



Finanziato  
dall'Unione europea  
NextGenerationEU



Ministero  
dell'Università  
e della Ricerca



Italiadomani

PIANO NAZIONALE  
DI RIPRESA E RESILIENZA



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



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

# Centro Nazionale di Ricerca in HPC, Big Data e Quantum Computing

*Lucio Anderlini, Giulio Bianchini, Diego Ciangottini, Federica Fanzago, Rosa Petrini, Massimo Sgaravatto, Daniele Spiga, Tommaso Tedeschi, Antonino Troja*

Prime esperienze nell'integrazione di risorse eterogenee con InterLink: stato e sviluppi futuri

Workshop sul Calcolo nell'INFN, Biodola (LI), 27/05/2025

# 1

## Introduzione

- *Integrazione di risorse eterogenee (HPC, HTC, Cloud)*
- *Provisioning trasparente per l'utente*
- *Gestione backend eterogenei con interLink*

# 2

## Esperienze d'uso

- ***Proof of Concept ICSC per integrazione CINECA - Testbed DataCloud WP6***
- ***Piattaforma AI INFN***
- ***Integrazione con Infrastrutture HPC - Vega, Julich***
- ***High Rate Analysis Platform***

# 3

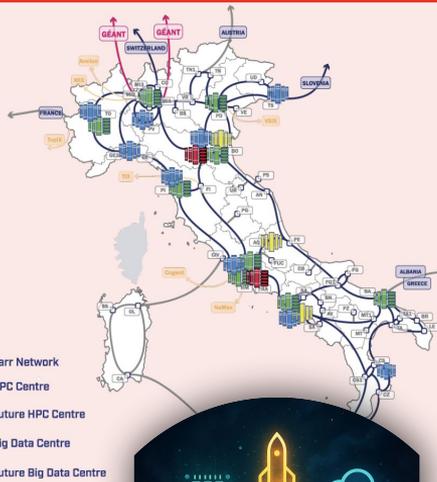
## Sviluppi futuri



The background of the slide is a deep blue color. On the left side, there is a vertical strip of abstract light effects. This strip consists of numerous thin, curved lines that appear to be light trails or fiber optic paths, all converging towards a central point. Interspersed among these lines are small, bright blue dots of varying sizes, some of which are slightly blurred, giving a sense of depth and motion. The overall effect is reminiscent of a digital or data network.

# Introduzione

## 0 SUPERCOMPUTING CLOUD INFRASTRUCTURE



High-level te  
the Spokes working gr



## HPC, HTC and Cloud

Integrare risorse eterogenee come High-Performance Computing (HPC), High-Throughput Computing (HTC) e **Cloud computing** da provider distribuiti è una sfida complessa e attuale.

*Come è possibile garantire un **uso efficiente delle risorse**, una **gestione uniforme dei workload** e la **trasparenza** dell'eterogeneità verso l'utente finale, integrando risorse che hanno modelli di provisioning differenti?*



