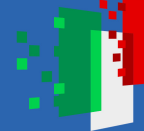




Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



Vision & Roadmap of the DataCloud Federation

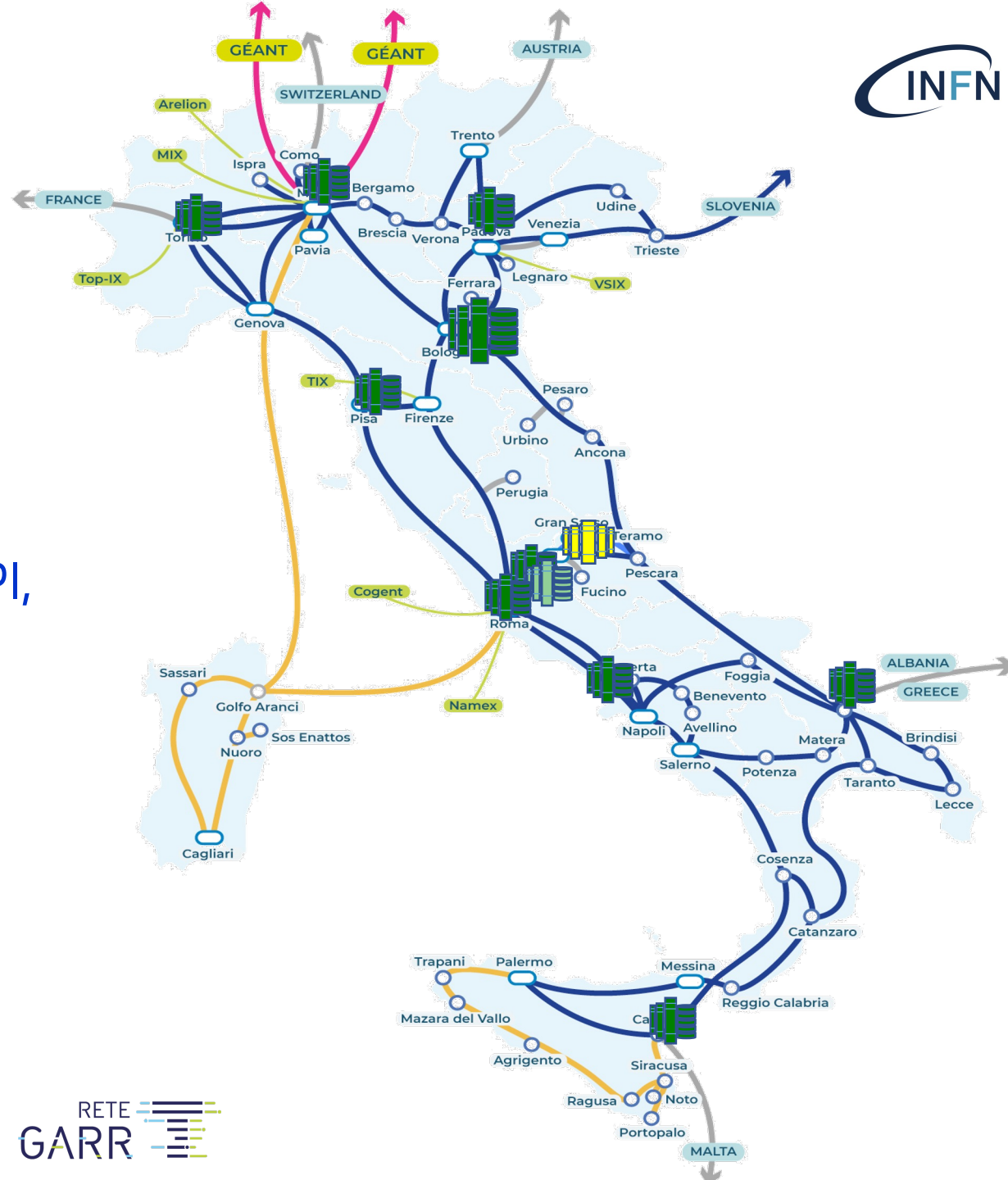
Claudio Grandi

Workshop sul calcolo INFN - Palau



DataCloud is the Infrastructure for INFN Scientific Computing

- Tier-1 (CNAF)
- Tier-2's (BA, CT, LNF, LNL/PD, NA, MI, PI, RM1, TO)
- INFN Cloud
 - Backbone and federated clouds
- HPC4DR (LNGS)
- (Tier-3)



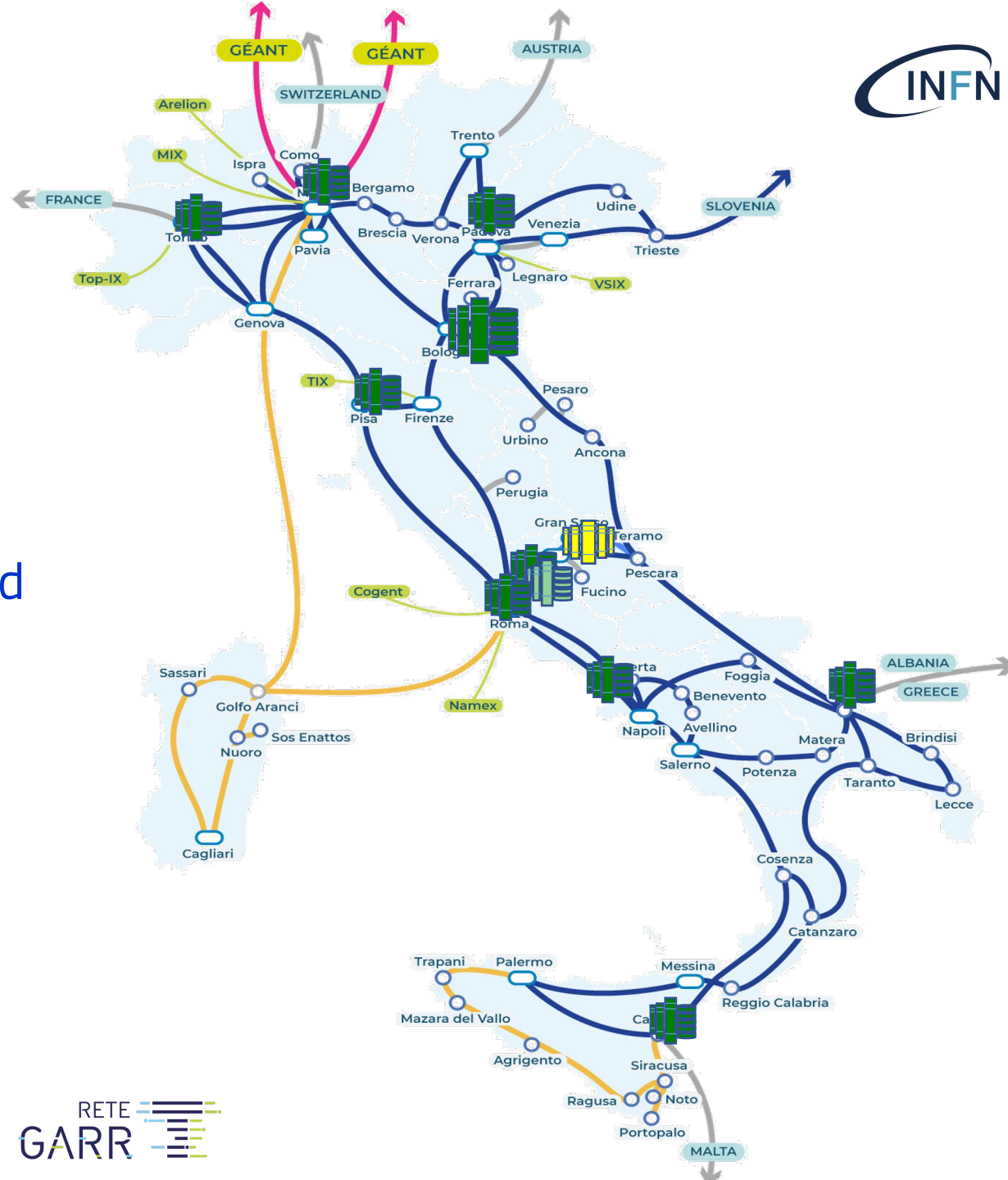


DataCloud addresses the needs of INFN research projects

- Internal projects: from CSN's
- External projects: regional, national and international projects, collaborations

