

AI_INF_N progress report

27/6/2024

Marco Verlato

Obbiettivi del progetto AI_INFN (2024-2026)

- Raccogliere l'eredità di ML_INFN in termini di comunità, hardware e software...
- ... per costruire un modello di calcolo in grado di soddisfare maggiore domanda e scalare agevolmente in vista di una più ampia disponibilità ed eterogeneità di risorse —> vedi ad es. **Progetto PNRR Terabit**
- Potenziare il supporto ai molti eventi di livello base organizzati dall'INFN e dalle Università, concentrando la propria azione sullo sviluppo di materiale audiovisivo (webinar) ed eventi di aggiornamento di tipo **Advanced Hackaton**
- Formare un nuovo WP dedicato allo studio di futuri acceleratori per le attività di ML, in particolare **FPGA** e **processori quantistici**

Anagrafica

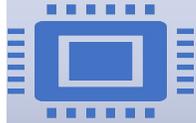
Strutture INFN

BA BO CNAF FE FI GE MIB NA PD PG PI ROMA1 ROMA3 TO

Ricercatori: 82 (12.25 FTE) - Tecnologi: 50 (4.45 FTE) - Tecnici: 3

Andreetto	Paolo		PD	G2	Dipendente	Tecnologo	Attivo	CCR	10%
Fanzago	Federica		PD	G2	Dipendente	Tecnologo	Attivo	CCR	5%
Gianelle	Alessio	5% AI_INF N sinergica a LHCb	PD	G2	Dipendente	Tecnologo	Attivo	CSN1	5%
Traldi	Sergio		PD	G2	Dipendente	Tecnologo	Attivo	CCR	10%
Zangrando	Lisa		PD	G2	Dipendente	Tecnologo	Attivo	CSN4	20%
Sgaravatto	Massimo		PD	G2	Dipendente	Primo Tecnologo	Attivo	CCR	10%
Verlato	Marco		PD	G2	Dipendente	Primo Tecnologo	Attivo	CCR	20%
Jaschke	Daniel		PD	G1	Associato	Scientifica Enti stranieri	Attivo	CSN4	5%
Ballarin	Marco		PD	G1	Associato	Scientifica Dottorandi	Attivo	CSN4	10%
Pagano	Alice		PD	G1	Associato	Scientifica Dottorandi	Attivo	CSN4	10%
Reiniz	Nora		PD	G1	Associato	Scientifica Assegni non INF N	Attivo	CSN4	10%
Sestini	Lorenzo	5% PNRR_ICSCS10 sinergico ad LHCb 5% AI_INF N sinergica ad LHCb 20% MUCOL sinergica a RD_MUCOL C3M - MC_C3M: 10 ore; RADIOLAB_C3M: 30 ore	PD	G1	Dipendente	Ricercatore	Attivo	CSN1	5%
Montangero	Simone		PD	G1	Associato	Incarico di Ricerca scientifica	Attivo	CSN4	5%

Terabit – HPC Bubbles (in arrivo a Luglio)



Nodo CPU

Min 112 core fisici (max 192)
RAM > 8GB/core DDR5
IB NDR 400G
20TBL + dischi di sistema



Nodo GPU

Come CPU + 4x NVIDIA H100 SXM5 con minimo 80GB e memoria HBM2e



Nodo FPGA

Min 32core
RAM > 512GB DDR4 o DDR5
IB NDR 440G
4 x XILINX U55C o 4 x TerasicP0701



Nodo Storage (CEPH Bricks)

Min 48core fisici
RAM >512GB DDR4 o DDR5
Almeno 360 TBL HDD + 12TBL SSD



Accessori

Switch IB, Switch ETH
Cavi IB, Cavi ETH
Transceiver vari
Assistenza 3+2

In particolare, lato Infrastruttura

- Necessita` di superare i limiti del modello basato sull'assegnazione di VM connesse a GPU in modo esclusivo a determinati progetti
- Un nuovo modello basato sui container e` allo studio per:
 - ✓ abilitare il "resizing" delle risorse assegnate ad un progetto
 - ✓ riassegnare GPU altrimenti inutilizzate per lunghi periodi
 - ✓ abilitare un tuning piu` dinamico del n. di CPU core per GPU
 - ✓ abilitare l'uso opportunistico di risorse da parte di code batch
 - ✓ implementare soluzioni di autoscaling (dispiegamento automatico di VM)
 - ✓ ridurre il numero di utenti forzati ad amministrare la propria VM