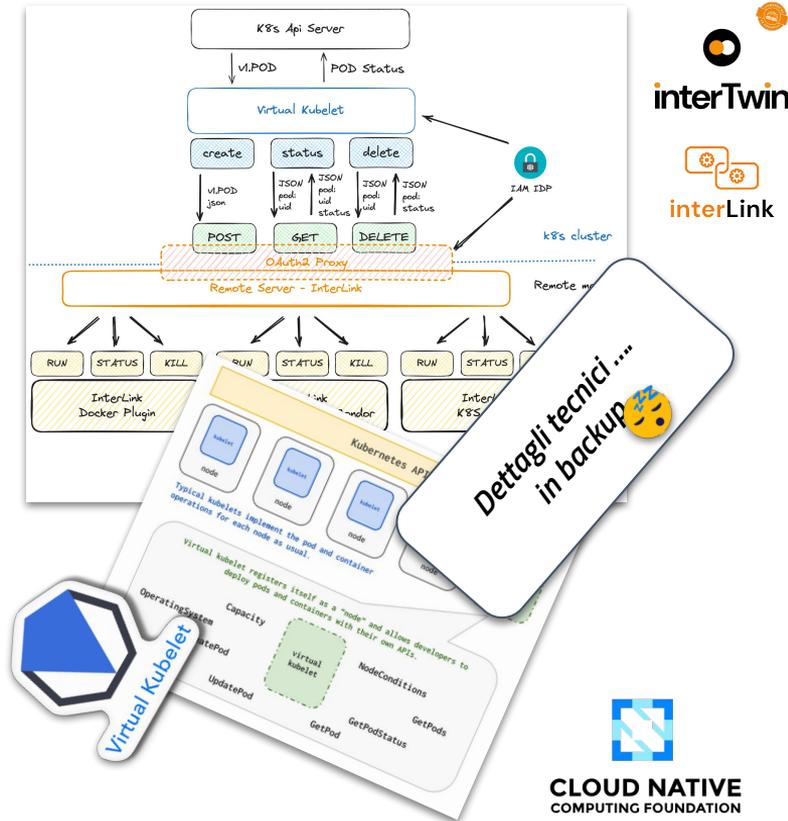
**interLink**

***consente l'esecuzione di qualsiasi container gestito da Kubernetes su backend eterogenei,**  
riducendo al minimo i requisiti per l'utente e per il provider.*

<https://interlink-hq.github.io/interLink/>

## interLink in a nutshell

- ✓ **Cosa**
  - **Gestione unificata di backend e modelli di provisioning diversi**
    - un unico set di API per utilizzare risorse eterogenee (Slurm, HTCondor, Kubernetes...) senza esporre la complessità all'utente finale.
- ✓ **Come**
  - **Esecuzione remota di POD su qualsiasi backend**
    - Estensione di Virtual Kubelet tramite un API layer generico per delegare l'esecuzione di workload su infrastrutture diverse.
- ✓ **Tecnologie abilitanti**
  - **VK Core**: simula un nodo e riceve richieste da Kubernetes.
  - **InterLink**: API server stateless, media le richieste tra VK e Plugin.
  - **Plugin**: esegue i container sull'infrastruttura e restituisce i risultati.
- ✓ **Alcuni highlights:**
  - Capace di gestire risorse eterogenee (GPU/CPU - FPGA)
  - Architettura modulare, personalizzabile tramite plugin.
  - Sui target provider nessuna dipendenza imposta da Kubernetes: estensione trasparente.



<https://www.cncf.io/projects/interlink/>

The background is a deep blue gradient. On the left side, there is a vertical column of light trails and dots that create a sense of depth and movement, resembling a digital or data visualization. The trails are composed of many thin, parallel lines that curve slightly, with small, bright blue dots scattered along them. The overall effect is futuristic and high-tech.

# Esperienze d'uso

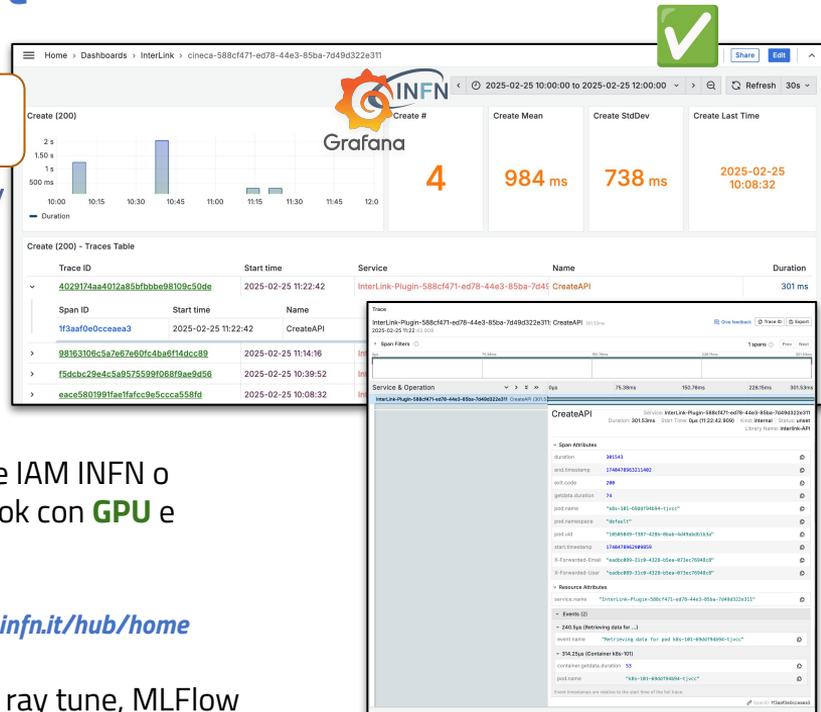
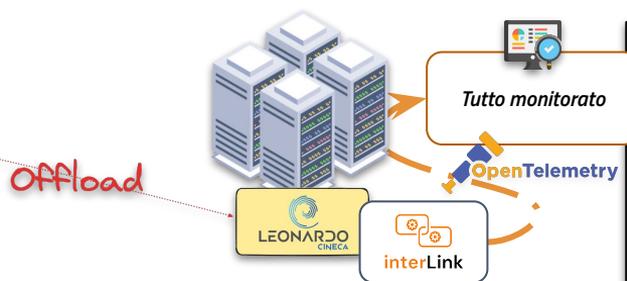
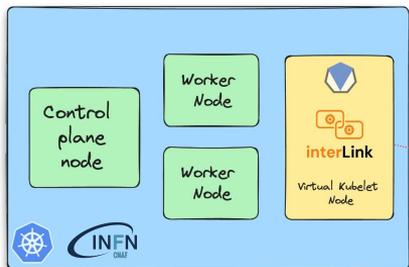


