastruttura di calcolo distribuito di INFN offre **~150.000 CPU cores, ~120PB di spazio disco e ~120 PB di spazio nastro**



DataCloud = GRID + INFN Cloud

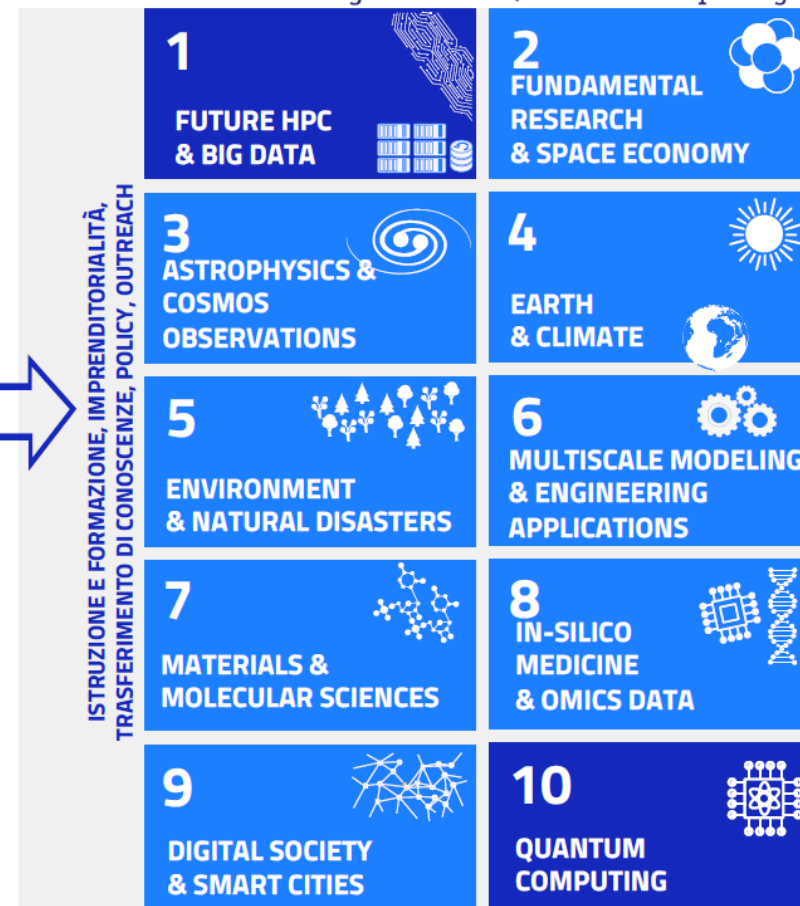
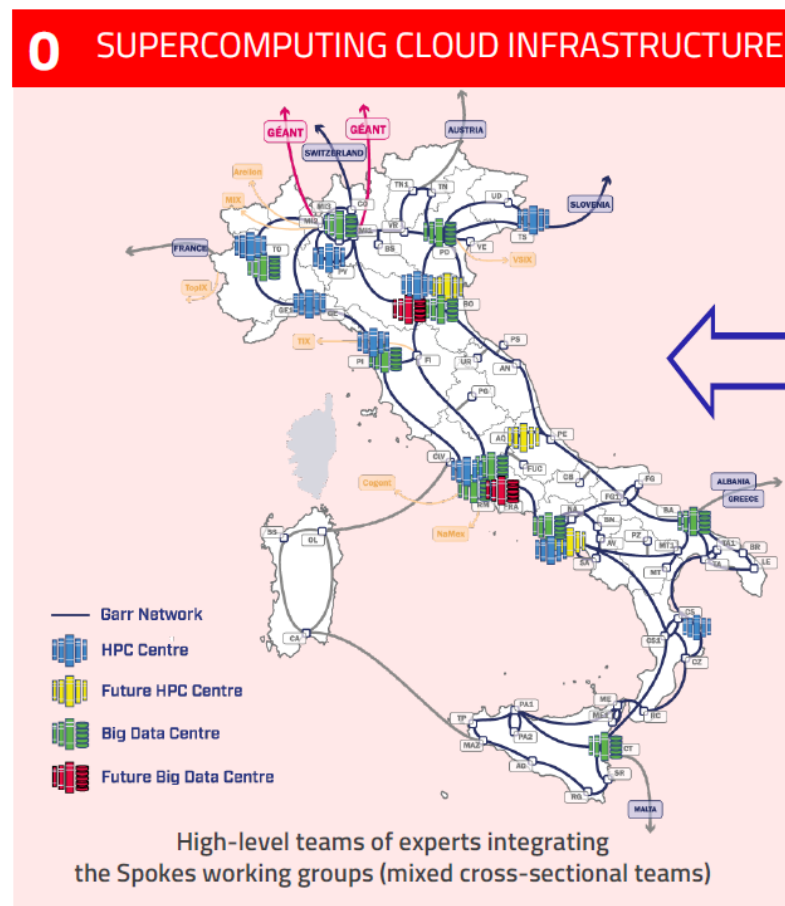
L'evoluzione di DataCloud attraverso i progetti PNRR