The competences developed in the past years have brought to INFN visibility at national and international level

External projects have become more important





DataCloud is evolving into a Cloud Federation

Following the INFN Cloud model, resources are being made available through Cloud interfaces

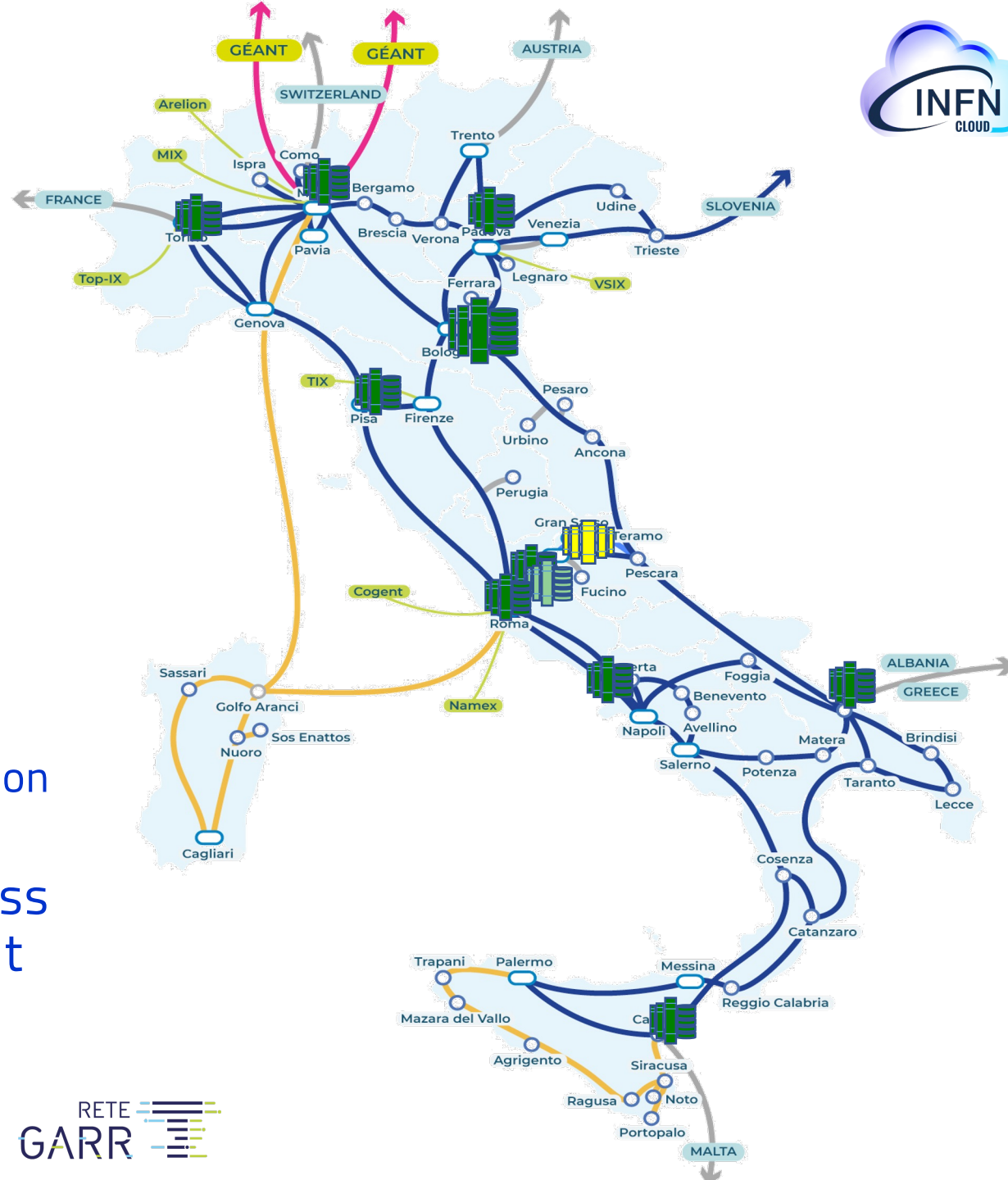
Inclusivity, through a lightweight federation model and the adoption of standards

Ease of use, through the PaaS orchestrator and dashboard

Flexibility, thanks to hybrid resource allocation mechanisms

Traditional (Grid and batch system) access remains as needed and when convenient

E.g. through Virtual Kubelets, ...





DataCloud is the basis for the Italian Cloud Federation

In the framework of the current NRRP projects, in particular ICSC and TeRABIT, INFN has a leading role in the creation of the Italian Cloud Federation

The goal is to access all Italian scientific computing resources through uniform interfaces

Main players: INFN, CINECA, GARR

But also: CMCC, ENEA, SISSA, IIT, UniTO, Sapienza, ...



Inclusivity

The federation will include data centres that are already in production, and part of international communities

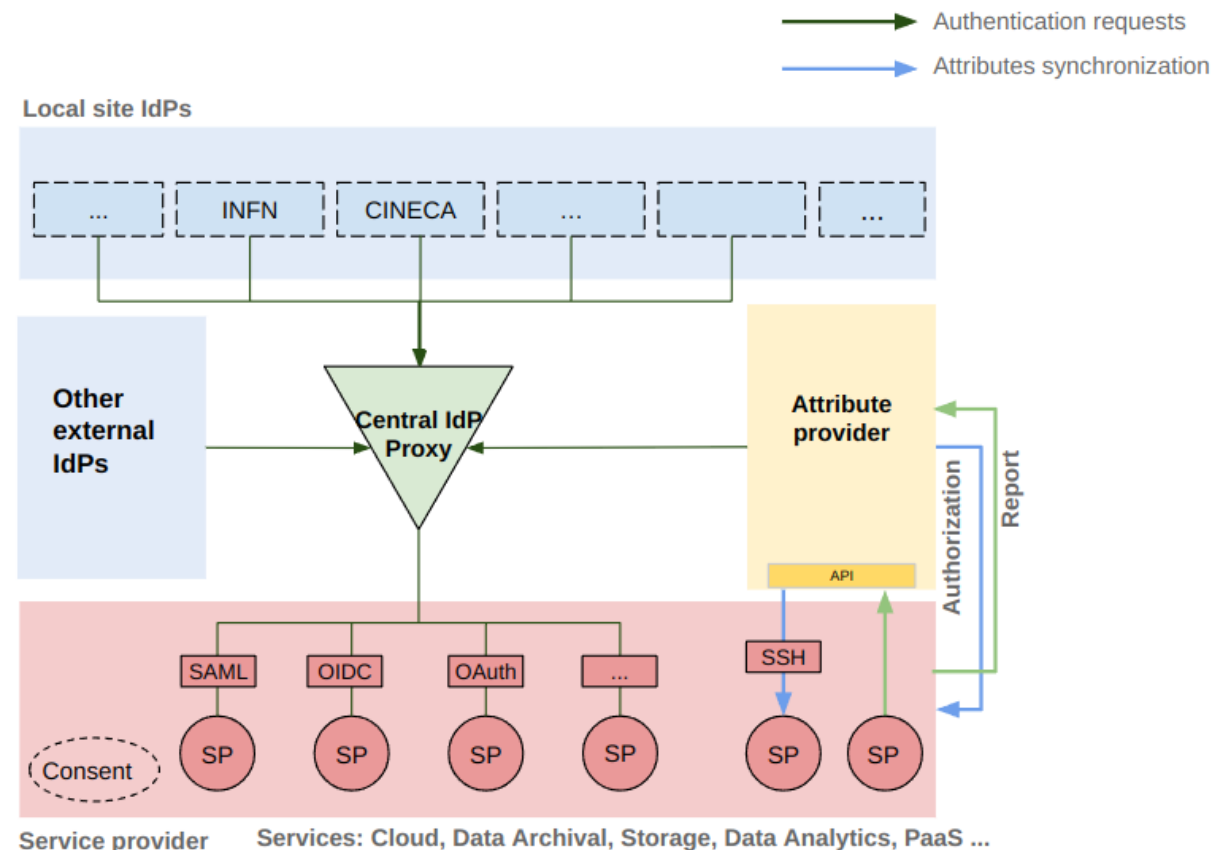
The procedures for joining the federation must be non-intrusive