## Disclaimer

Attualmente sono in corso molteplici iniziative dove si sta testando l'integrazione di interLink. In questa presentazione verranno discussi solo alcuni. Altri sono parte di interventi al WorkShop.

- **[Advanced Tracking Analysis in Space Experiments with Graph Neural Networks](#) F. Cuna**
- **[INFN Cloud Kubeflow as a Platform \(KaaP\) and a ChatBot use case](#) M. Gattari**
- **[ICSC HaMMoN: Intelligenza Artificiale e Cloud Computing per il Monitoraggio dei Rischi Climatici](#) A. Casale**
- **Studi su mpi (CNAF)**
- **interTwin [<https://github.com/interTwin-eu/itwinai>]**

## PoC ICSC/TeRABIT per implementazione DataLake Nazionale INFN-CINECA HPC



**Ambiente di test WP6** progettato per sfruttare il meccanismo di offloading tramite InterLink, al fine di accedere alle risorse HPC di LEONARDO/Cineca (**beta users sono BENVENUTI !**).

Da semplici POD ...  
Fino a JHUB (con autenticazione IAM INFN o PoC-ICSC) per spawn di notebook con **GPU** e **Desktop remoto**.

<https://jhub.131.154.99.68.myip.cloud.infn.it/hub/home>

Qualsiasi altra integrazione (i.e. ray tune, MLFlow etc. supportabile come R/D)

[https://monitoring.cloud.infn.it:3000/goto/bKSu\\_cbxNR?orgId=1](https://monitoring.cloud.infn.it:3000/goto/bKSu_cbxNR?orgId=1)

## PIATTAFORMA AI\_INFN

In fase di test (anche) con HPC bubble di Padova

use case selezionato anche per il PoC ICSC/TeraBit

- Tramite la piattaforma [AI\\_INFN](#) (dettagli -> [talk L. Aderlini](#) & [talk R. Petrini](#)) l'utente avvia un JupyterLab dal quale **sottomette un job Modulus via Snakemake**
- La sottomissione del Job avviene tramite il **vk-dispatcher**
- Il Job viene schedulato ad un nodo virtuale **InterLink** che lo inoltrerà verso la risorsa HPC esterna
- Il job inoltrato al sistema HPC (LEONARDO) e' sottomesso allo SLURM workload manager tramite il **plugin SLURM**
- Alcuni metadati vengono scritti su **JuiceFS** ed accessibili all'utente dal suo JupyterLab notebook

The screenshot displays a JupyterLab notebook interface with several key components and annotations:

- Top Panel:** Shows the terminal output of a Snakemake job submission. A green arrow labeled '1' points to the job ID '0'. A green arrow labeled '2' points to the 'vkd' command used for submission.
- Right Panel:** Displays the 'AI\_INFN Kubernetes cluster' configuration, including context, cluster, user, and resource requirements (K9s Rev: v0.40.8, K8s Rev: v1.31.6+rke2r1, CPU: 0%, MEM: 0%). A blue arrow labeled '3' points to the 'interlink' command used for job forwarding.
- Bottom Panel:** Shows the output of the job on the remote HPC node (Leonardo). A green arrow labeled '4' points to the 'JOBID PARTITION NAME USER ST TIME NODES NODELIST(REASON)' table. A green arrow labeled '5' points to the 'INFO:matplotlib.animation.Animation.save using' output, indicating the use of JuiceFS.
- Bottom Right:** A green box labeled 'Leonardo login node' is visible.