Il centro nazionale HPC, BD e QC



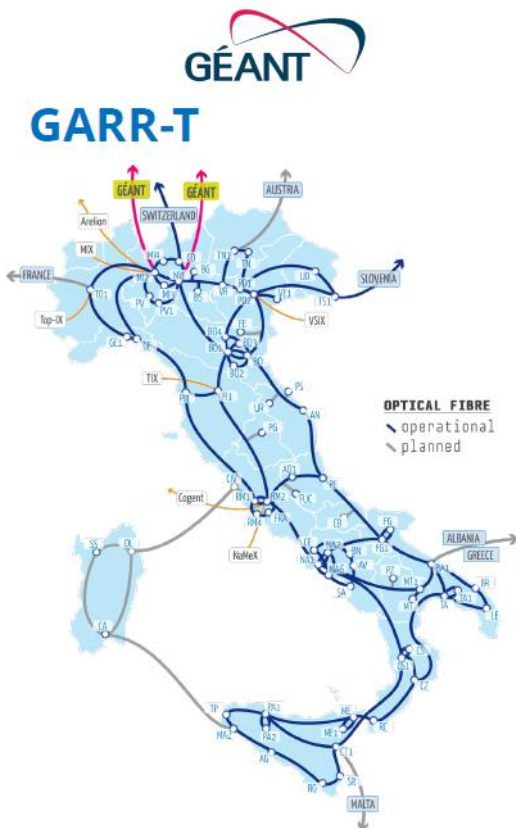
ICSC include
10 spokes tematici
1 spoke infrastruttura



L'evoluzione di DataCloud attraverso i progetti PNRR



Le infrastrutture di ricerca di TeRABIT



PRACE-Italy



Galileo100 – cluster HPC
Ospitato al CINECA - Bologna

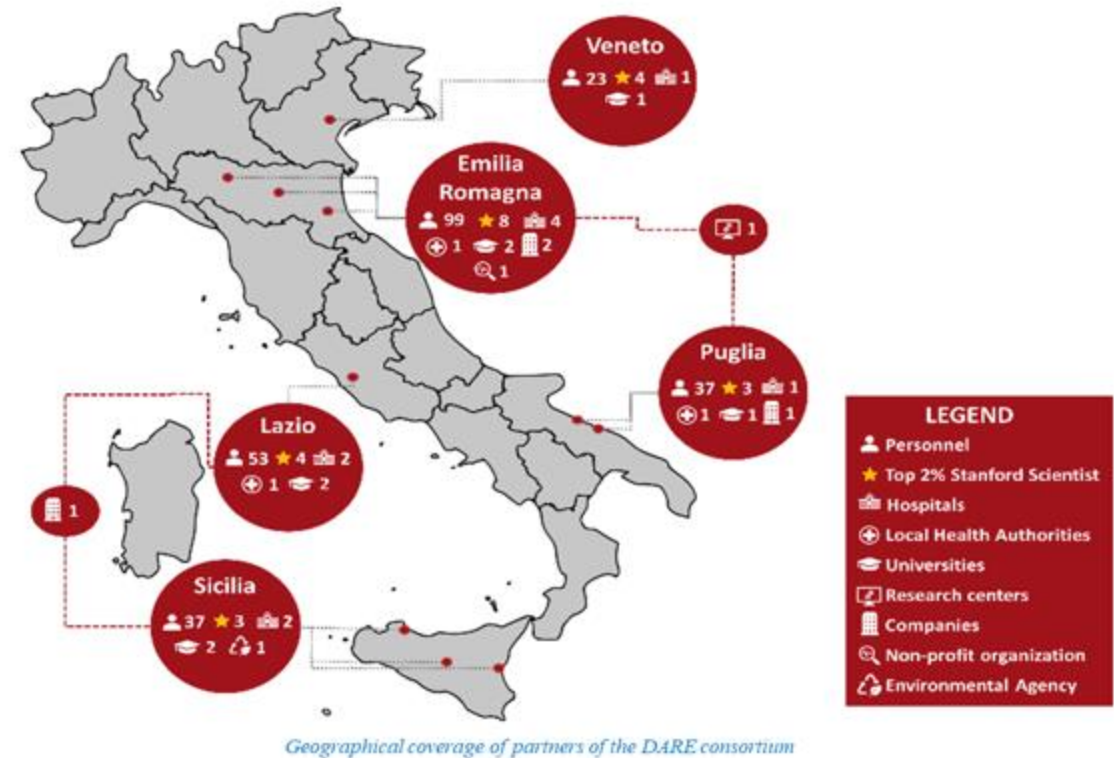
HPC-BD-AI - INFN Cloud



L'evoluzione di DataCloud attraverso i progetti PNRR



DARE: A public-private partnership aimed at studying, coordinating and implementing innovative digital health solutions in Italy with a specific focus on primary, secondary and tertiary prevention: from mitigation of health risks to disease monitoring