Standard must be used whenever possible, and developed when missing

The federation will serve users of several fields and organizations

The procedures for user's onboarding must be as simple as possible

E.g.: use of Identity Federations

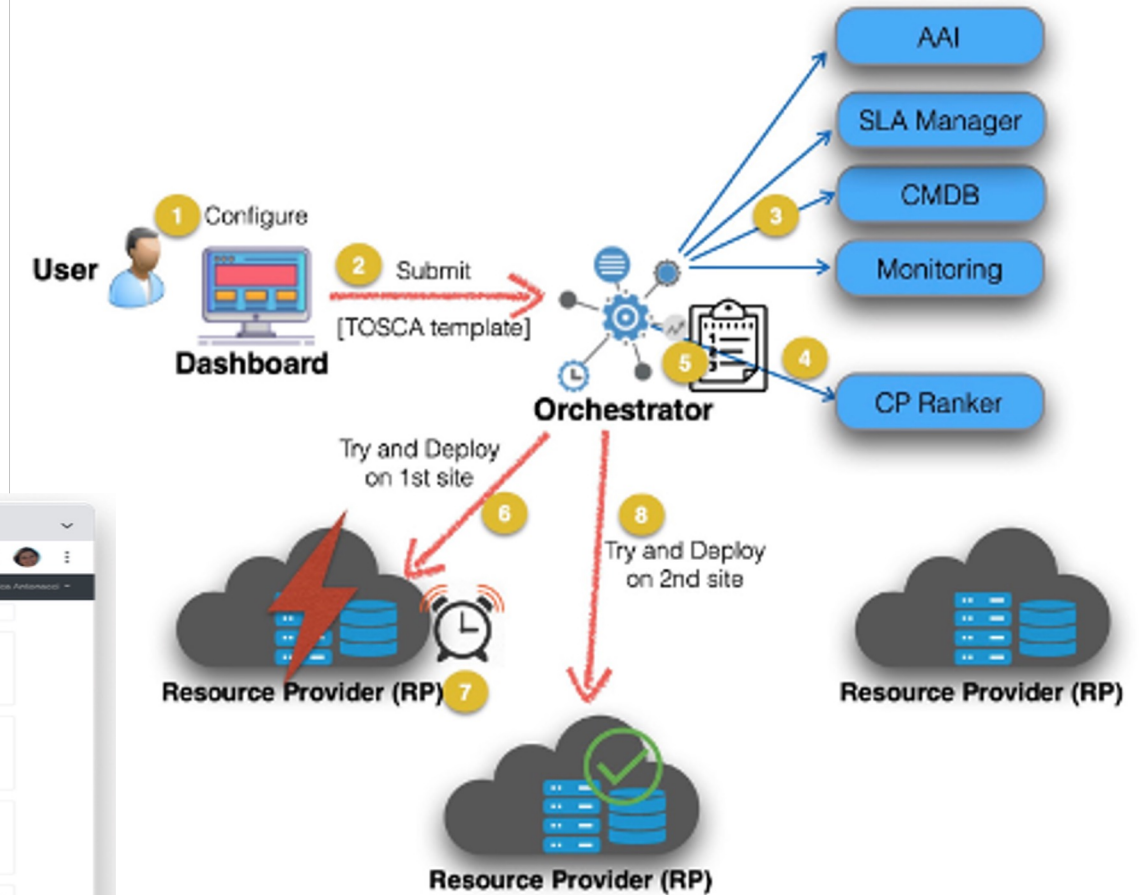
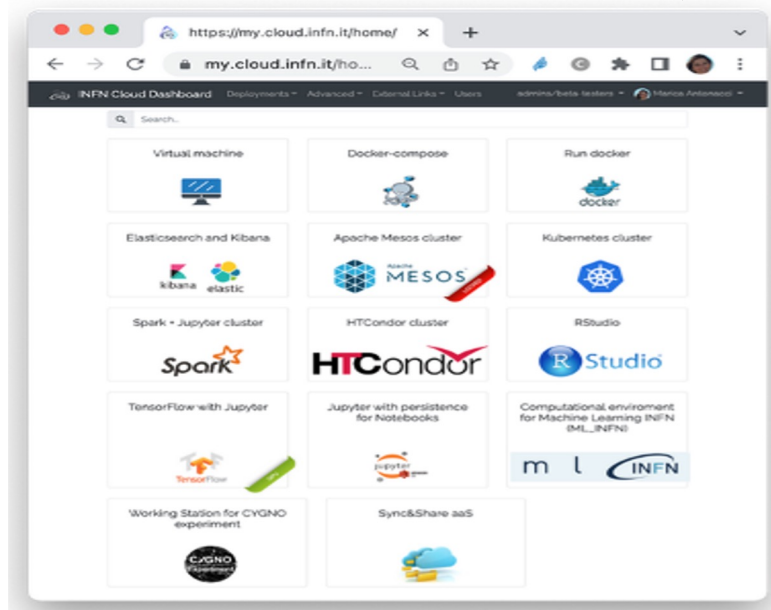


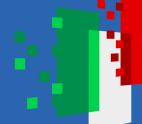
Ease of use

The federation will serve users with different computing competences

Complexity of the underlying infrastructure hidden to the end user

Support field experts in developing platforms that enable the effective exploitation of the infrastructure through composition of services and resources