## Analisi interattiva su risorse distribuite (grid, bubble, hpc...)

### ✓ Sinergia con use case Spoke2/3

- High rate analysis
- Scientific hub

### ✓ Scale out per "Analysis Facility" di CMS

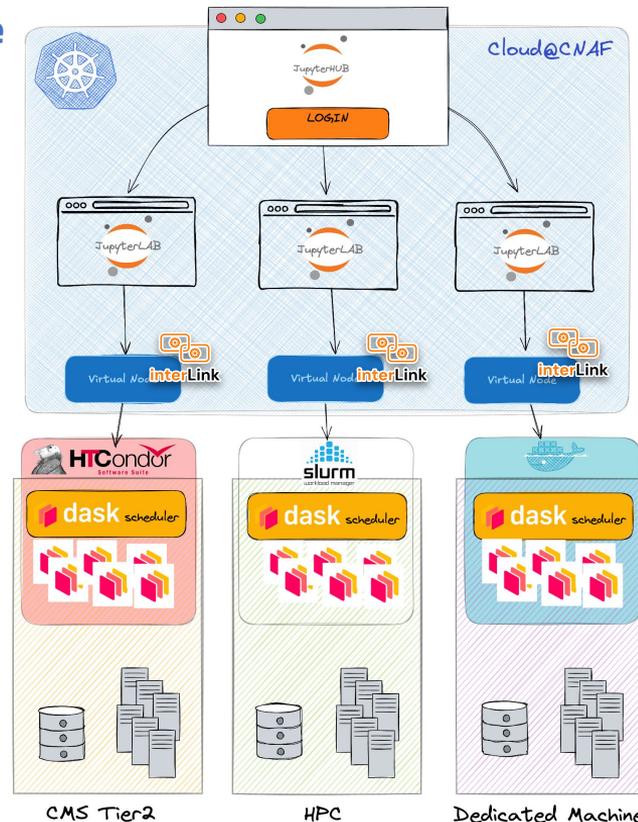


- Integrazione delle risorse HTC Tier2 CMS-IT
- Integrazione con bubble HPC (in modo opportunistico) in collaborazione con infn-pd
  - studi in vista di High-Luminosity Large Hadron Collider

### ✓ Attività in sinergia con il CERN , in valutazione per

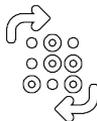


- Analysis facility a CERN (Swan)
- Sistema di provisioning GPU

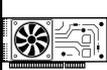


# Highlights from the EU context (VEGA - Slovenian EuroHPC)

## Integrazione con framework ML (sinergia progetto EU interTwin)



Caso d'uso basato su un algoritmo GAN 3D sviluppato al CERN che e' stato integrato con *mlflow* per il monitoraggio e la gestione degli esperimenti ML



Le GPU utilizzate messe a disposizione dai sistemi HPC utilizzati lungo l'intera chain grazie ad InterLink



Le annotazioni sui POD sono un aspetto chiave e uno strumento molto potente: tramite queste e' possibile inviare le richieste specifiche al sistema remoto



In tutto questo, dove si trova il Virtual Kubelet? Il VK fa parte di un cluster k8s fornito dalle risorse cloud dell'INFN

```

1 apiVersion: v1
2 kind: Pod
3 metadata:
4   name: itwinai-cern-use-case-training-run0
5   annotations:
6     slurm-job.vk.io/flags: "--gres=gpu:4"
7     --job-name=itwinai \
8     --output=.local/interlink/jobs/itwinai/job.out \
9     --account=interTwin \
10    --mail-type=ALL \
11    --time=20:00:00"
12 spec:
13   restartPolicy: Never
14   containers:
15     - image: /p/project/interTwin/zoechbauer1/T6.5-AI-and-ML/use-cases/3dgan/containers/cern.sif
16     command: ["/bin/bash"]
17     args: ["-c", "'cd /workspace/T6.5-AI-and-ML/use-cases/3dgan && python train.py -p pipeline.yaml'"]
18     imagePullPolicy: Always
19     name: ai
20     ports:
21       - containerPort: 8080
22   nodeSelector:
23     kubernetes.io/hostname: jul-vk
24   tolerations:
25     - key: virtual-node.interlink/no-schedule
26     operator: Exists