DataCloud è la piattaforma di calcolo che metterà a sistema le risorse DARE dislocate nelle varie regioni

ARCHITECTURAL FOUNDATIONS



NO VENDOR LOCK-IN

Open-source,
vendor-neutral
architecture



FEDERATION

of existing Cloud
infrastructures for
both compute
and data



DYNAMIC ORCHESTRATION

of resources via
the INDIGO PaaS
Orchestrator



CONSISTENT AUTHN/AUTHZ

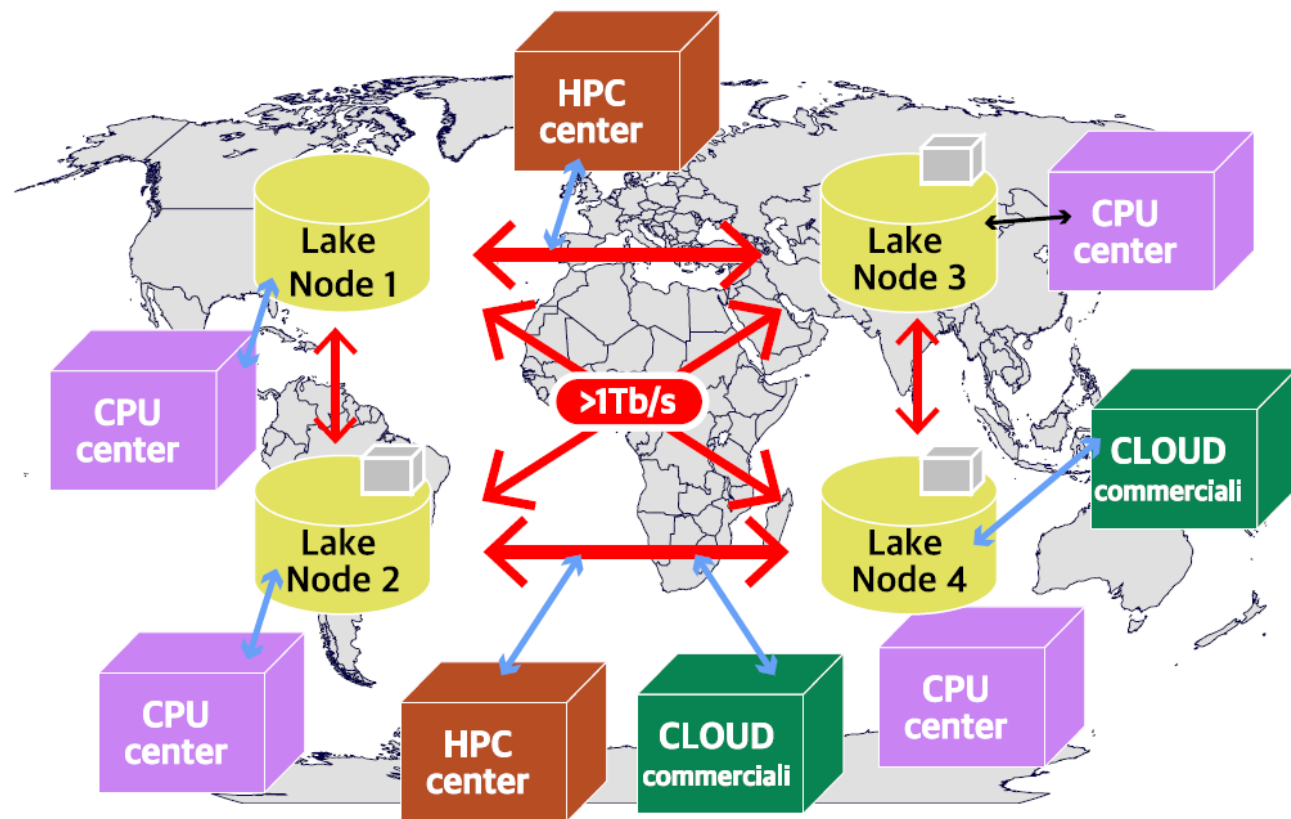
at all cloud levels
via OpenID-
Connect/OAuth2