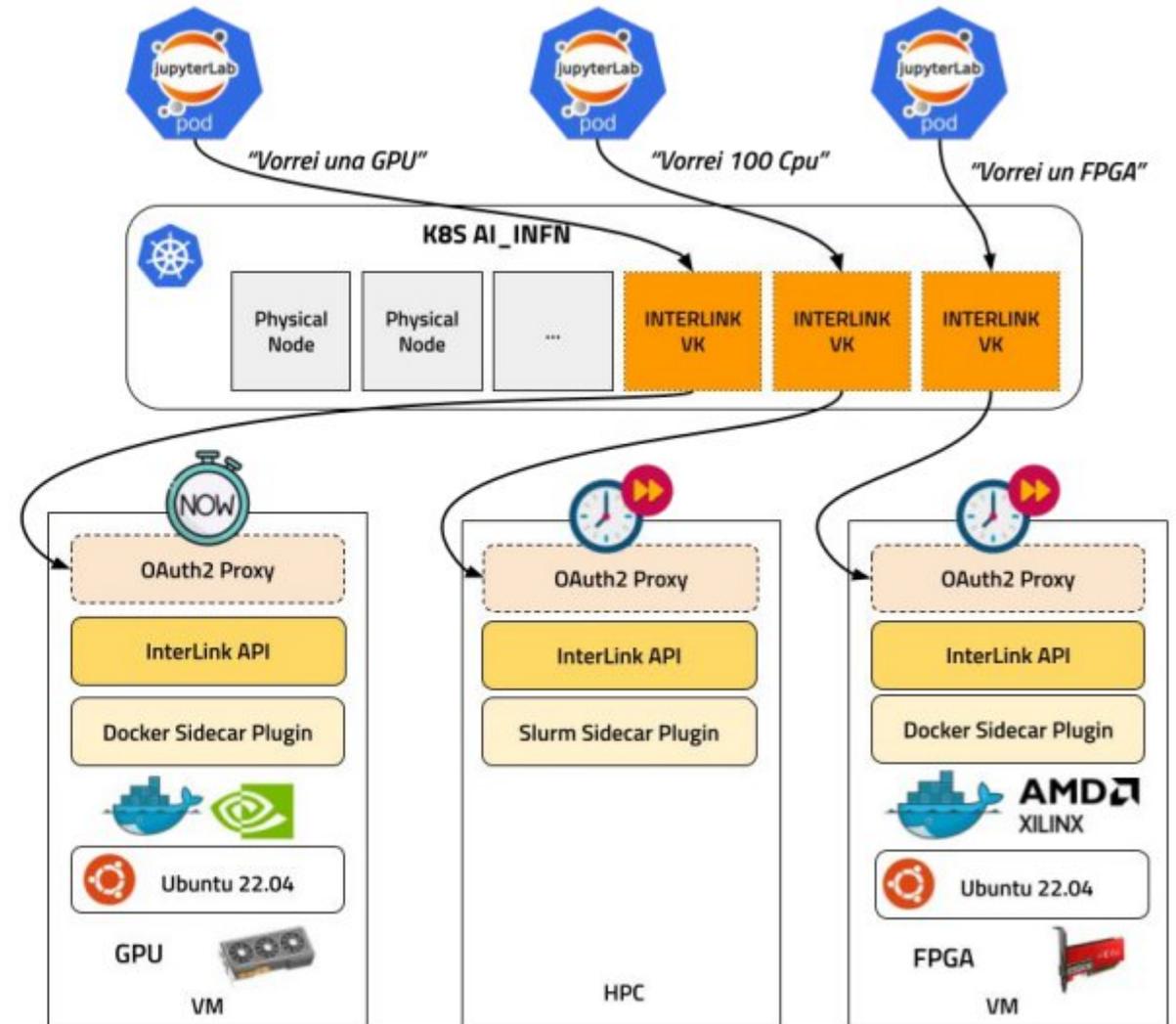
Sinergie con attivita' R&D di INFN-DataCloud

- Rendere trasparente e agile l'uso di risorse disponibili con diversi backend
- Abbiamo adottato il sistema interLink (sviluppato nell'ambito del progetto interTwin) che consente un "**offload trasparente**" di payload containerizzati utilizzando le primitive API Kubernetes, verso un qualsiasi tipo di backend
 - estensione open source del concetto di Virtual-Kubelet con un design che punta ad un'astrazione comune su backend eterogenei e distribuiti
- Sinergia con AI-INFN nello sviluppo di sistemi per testare ed eventualmente validare sia il modello che l'implementazione