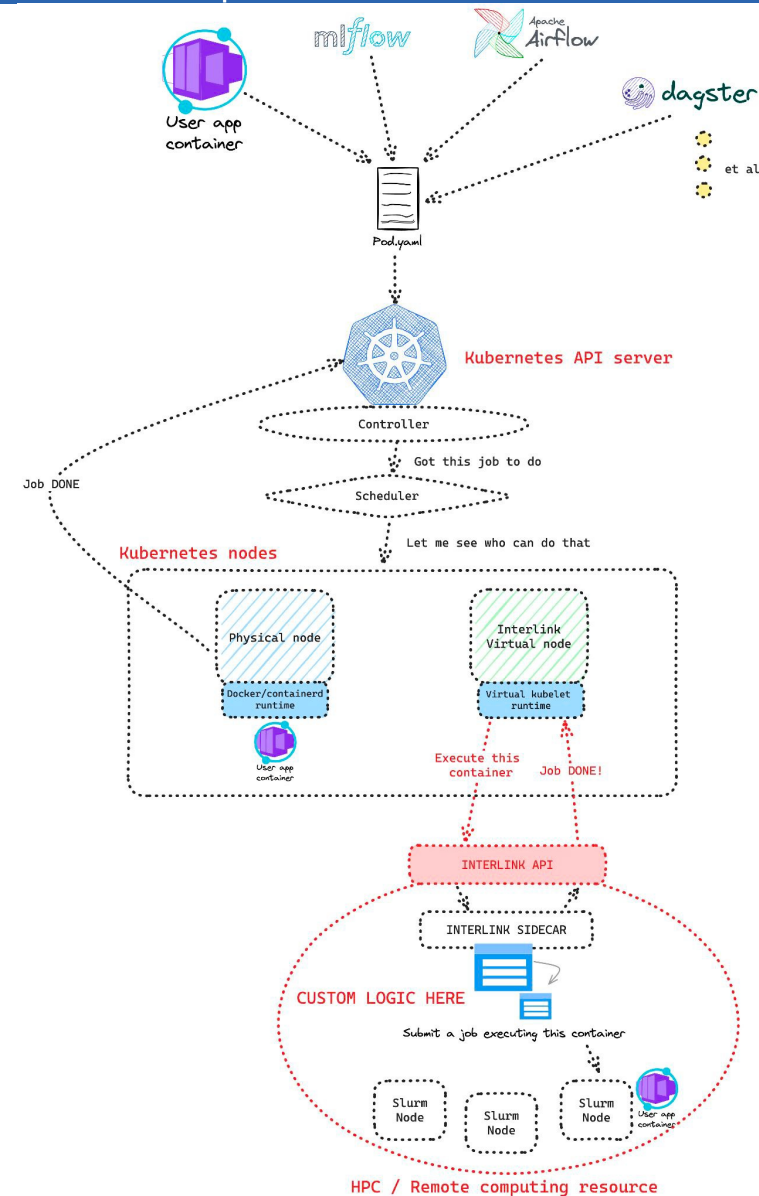
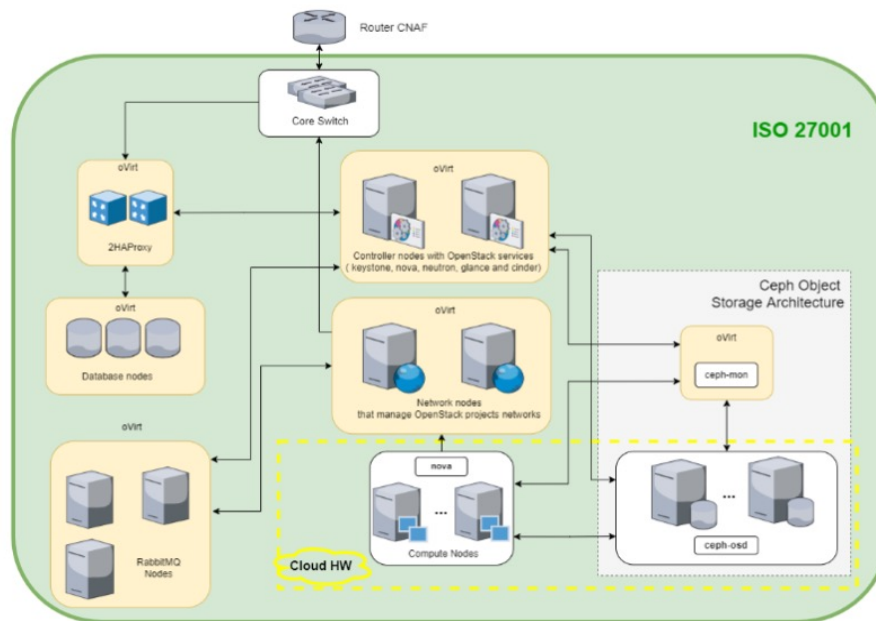


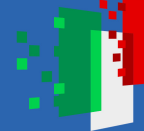
Flexibility

Support multiple access methods to the resources, oriented to:

- a. Transparency and ease of use
- b. Efficiency and effectiveness

Support application-specific requirements
E.g. enhanced privacy





A data lake for research

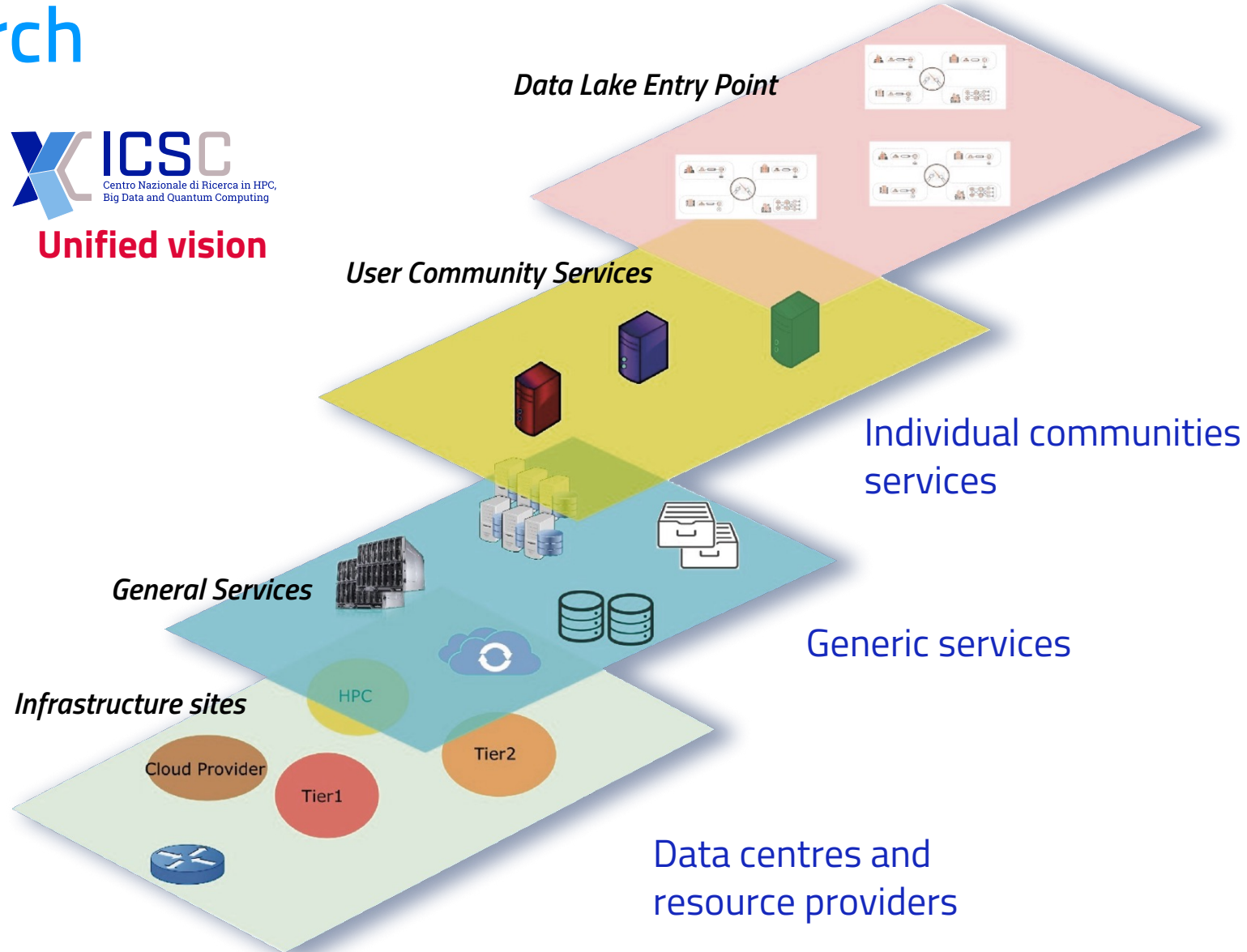
Existing infrastructures aggregation, upgraded and made available to scientific domains

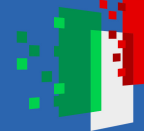
A dynamic model, where infrastructures and domains can also be temporary