Modello data-centrico

Disaccoppiamento di storage e CPU

Nodi storage interconnessi tramite una rete a banda larga

Nodi eterogenei possono accedere ai dati in modo indipendente dalla locazione

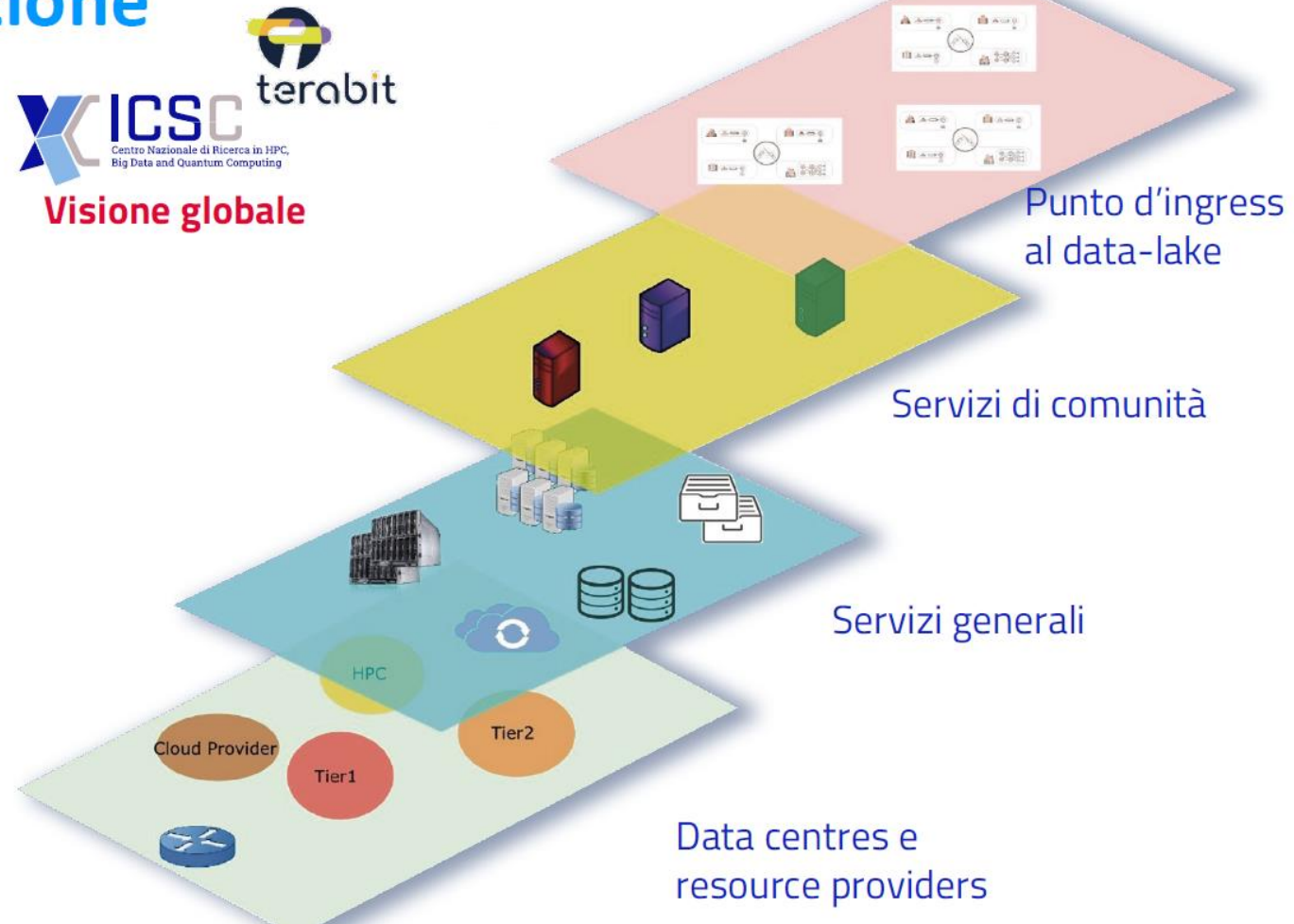


I principi della federazione

Inclusività, attraverso una federazione "leggera" e l'adozione di standard

Facilità d'uso, attraverso un orchestratore PaaS e una dashboard

Flessibilità, grazie a meccanismi ibridi di allocazione delle risorse



Inclusività



La federazione includerà data centres che sono già in produzione, e parte di comunità internazionali

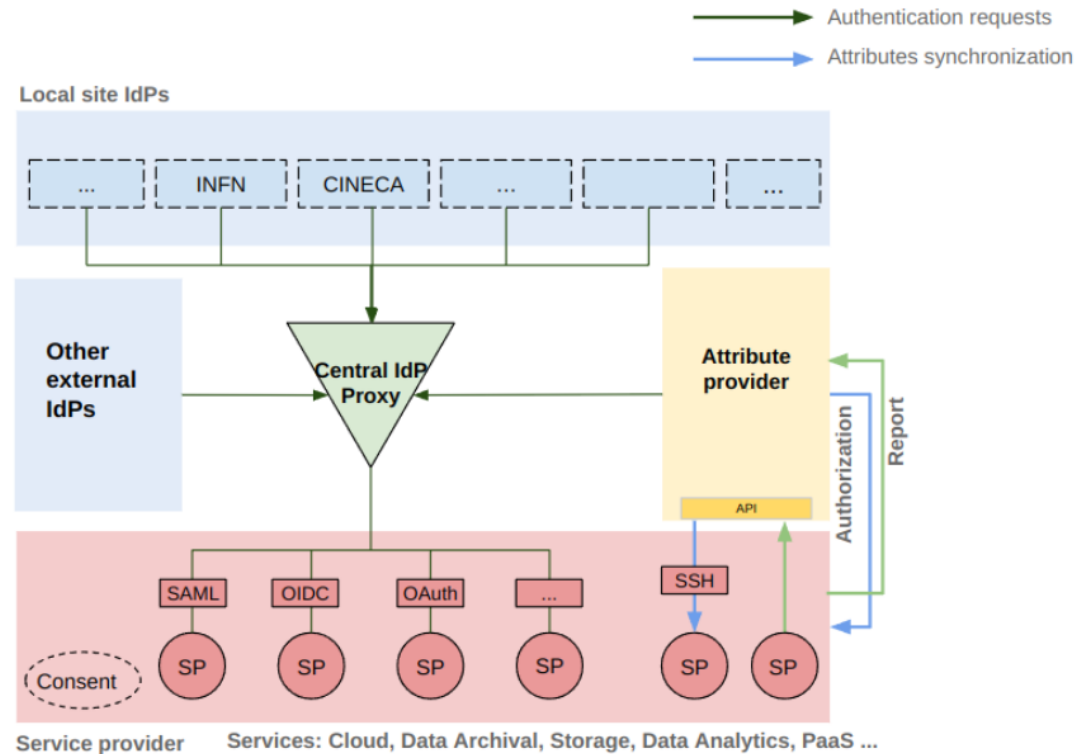
Le procedure per entrare nella federazione devono essere non intrusive