Integrazione nella piattaforma AI_INFN

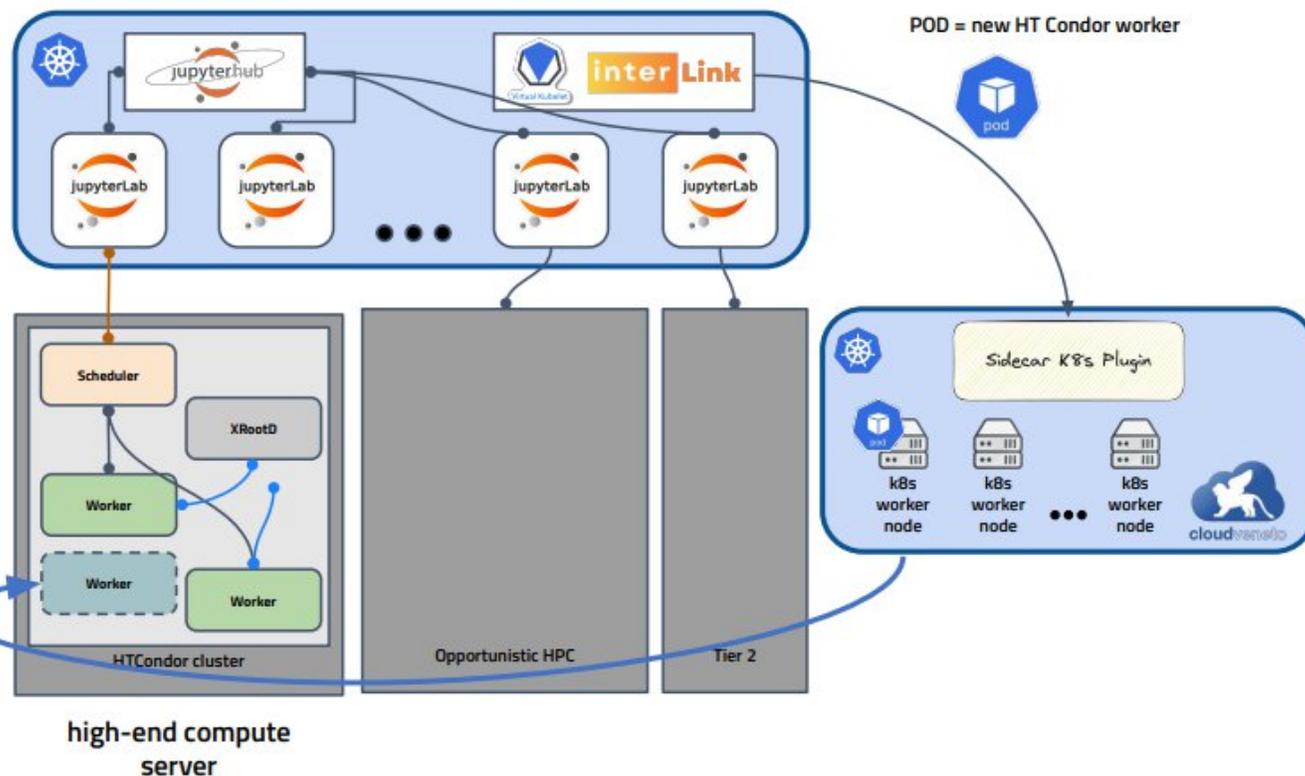
- La piattaforma **AI_INFN** offre uno **use-case cloud-native complesso** per mettere alla prova interLink con:
 - accesso interattivo tramite offloading
 - computing eterogeneo (CPU, GPU, FPGA...)
- Stiamo procedendo in parallelo con lo sviluppo di due plugin (**docker e kueue**) per dimostrare il disaccoppiamento dal backend.
- **Il plugin docker** è già stato validato per il provisioning di GPU e **può supportare analogamente il provisioning di FPGA**.
- In azioni concrete, la piattaforma AI_INFN potrebbe già sfruttare questo sistema di offloading per lo spawn di istanze JupyterLab con alcune limitazioni (NFS)



Integrazione col servizio CaaS di CloudVeneto

E' stato sviluppato il **plugin sidecar kubernetes**. Il POD sottomesso al VK del cluster k8s della Analysis Facility diventa un worker node che si aggiunge al pool centrale HT Condor sfruttando le risorse di un cluster K8s che utilizza risorse di CloudVeneto.