```

JOBID	PARTITION	NAME	USER	ST	TIME	NODES	REASON
68936	batch	Prototyp	remol	GF	1:28	1	hdflc493



```

Predicting: | | 0/? [00:00<?, ?it/s]
Predicting: 0% | | 0/2 [00:00<?, ?it/s]
Predicting DataLoader 0: 0% | | 0/2 [00:00<?, ?it/s]
Predicting DataLoader 0: 50% | | 1/2 [00:02<00:02, 0.34it/s]
Predicting DataLoader 0: 100% | | 2/2 [00:03<00:00, 0.58it/s]
Predicting DataLoader 0: 100% | | 2/2 [00:04<00:00, 0.49it/s]

```



Flusso di lavoro ML che può essere eseguito anche a Cineca



## Dal punto di vista dei provider di risorse...

 <https://github.com/interlink-hq>

- Il modello a **plugin** offre l'opportunità di integrare diversi sistemi di gestione di risorse di calcolo.
- Non tutti i plugin attualmente sono allo stesso livello di maturità'
- esistono plugin personalizzati da utenti finale (punto di forza)



<https://github.com/interlink-hq/interlink-htcondor-plugin>

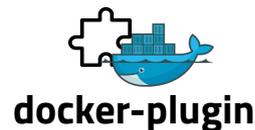


<https://github.com/interlink-hq/interlink-kueue-plugin>



<https://github.com/interlink-hq/interlink-slurm-plugin>

**Soluzioni attuali**  
cosa si può integrare oggi



<https://github.com/interlink-hq/interlink-docker-plugin>



<https://baltig.infn.it/atroja/k8sidecar>  
<https://github.com/interlink-hq/interlink-kubernetes-plugin>  
<https://baltig.infn.it/mgattari/interlink-kubernetes-plugin>



<https://github.com/interlink-hq/interlink-unicore-plugin>

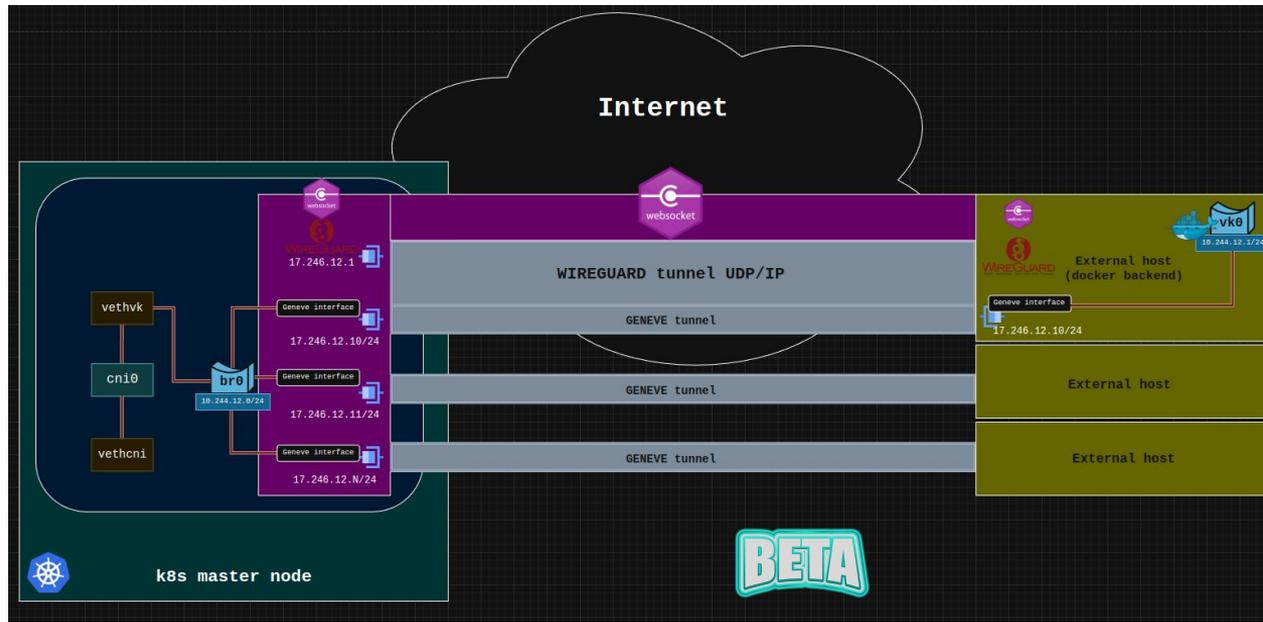
The background is a deep blue gradient. On the left side, there is a vertical column of light trails and dots that create a sense of depth and movement, resembling a digital or data stream. The trails are composed of many thin, parallel lines that converge towards the center, with small, bright blue dots scattered along them. The overall effect is futuristic and technological.