A clear separation between the physical and the logical levels

A high-speed network interconnection to hide the actual resource locations

A unified vision (when needed) of an Italian research data-lake



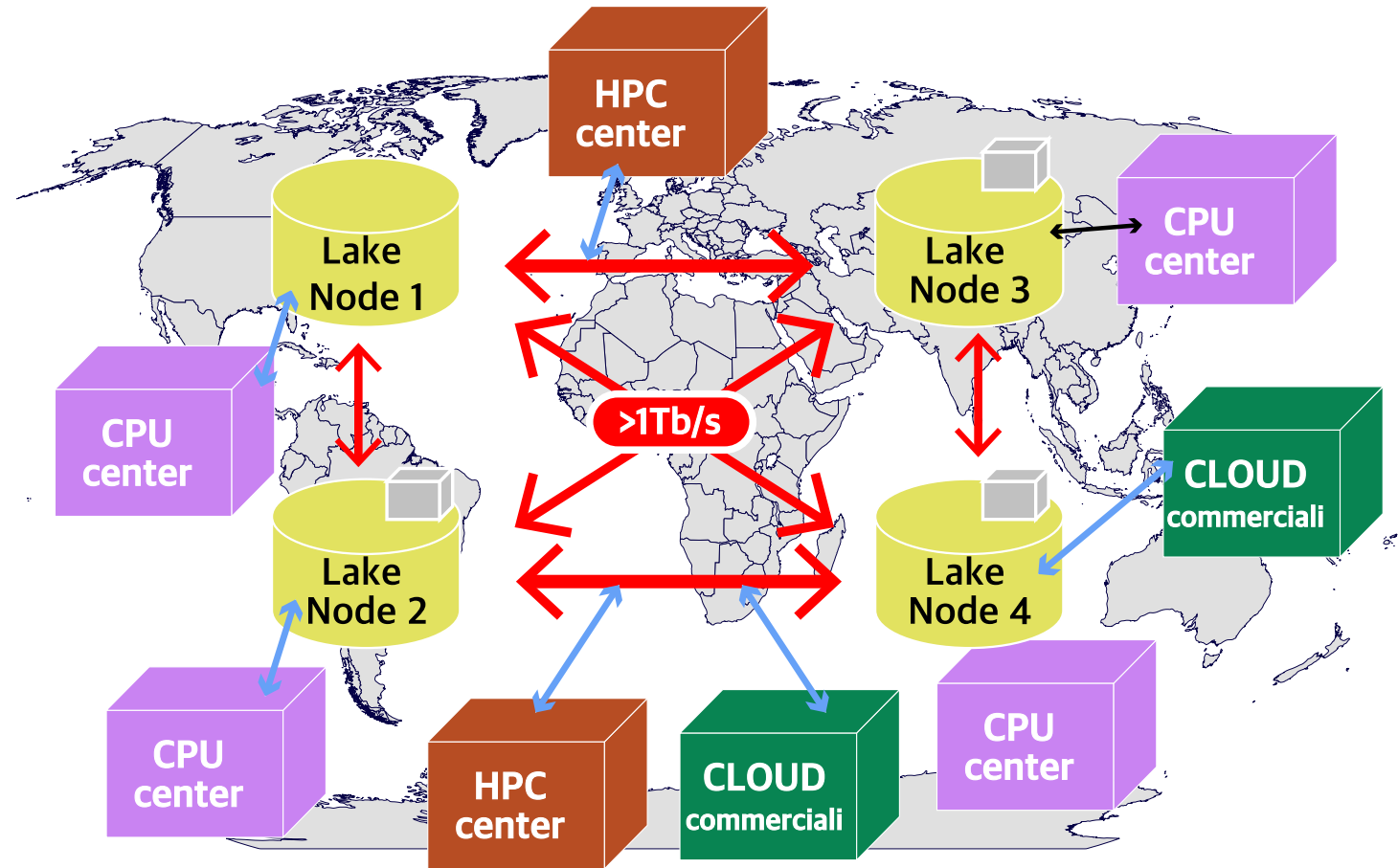


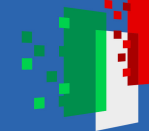
Data-centric model

Decouple storage and CPU

Storage nodes interconnected with high bandwidth network

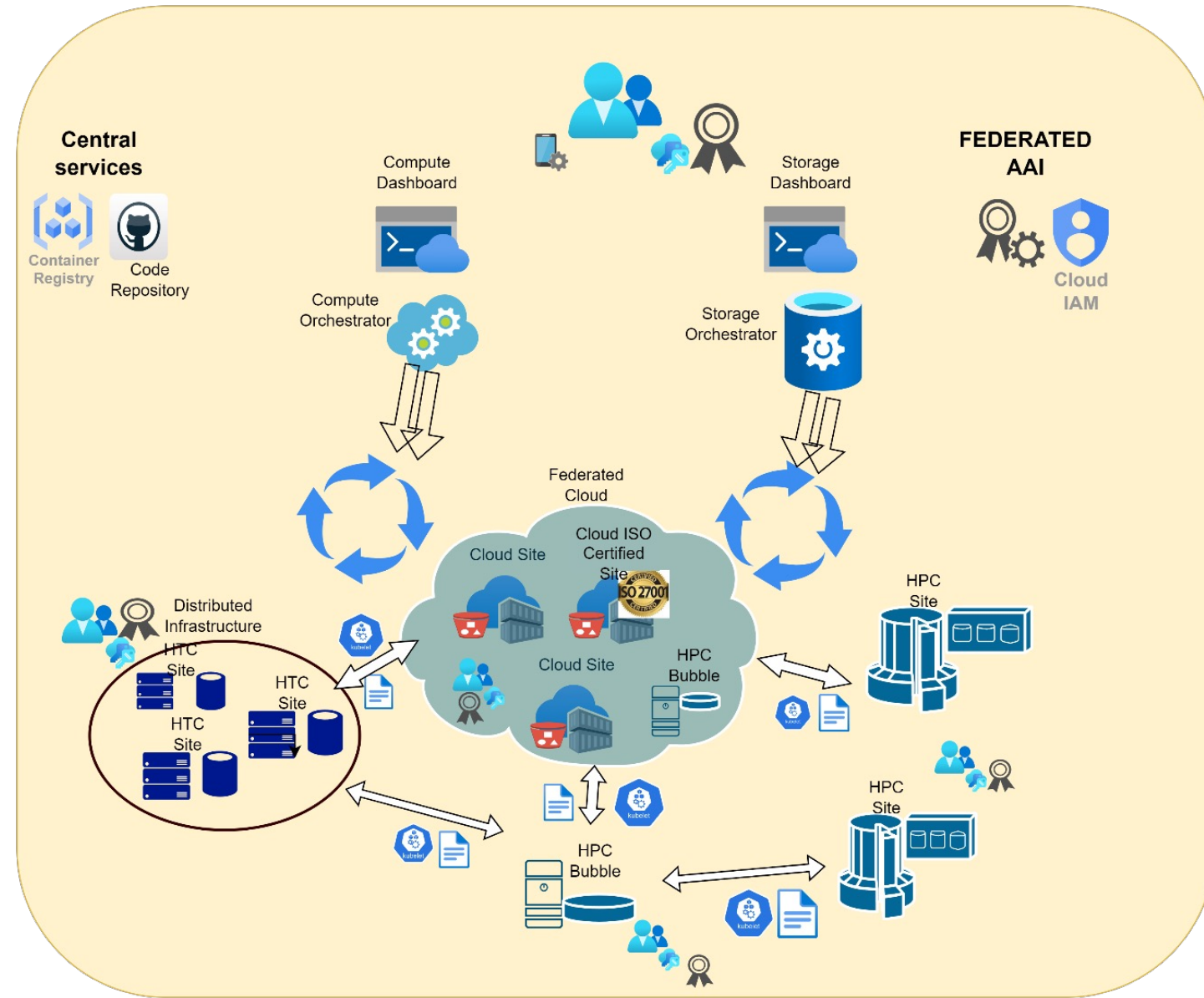
Heterogeneous computing nodes can access data wherever they are