Stato



- ✓ **Workflow completamente funzionante.**
 - nel pool condor vengono aggiunti dei nodi CloudVeneto che possono essere utilizzati per mandare job condor
 - Molteplici VK coesistono nello stesso k8s di AF e fanno offloading su diversi provider:
 - VK dedicato al tier2 di Legnaro, Bari e Pisa (HTCondor sidecar)
 - VK dedicato al tier2 di Roma (ARC sidecar)
- ✓ **Sidecar plugin k8s**
 - Supporta provisioning di CVMFS

Stato dell'implementazione della piattaforma

Feature	Proof of concept	Beta-tested in hub.ai	Available for all users*	Ready for DataCloud
Interactive development (GPU)	2023-05-18	2023-12-13	2024-03-08	
User-defined environments	2023-05-18	2024-04-22	2024-06-03	
Interactive develop. (QC/FPGA)	<i>QC coming soon</i>			
Monitoring & Accounting	2024-03-18	2024-04-22	2024-05-13	
Group-specific resources	2024-04-22	2024-04-22	<i>upon review</i>	
Batch job submission	2023-12-19	2024-04-18	<i>coming soon</i>	
Offloading towards Kueue	2024-05-16	2024-05-27		
Offloading to Docker (GPU)	<i>coming soon</i>			

Documentazione in: <https://ai-infn.baltig-pages.infn.it/wp-1/docs/>

1° AI_INFN User Forum 11-12 Giugno a Bologna

Welcome and introduction <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Lucio Anderlini</i> 14:30 - 14:40
INFN-CNAF: status and perspectives <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Luca Dell'Agnello</i> 14:40 - 15:00
Enhancing Nodule segmentation Utilizing Attention U-Net: Insights from LUNA-16 Dataset <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Arman Zafarani</i> 15:00 - 15:25
Use of a UNet network for the identification of cavities inside mines <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Mr Andrea Paccagnella</i> 15:25 - 15:50
Benchmarking image segmentation on AMD-Xilinx FPGAs <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Ms Valentina Sisini</i> 15:50 - 16:15
Coffee break <i>Room BP-1A, Plesso Berti Pichat</i>	16:15 - 16:45
HERD data classification <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Luca Tabarroni</i> 16:45 - 17:10
Virtual Painting recoloring using Vision Transformer on Deep Embedded X-Ray Fluoresce synthetic dataset <i>Alessandro Bombini</i>	
First Stages on Spectral Classification using Synthetic Datasets. <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Fernando Garcia-Avello Bofias</i> 17:35 - 18:00
AI-based approach for provider selection in the INDIGO PaaS Orchestration System of INFN Cloud <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Luca Giommi</i> 09:15 - 09:40
Leveraging RAG Architecture for Effective Email Response Automation: a CNAF Tier-1 User Support use case <i>Alberto Trashaj</i>	
Transformer-based models for scientific text classification <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Giovanni Zurlo</i> 10:05 - 10:30
Coffee break <i>Room BP-1A, Plesso Berti Pichat</i>	10:30 - 11:00
Multi-scale cross attention transformer encoder for \$tau\$ lepton pair invariant mass reconstruction <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Valentina Camagni</i> 11:00 - 11:25
Hyperparameter Optimization for Deep Learning Models Using High Performance Computing <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Muhammad Numan Anwar</i> 11:25 - 11:50
Quantum Machine learning frameworks for charged particle tracking <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Laura Cappelli</i> 11:50 - 12:15
Final remarks and closing <i>Room BP-2B, Plesso Berti Pichat</i>	<i>Elisabetta Ronchieri</i> 12:15 - 12:30

<https://agenda.infn.it/event/40489/overview>

1° Advanced Hackaton a Padova in Novembre?

Day 1

Day 2

Day 3

Day 4

Lectures:
Introduction to
Autoencoders

Lectures:
Solving diff. eq.
with ML

Lectures:
ML in Medicine,
QML in HEP, etc.

Lectures:
Introduction to
Transformers

Lectures:
Introduction to
CNN and U-Net

**Lectures and
closing**

----- Lunch Break -----

Lectures:
Infrastructure and
AI_INFN platform

Hackatons
1 & 2

Hackatons
3 & 4

**Lectures/
Hands-on**

Hackatons
1 & 2

Hackatons
3 & 4

Formato simile a quello di Pisa tenutosi nel
Novembre 2023:

<https://agenda.infn.it/event/37650/overview>