Bisogna usare standard quando possibile, e sviluppati quando non ci sono

La federazione servirà utenti di diverse organizzazioni in diversi campi

Le procedure di ingress devono essere il più semplice possibile

Ad esempio tramite l'uso di federazioni di identità

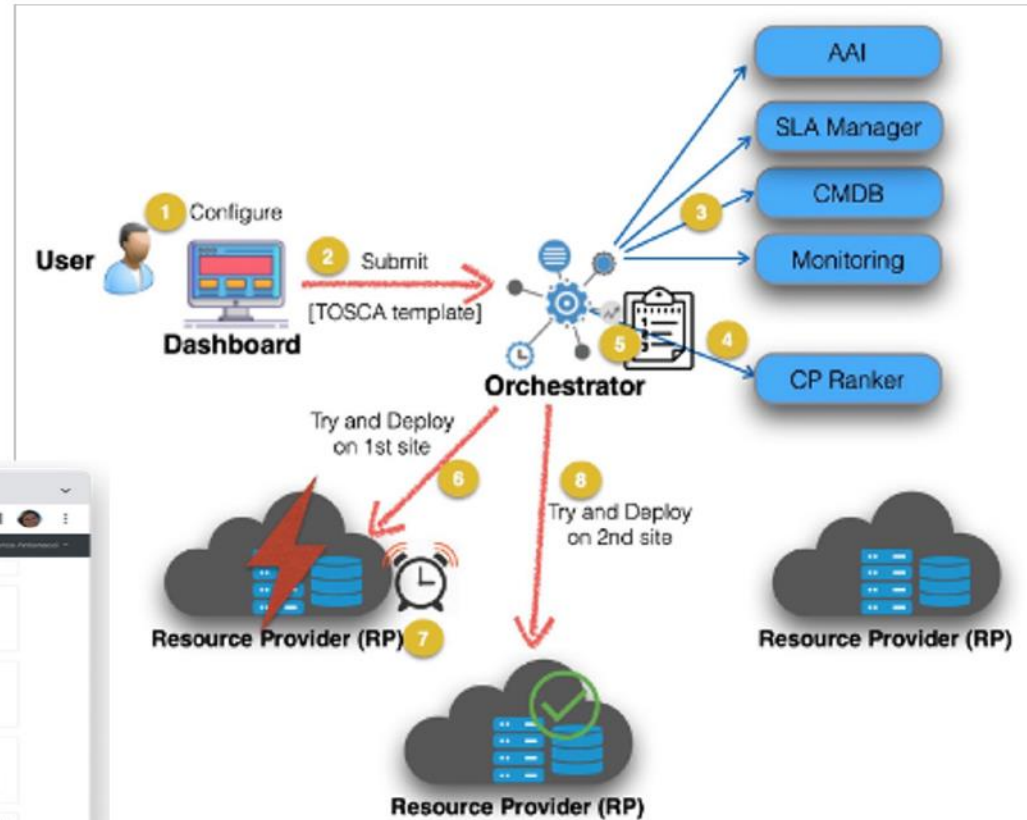
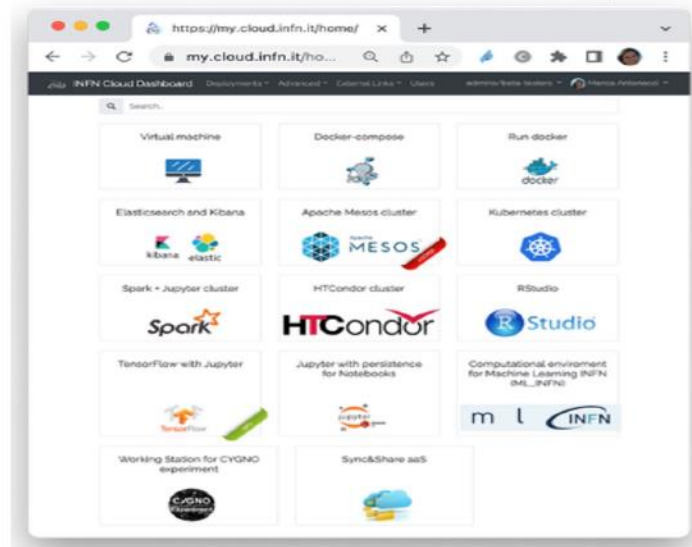