Etherogeneity

Integration of a diverse set of resources, providers, and solutions



A new organization for DataCloud - history

In July 2023, the INFN Steering Committee of the National Computing Coordination (C3SN) created a working group:

Daniele Cesini, Giacinto Donvito, Claudio Grandi, Barbara Martelli, Daniele Spiga

With the mandate to write a document with:

The description of the DataCloud targets

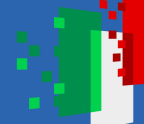
The description of the internal organization, deliverables and milestones

The links to other actors of INFN computing

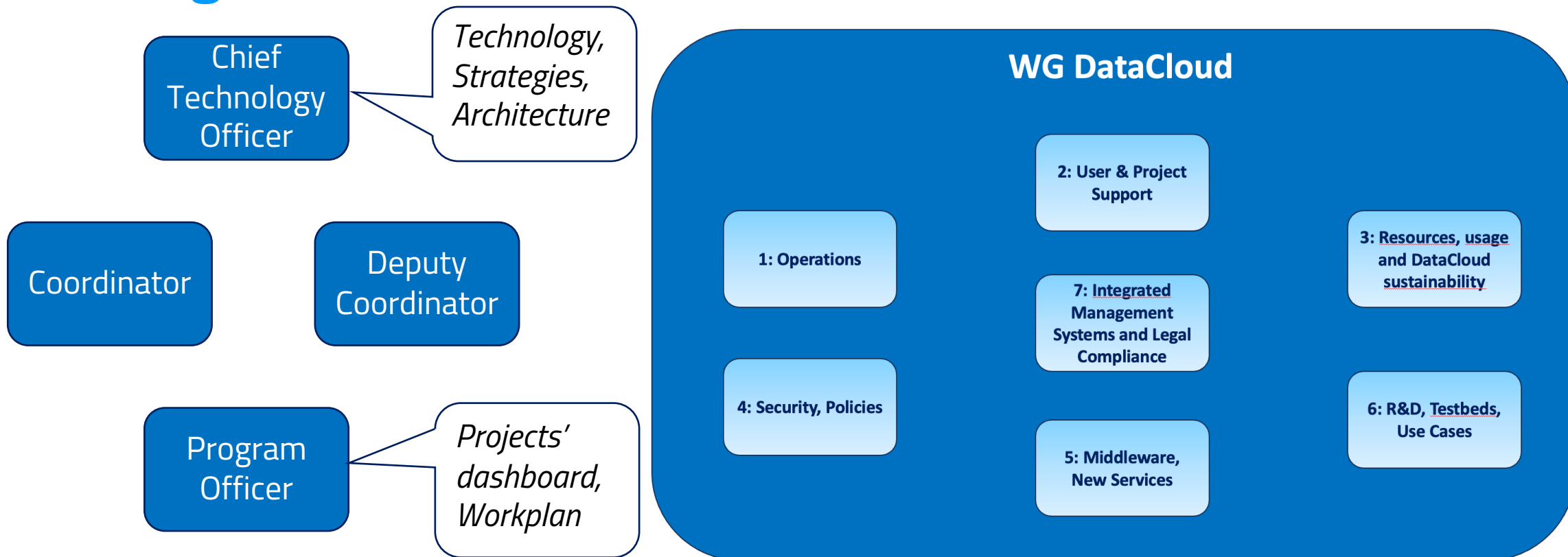
The management structure

The document has been presented to C3SN in October 2010, modified according to requirements and approved by INFN management in December

The new organization is active since mid-January 2024



A new organization for DataCloud - structure



WG organization did not change. See: D.Salomoni, *WG Infrastruttura: mandato, struttura generale e operativa* DataCloud Kick-off, October 17th -18th, 2022

<https://agenda.infn.it/event/32364/contributions/178083/attachments/98155/135704/WG%20DataCloud%20221017%20finale.pdf>



Management of external projects

Exploit synergies with DataCloud

Maintain compatibility

Build common work programs

Project	Deliverable/Milestone (full table here)	Original Date	Real Date	Type	Breakdown	Note
ICSC-S0	Activation of a PoC for resource federation	30 June 2024	30 June 2024	Integration/devel	BD1	
ICSC-S0	HAMMON D1.3 Implementation of the first PoC of the Cloud Platform	30 June 2024	TBD	Integration/devel	BD_H_1	
ICSC-S0	HAMMON D1.4 Implementation of the first integrated version of the Cloud Platform	31 October 2024	TBD	Integration/devel	BD_H_2	
ICSC-S0	End of project: infrastructure in production	31 August 2025	31 August 2025	Integration/devel	BD2	
ICSC-S0	HAMMON D1.5 Implementation of the fully featured high-available Cloud Platform	31 August 2025	TBD	Integration/devel	BD_H_3	
TeRABIT	D4.3 - Report on the deployment of PoC applications over multiple HPC Bubbles	31 December 2023	30 June 2024	Integration/devel	Vedi BD1	Come BD1. Le Bubbles si emulano se non ci sono
TeRABIT	D4.4 - Report on the deployment of a PoC for dynamic caches via the PaaS Orchestrator	31 March 2024	30 June 2024	Integration/devel	BD_T_1	
TeRABIT	End of project: infrastructure in production	30 June 2025	30 June 2025	Integration/devel	Vedi BD2	Come BD2
ICSC-S8	Migrare UISS da piattaforma bare metal a piattaforma cloud certificata e certificarlo come SaaS.	28 February 2024	30 August 2024	Integration/devel/data governance	BD_S8_1	
ICSC-S8	Lavorare sull'integrazione dei servizi cloud INFN e CINECA (in sinergia con Spoke0) a livello tecnico (IAM, Orchestratore, data management) e organizzativo (armonizzazione delle policy)	30 October 2024	30 October 2024	Integration/devel/data governance	BD_S8_2	
ICSC-S8	Avviare e concludere una PoC per replicare il lavoro svolto da A.Cavalli con Humanitas e Engineering sulla piattaforma ICSC. Avremo una prima call di avvio con il rettore di Humanitas il 5 marzo dove inizieremo a pianificare le attività. Obiettivo chiave della PoC sarà dimostrare che la piattaforma cloud ICSC è in grado di gestire le noli di dati gestite attualmente da humanitas/engineering con performance equivalenti.	30 August 2025	30 August 2025	Integration/devel/data governance	BD_S8_3	per questo task dovremmo partire dalla PoC della riga precedente e dimostrare di essere in grado di gestire il trasferimento dei dati genomici dalla fonte (sequenziatori del progetto 5000 genomi di Aosta) alla piattaforma ICSC (CINECA+INFN). Il fatto che CINECA e INFN siano entrambi coinvolti e raggiungibili da un unico entry point DataCloud è un requisito. L'idea è di utilizzare RUCIO/ETS.
DARE	Integrare la piattaforma AlmaHealthDB in ICSC DataCloud (la relazione tra DARE e ICSC è esplicitata nel progetto) e certificarla come SaaS.		30 October 2024	Integration/devel/data governance	BD_M_1	SaaS è inteso come service level nello shared responsibility model standard. In pratica significa che noi ci occupiamo e siamo responsabili (o co-responsabili insieme a partner di UniBO e CINECA) del "contenuto del tenant". Più dettagli tecnici li avremo dopo il 23 aprile, data di una riunione tecnica Spoke8/DARE dove definiremo la proposta. Parlo di Spoke8/DARE perché le persone che lavorano su questo sono le stesse in entrambi i contesti
DARE	Integrare i progetti pilota DARE nella piattaforma ICSC e certificarli come SaaS - Alcuni progetti si riducono all'utilizzo di Ansys su cloud, altri necessitano di workflow manager come NextFlow (expertise di Gasparetto maturata con Sant'Orsola) Sara' probabilmente necessario integrare anche RedCap, gestito da un consorzio di cui fa parte anche il CINECA. Dovremo decidere come porci (entrare nel consorzio? integrare il software come "esterni"?)			Integration/devel/data governance	BD_M_2	
S. Orsola	Costruire una piattaforma di genomica computazionale cloud based e certificata come SaaS. La piattaforma integrerà: Il software dell'ecosistema Elixir (Galaxy &C.) IAM + PaaS Orchestrator + Rucio & Metadati	Settembre '25	Settembre '25	Integration/devel/data governance	BD_M_3	
Health Big Data	Altri elementi che dovremo integrare in HBD sono RedCap (vedi DARE), e XNAT (usato anche in THE, su questo stiamo impostando una collaborazione con A. Retico) Aggiungo alcune informazioni ricevute da ACC (piattaforma principale in produzione su HBD) per l'evoluzione			Integration/devel/data governance	vedi BD_S8_2	



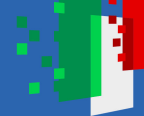
Management of external projects

Exploit synergies with DataCloud

Maintain compatibility

Build common work programs

Project	Deliverable/Milestone (full table link)	Original Date	Real Date	Type	Breakdown	Note
ICSC-S0	Activation of a PoC for resource federation	30 June 2024	30 June 2024	Integration/devel	BD1	
ICSC-S0	HAMMON D1.4 Implementation of the first integrated version of the Cloud Platform	31 October 2024	TBD	Integration/devel	BD_H_2	
ICSC-S0	End of project: Infrastructure in production	31 August 2025	31 August 2025	Integration/devel	BD2	
ICSC-S0	HAMMON D1.5 Implementation of the fully featured high-available Cloud Platform	31 August 2025	TBD	Integration/devel	BD_H_3	
TeRABIT	D4.3 - Report on the deployment of PoC applications over multiple HPC Bubbles	31 December 2023	30 June 2024	Integration/devel	Vedi BD1	Come BD1. Le Bubbles si emulano se non ci sono
TeRABIT	D4.4 - Report on the deployment of a PoC for dynamic caches via the PaaS Orchestrator	31 March 2024	30 June 2024	Integration/devel	BD_T_1	
TeRABIT	End of project: infrastructure in production	30 June 2025		Integration/devel	Vedi BD2	Come BD2
ICSC-S8	Migrare UISS da piattaforma bare metal a piattaforma cloud certificata e certificarla come SaaS.	28 February 2024	30 August 2024	Integration/devel/data governance	BD_S8_1	
ICSC-S8	Lavorare sull'integrazione dei servizi cloud INFN e CINECA (in sinergia con Spoke0) a livello tecnico (IAM, Orchestratore, data management) e organizzativo (armonizzazione delle policy)	30 October 2024	30 October 2024	Integration/devel/data governance	BD_S8_2	
ICSC-S8	Avviare e concludere una PoC per replicare il lavoro svolto da A.Cavalli con Humanitas e Engineering sulla piattaforma ICSC. Avremo una prima call di avvio con il rettore di Humanitas il 5 marzo dove inizieremo a pianificare le attività. Obiettivo chiave della PoC sarà dimostrare che la piattaforma cloud ICSC è in grado di gestire le noli di dati gestite attualmente da humanitas/engineering con performance equivalenti.	30 August 2025	30 August 2025	Integration/devel/data governance	BD_S8_3	per questo task dovremmo partire dalla PoC della riga precedente e dimostrare di essere in grado di gestire il trasferimento dei dati genomici dalla fonte (sequenziatori del progetto 5000 genomi di Aosta) alla piattaforma ICSC (CINECA+INFN). Il fatto che CINECA e INFN siano entrambi coinvolti e raggiungibili da un unico entry point DataCloud è un requisito. L'idea è di utilizzare RUCIO/ETS.
DARE	Integrare la piattaforma AlmaHealthDB in ICSC DataCloud (la relazione tra DARE e ICSC è esplicitata nel progetto) e certificarla come SaaS.		30 October 2024	Integration/devel/data governance	BD_M_1	SaaS è inteso come service level nello shared responsibility model standard. In pratica significa che noi ci occupiamo e siamo responsabili (o co-responsabili insieme a partner di UniBO e CINECA) del "contenuto del tenant". Più dettagli tecnici li avremo dopo il 23 aprile, data di una riunione tecnica Spoke8/DARE dove definiremo la proposta. Parlo di Spoke8/DARE perché le persone che lavorano su questo sono le stesse in entrambi i contesti
DARE	Integrare i progetti pilota DARE nella piattaforma ICSC e certificarli come SaaS - Alcuni progetti si riducono all'utilizzo di Ansys su cloud, altri necessitano di workflow manager come NextFlow (expertise di Gasparetto maturata con Sant'Orsola) Sara' probabilmente necessario integrare anche RedCap, gestito da un consorzio di cui fa parte anche il CINECA. Dovremo decidere come porci (entrare nel consorzio? integrare il software come "esterni"?)			Integration/devel/data governance	BD_M_2	
S. Orsola	Costruire una piattaforma di genomica computazionale cloud based e certificata come SaaS. La piattaforma integrerà: Il software dell'ecosistema Elixir (Galaxy &C.) IAM + PaaS Orchestrator + Rucio & Metadati	Settembre '25	Settembre '25	Integration/devel/data governance	BD_M_3	
Health Big Data	Altri elementi che dovremo integrare in HBD sono RedCap (vedi DARE), e XNAT (usato anche in THE, su questo stiamo impostando una collaborazione con A. Retico) Aggiungo alcune informazioni ricevute da ACC (piattaforma principale in produzione su HBD) per l'evoluzione			Integration/devel/data governance	vedi BD_S8_2	



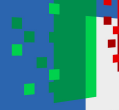
Management of external projects

Exploit synergies with DataCloud

Maintain compatibility

Build common work programs

Project	Deliverable/Milestone	Involved Areas	Involved Services/Components	Fine grained activities	Date	Priority	Owner	Contributors	NOTE
ICSC-S0	Activation of a PoC for resource federation								
Activation of a PoC for resource federation 1- Federazione minimale di un provider Openstack 2- Federazione a livello di applicazione via offloading									
ICSC-S0	HAMMON D1.4 Implementation of the first integrated version of it								
ICSC-S0	End of project: infrastructure in production								
ICSC-S0	HAMMON D1.5 Implementation of the fully featured high-availability								
TeRABIT	D4.3 - Report on the deployment of PoC applications over multiple								
TeRABIT	D4.4 - Report on the deployment of a PoC for dynamic caches via it								
TeRABIT	End of project: infrastructure in production								
ICSC-S8	Migrare UISS da piattaforma bare metal a piattaforma cloud e lavorare sull'integrazione dei servizi cloud INFN e CINECA (il Orchestratore, data management) e organizzativo (armonizz								
ICSC-S8		Compute Federation	PaaS Orchestrator - 1. Openstack federation, - 2. Service level offloading	Deployment of multitenancy dashboard Orchestrator configuration (i.e. federation registry) - Option 1			WP5 WP1	WP2/1 WP5	
ICSC-S8	Avviare e concludere una PoC per replicare il lavoro svolto da piattaforma ICSC. Avremo una prima call di avvio con il rettor pianificare le attività. Obiettivo chiave della PoC sarà di most gestire le noli di dati gestite attualmente da humanitas/engine			Federation validation (i.e. high level service deployment)			WP6		Per la federazione via Orchestratore vanno definiti i servizi high level da "usare" per il PoC
ICSC-S8				High level service deployed by Orchestrator configured to offload toward external Provider (batch and other: vm or k8s)			WP6	WP5	Selezionando un servizio k8s puo' essere abilitato l'offloading nelle seguenti condizioni target. - Nodo (fat node istanziato su un qualunque OpenStack) - batch - k8s
DARE	Integrare la piattaforma AlmaHealthDB in ICSC DataCloud (e progetto) e certificarla come SaaS.	Data Federation	RUCIO FTS	Actual Rucio(s) Deployment FTS Deployment (if needed, might reuse the existing one)			WP6 WP6	WP1 WP1	Puo' essere un servizio dedicato sull'infrastruttura di WP6 Si usa l'istanza INFN.
DARE	Integrare i progetti pilota DARE nella piattaforma ICSC e certificarla come SaaS - Alcuni progetti si riducono all'utilizzo di Ansys su cloud, altri necessitano di workflow manager come NextFlow (expertise di Gasparetto maturata con Sant'Orsola) Sara' probabilmente necessario integrare anche RedCap, gestito da un consorzio di cui fa parte anche il CINECA. Dovremo decidere come porci (entrare nel consorzio? integrare il software come "esterni"?)			Storage endpoint configuration at site (depend also on external providers specific setup)			WP6	WP1	Vanno definiti gli endpoint da federare. Quali e Quanti
S. Orsola	Costruire una piattaforma di genomica computazionale cloud based e certificata come SaaS. La piattaforma integrerà: Il software dell'ecosistema Elixir (Galaxy &C.) IAM + PaaS Orchestrator + Rucio & Metadati		Settembre '25	Settembre '25	Integration/develop/data governance				BD_M_2
Health Big Data	Altri elementi che dovremo integrare in HBD sono RedCap (vedi DARE), e XNAT (usato anche in THE, su questo stiamo impostando una collaborazione con A. Retico) Aggiungo alcune informazioni ricevute da ACC (piattaforma principale in produzione su HBD) per l'evoluzione				Integration/develop/data governance				vedi BD_SR_2