Facilità d'uso

La federazione servirà utenti con diverse competenze sul calcolo

Nascondere all'utente finale la complessità dell'infrastruttura sottostante

Esperti dei diversi campi sviluppano piattaforme che consentono l'uso efficace delle infrastrutture attraverso la composizione di servizi e risorse



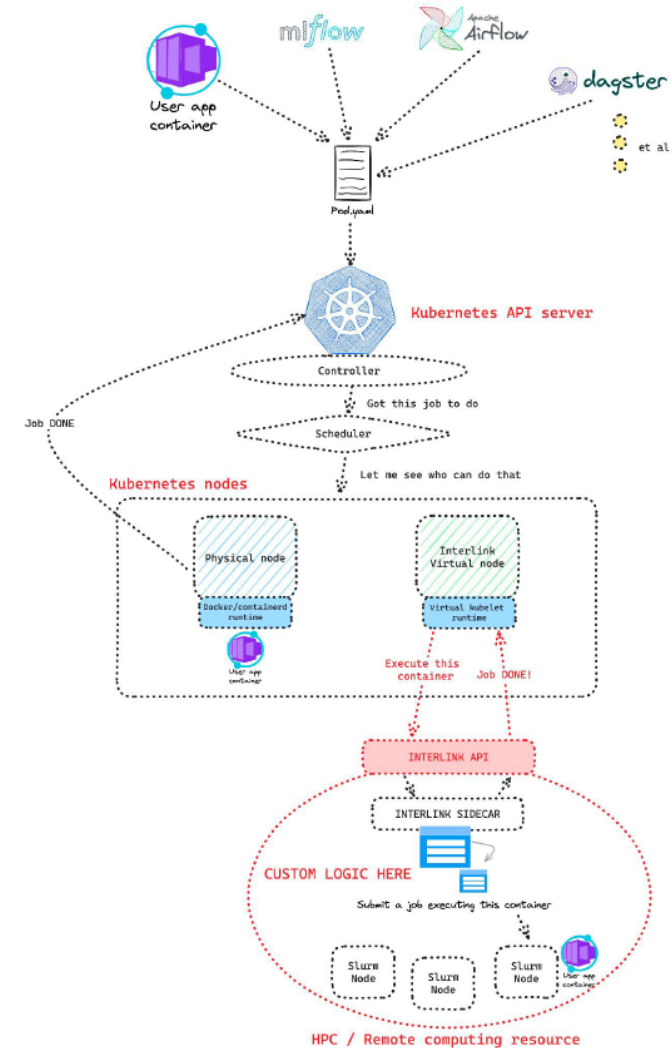
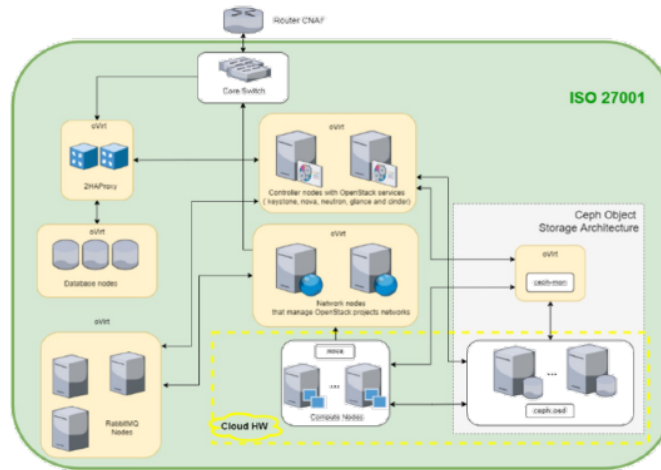
Flessibilità

Supporto a diversi metodi di accesso alle risorse, con attenzione a:

- a. Trasparenza e facilità d'uso
- b. Efficienza ed efficacia

Supporto a requisiti specifici delle applicazioni

Ad es. Piattaforme con requisiti particolari sulla privacy

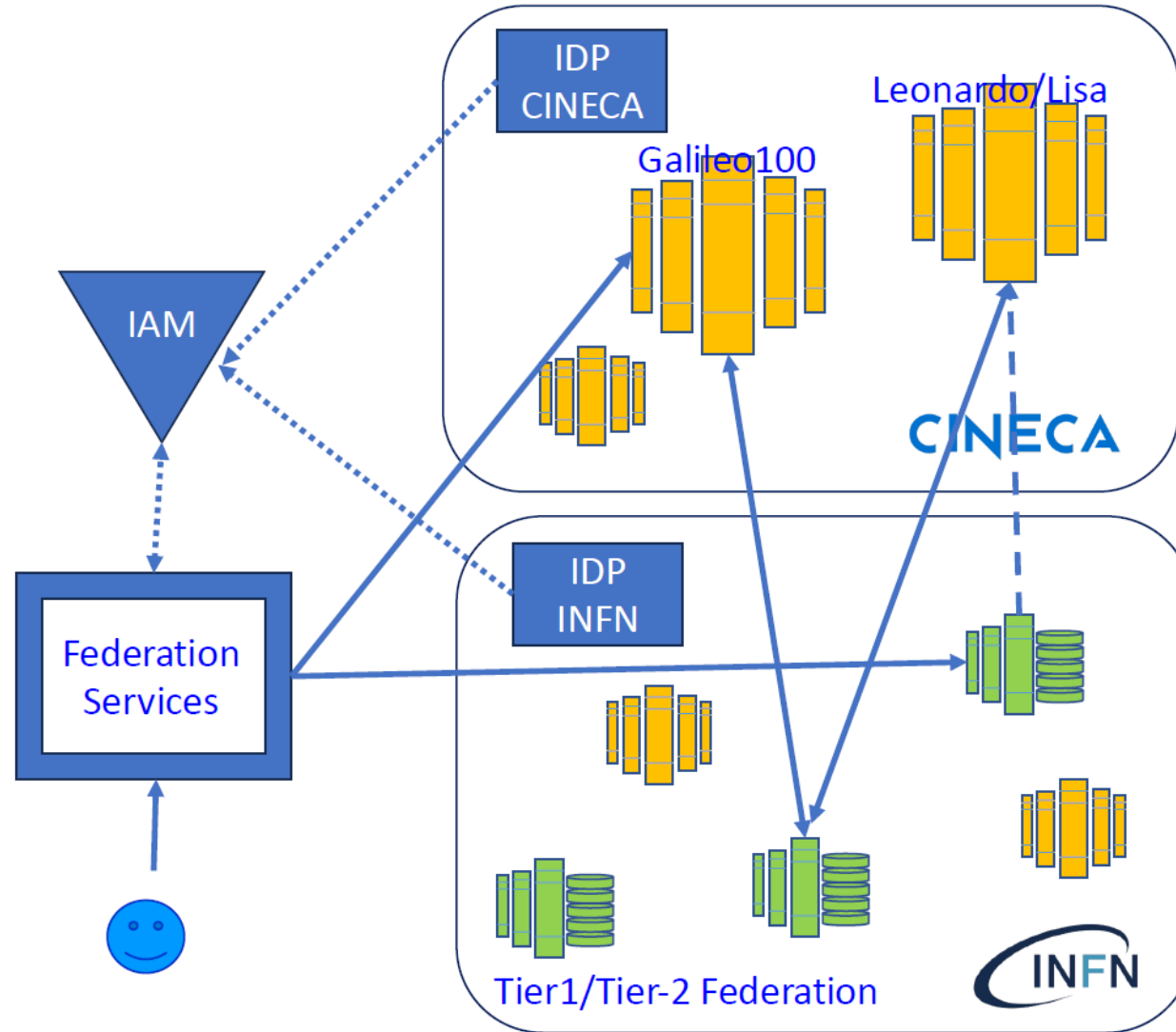


I piani

Giugno 2024: PoC con funzionalità limitate e utenti «cavia»

Giugno 2025: architettura adottata sull'infrastruttura di produzione con funzionalità via via crescenti

Nel frattempo: accesso alle risorse con meccanismi specifici delle infrastrutture



Le sfide

Autenticazione e autorizzazione federata

Costruire il *trust* fra istituzioni e con le federazioni di identità

Definire i meccanismi per comunicare il livello di *assurance* delle operazioni di identificazione

REFEDS Assurance Framework

Adottare sugli strumenti della federazione e nelle istituzioni questi meccanismi

Accesso alle risorse di calcolo

Utilizzo di soluzioni standard e supporto da parte degli strumenti di federazione

OpenStack, Kubernetes, ...

Accesso allo storage e gestione dei dati

Supporto a diversi tipi di storage (S3, ...) oggi in produzione

Mantenere efficienza nella comunicazione fra storage e CPU

Tutto questo su infrastrutture in produzione!



EPIC Cloud

Enhanced Privacy and Compliance Cloud



Enhanced Privacy and Compliance Cloud is an ISO certified cloud platform

A region of INFN Cloud with a certified Information Security Management System



EPIC Cloud offers an IaaS Community Cloud for the communities of

Biomedical and genomic researchers
Industrial researchers



Site locations: CNAF (active now), Bari and Catania sites will be added in Feb 2025 enabling for high availability and disaster recovery



Resource available today: about 1,5 PB HDD 600 TB SSD of storage, 1440 cores, 10 TB RAM, 6 GPU A100
On going expansion with 3M euro of NRRP resources and 4M euro of funds from other projects



Lo scope del certificato EPIC Multisito



Co-progettazione, sviluppo e manutenzione di soluzioni software DataCloud per il settore della ricerca.