Management of external projects

Exploit synergies with DataCloud

Maintain compatibility

Build common work programs

Activation of a PoC for resource federation

[Access and Identity Management](#)

[Workplan](#)

[Stima dei tempi:](#)

[Compute Resource Federation](#)

[Workplan:](#)

[Federazione di istanze Cloud](#)

[Service Level Offloading](#)

[Stima dei tempi:](#)

[Data Federation](#)

[Workplan](#)

[Stima dei tempi:](#)

Project	Deliverable/Milestone
ICSC-S0	Activation of a PoC for resource federation
ICSC-S0	HAMMON D1.4 Implementation of the first integrated version of IT
ICSC-S0	End of project: infrastructure in production
ICSC-S0	HAMMON D1.5 Implementation of the fully featured high-availability
TeRABIT	D4.3 - Report on the deployment of PoC applications over multiple
TeRABIT	D4.4 - Report on the deployment of a PoC for dynamic caches via IT
TeRABIT	End of project: infrastructure in production
ICSC-S8	Migrare UISS da piattaforma bare metal a piattaforma cloud e lavorare sull'integrazione dei servizi cloud INFN e CINECA (il Orchestratore, data management) e organizzativo (armonizzati)
ICSC-S8	Avviare e concludere una PoC per replicare il lavoro svolto da piattaforma ICSC. Avremo una prima call di avvio con il rettor pianificare le attività. Obiettivo chiave della PoC sarà dimostrarci gestire le nubi di dati gestite attualmente da humanitas/engineering
DARE	Integrare la piattaforma AlmaHealthDB in ICSC DataCloud (e progetto) e certificarla come SaaS.
DARE	Integrare i progetti pilota DARE nella piattaforma ICSC e certificarli come SaaS. - Alcuni progetti si riducono all'utilizzo di Ansys su cloud, altri necessitano di workflow manager come NextFlow (expertise di Gasparetto maturata con Sant'Orsola) Sarà probabilmente necessario integrare anche RedCap, gestito da un consorzio di cui fa parte anche il CINECA. Dovremo decidere come porci (entrare nel consorzio? integrare il software come "esterni"?)
S. Orsola	Costruire una piattaforma di genomica computazionale cloud based e certificata come SaaS. La piattaforma integrerà: Il software dell'ecosistema Elixir (Galaxy &C.) IAM + PaaS Orchestrator + Rucio + Metadati
Health Big Data	Altri elementi che dovremo integrare in HBD sono RedCap (vedi DARE), e XNAT (usato anche in THE, su questo stiamo impostando una collaborazione con A. Retico) Aggiungo alcune informazioni ricevute da ACC (piattaforma principale in produzione su HBD) per l'evoluzione

Involved Areas	Involved Services/Components	Fine grained activities	Date	Priority
AAI		IAM PoC/Test instance Deployment IAM(s) Operations		
Compute Federation	PaaS Orchestrator - 1. Openstack federation, - 2. Service level offloading	Deployment of multitenancy dashboard Orchestrator configuration (i.e. federation registry) - Option 1 Federation validation (i.e. high level service deployment) High level service deployed by Orchestrator configured to offload toward external Provider (batch and other: vm or k8s)		
Data Federation	RUCIO FTS	Actual Rucio(s) Deployment FTS Deployment (if needed, might reuse the existing one) Storage endpoint configuration at site (depend also on external providers specific setup)		
		Integration/develop/data governance	BD_M_2	
	Settembre '25	Integration/develop/data governance	BD_M_3	
	Settembre '25	Integration/develop/data governance	vedi BD_SR_2	

Activation of a PoC for resource federation
1- Federazione minimale di un provider Openstack
2- Federazione a livello di applicazione via offloading

Questo è un documento sintetico e operativo per la realizzazione del primo PoC per la federazione delle risorse ICSC. Tutte le scelte implementative qui presentate sono in linea con l'obiettivo e l'architettura finale ma, ovviamente, tengono conto dei tempi strettissimi del deliverable del 30 Giugno 2024. Pertanto sono state identificate molteplici configurazioni semplificate che saranno evolute in un secondo momento.

Le attività tecniche relative al cronoprogramma descritto in questo documento sono **quasi totalmente parallelizzabili**, ovvero sono attività atomiche e autoconsistenti e non hanno interdipendenze. Lato INFN si possono identificare quindi pool di persone diverse che devono supportare la fase realizzativa.

L'unica vera interdipendenza è la parte di Access e Identity Management, utilizzata sia nella federazione "compute" che in quella "data".

NOTA: In rosso sono indicate le azioni operative che prevediamo siano a carico di CINECA e che quindi vanno concordate con loro (oltre ovviamente a concordare con loro il piano globale). Queste saranno fuori del controllo diretto realizzativo dell'INFN.

Access and Identity Management

In questa prima fase, la proposta è di avere una istanza IAM unica (dedicata al PoC) che federi gli Identity Providers del CINECA e dell'INFN (INFN-AAI) e altri se necessario.

Workplan

- Deployment di uno IAM comprensivo del servizio VOMS-AA

Short term activities

Access to ICSC(++) users

Limited amount of resources

14 allocations by RAC, 5 more in evaluation

Procedures: user registration, verification of user's training

Allocation of resources: cloud when requested, grid when possible, local when needed

User support

PoC of the ICSC/TeRABIT Cloud Federation

Identity federations, authorization

Access to Galileo100 OpenStack

Storage federation via Rucio, access to CINECA S3 storage, offloading to Leonardo

Problems are more related to policies than technical

Mid-term goals

Deployment of procured hardware

Consolidation of middleware

IAM, PaaS Orchestrator, Dashboard

Complete the INFN Cloud federation

Italian Federation in production

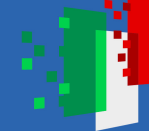
PoC architecture becomes the standard operational environment

Need to build trust mechanisms

Certification of critical components?

Inclusion of other providers (CMCC, ...)

Consolidation of a common User Support



Summary

Extremely challenging task: new approach to resources management, maintenance of systems in production

Central role of INFN in the Italian (and not only) context

Unprecedented quantity of resources and people