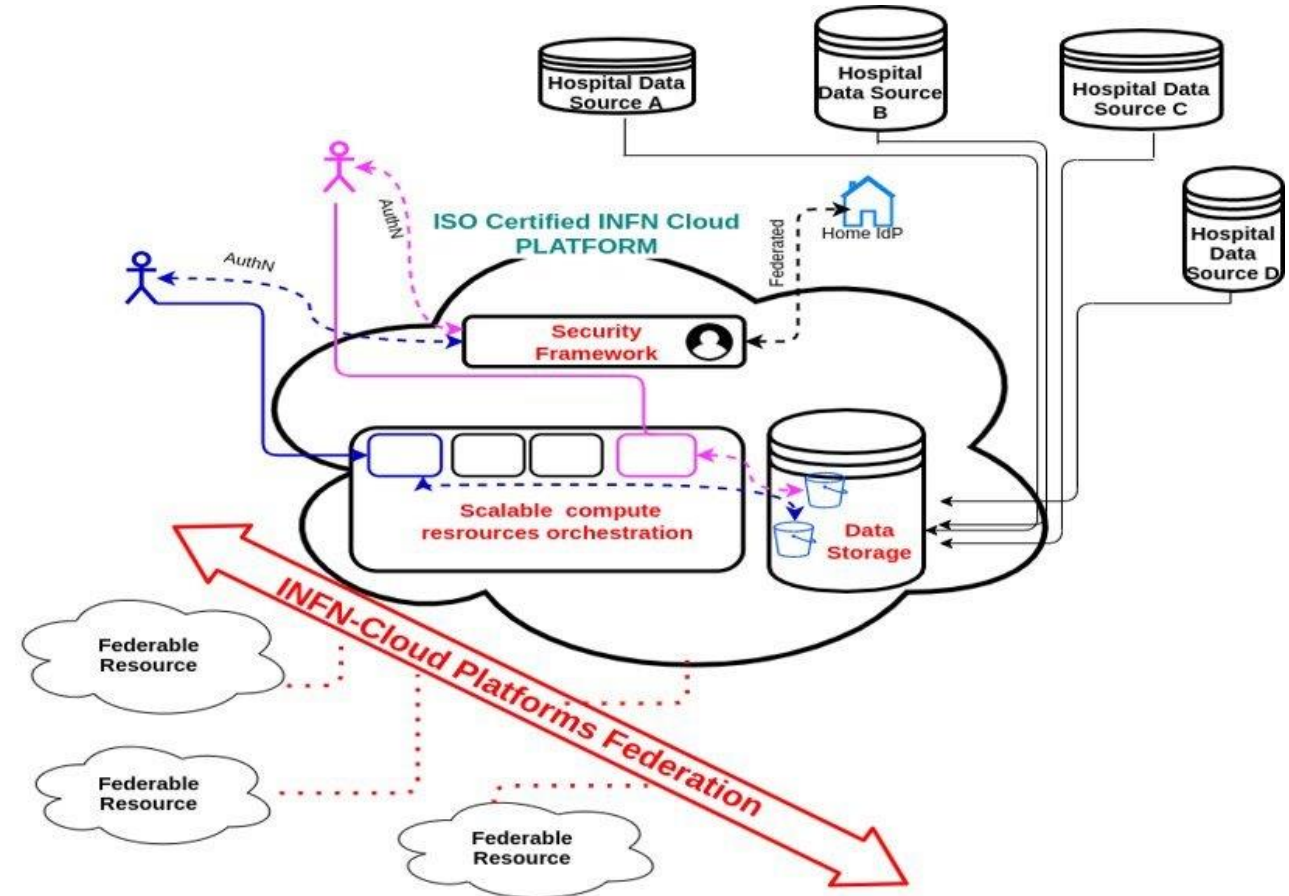
Erogazione di servizi DataCloud IaaS, SaaS e PaaS con modello di deployment Community Cloud.

Il DataLake ed i servizi di EPIC Cloud Multisito



Cosa stiamo costruendo:

- Una piattaforma di **genomica computazionale** basata su BOSCO di IRCCS AOU di Bologna (Sant'Orsola)
- Una piattaforma **Molecular Tumor Board** sviluppata da Health Big Data
- Una piattaforma di **radiomica** sviluppata da INFN Next-AIM
- Prototipi di **HPC-Bubble certificate** per utilizzare l'HPC via cloud
- Prototipi di **piattaforme IoT** per la raccolta di dati da sensori e dispositivi *edge*



[https://www.physicamedica.com/article/S1120-1797\(21\)00320-3/fulltext](https://www.physicamedica.com/article/S1120-1797(21)00320-3/fulltext)

Buon Lavoro!