# Sviluppi futuri

## Sviluppi futuri - rete



Abilitare InterLink all'uso di una rete pod-to-pod in ambienti eterogenei e distribuiti (HPC, HTC, cloud) tutto a **user level!**



### Motivazioni principali



#### Comunicazione diretta tra pod

Permetterebbe la comunicazione diretta tra pod come in un unico cluster.



#### Supporto ad applicazioni distribuite

Supporto ad applicazioni con componenti distribuite, come il training distribuito.

## Sviluppi futuri - architetture eterogenee (non solo x86)



AI VK è possibile **specificare le risorse disponibili, comprese risorse specializzate, come GPU e FPGA**, permettendo il loro corretto scheduling e gestione tramite plugin remoti, che associano le risorse richieste ai container e ne instradano il carico sull'infrastruttura.

```
virtualNode:
  #image: ghcr.io/dciangot/interlink/vk:v0.17-debug
  image: biancoj/vk_0.39
  resources:
    CPUs: 8
    memGib: 49
    pods: 100
    accelerators:
      #- resource_type: nvidia.com/gpu
      # model: t4
      # available: 1
      - resource_type: xilinx.com/fpga
        model: u55c
        available: 2
  HTTPProxies:
    HTTP: null
    HTTPS: null
  HTTP:
    Insecure: true
  kubeletHTTP:
    Insecure: true
```

```
apiVersion: v1
kind: Pod
metadata:
  name: podfpga0
  namespace: interlink
spec:
  restartPolicy: Never
  containers:
  - image: my_fpga_image:latest
    imagePullPolicy: Always
    name: first-fpga-test
    command: ["/bin/bash", "-c"]
    args: ["source /opt/xilinx/xrt/setup.sh; xbtutil examine"]
    resources:
      limits:
        xilinx.com/fpga: 1
  dnsPolicy: ClusterFirst
  nodeSelector:
    kubernetes.io/hostname: my-fpga-node
  tolerations:
  - key: virtual-node.interlink/no-schedule
    operator: Exists
    effect: NoSchedule
  - key: accelerator-FPGA
    operator: Equal
    value: U55C
    effect: NoSchedule
```



### Server Options

Select your desired image:

- biancoj/lab-fpga0.2  
JupyterLab with FPGA tools\*
- jupyter/xcipy-notebook:latest  
Jupyter Notebook with Python 3, R, and Julia\*
- jupyter/tensorflow-notebook:3.0.16  
Jupyter Notebook with Python 3 and TensorFlow\*

CPUs

1 2 4 8

RAM

2GB 4GB 8GB 16GB 32GB 64GB

FPGA Model	Total FPGAs	Used FPGAs	Available FPGAs
U55C	2	0	2

Available accelerators:

Select your desired number of FPGAs:

Enable Offloading to:

Start

Docker plugin

Slurm plugin

HTCondor plugin

- Il plugin Docker supporta già il provisioning delle FPGA.
- È in corso lo sviluppo per abilitare anche gli altri plugin alla gestione di queste risorse.
- **L'obiettivo a lungo termine è estendere il supporto ad altri acceleratori hardware e rendere il sistema compatibile anche con architetture diverse da x86 (RISC-V)**

## Summary & conclusions

- Progetto **InterLink** nato circa 2 anni fa nell'ambito dei progetti [ICSC](#) e [interTwin](#) (EU funded)
- Soluzione **lightweight** per l'integrazione di risorse eterogenee via Kubernetes
- In fase di test con molteplici use case:
  - AI\_INFN, analisi ad alto rate, 3D GAN training, ...
  - sfruttamento GPU, scale-out su HPC per training/inferenza etc...



- Uno dei 4 elementi per l'Integrazione tra risorse di



- Attuali risultati sembrano promettenti → fase di consolidamento

- Accettato nella [CNCF Sandbox](#): passo chiave per la sostenibilità a lungo termine



**CLOUD NATIVE**  
COMPUTING FOUNDATION

**Grazie per l'attenzione**

# Backup

Come:

HOW?

Estendendo la soluzione **Virtual Kubelet** realizzando un primo draft di un generico API layer per delegare l'esecuzione di POD su **QUALSIASI** backend remoto.

WHAT?



Virtual Kubelet

Estendere k8s senza imporre dipendenze specifiche di k8s



**VK core**

Un pod che si maschera da nodo e prende in carico le richieste di POD dallo scheduler di K8S

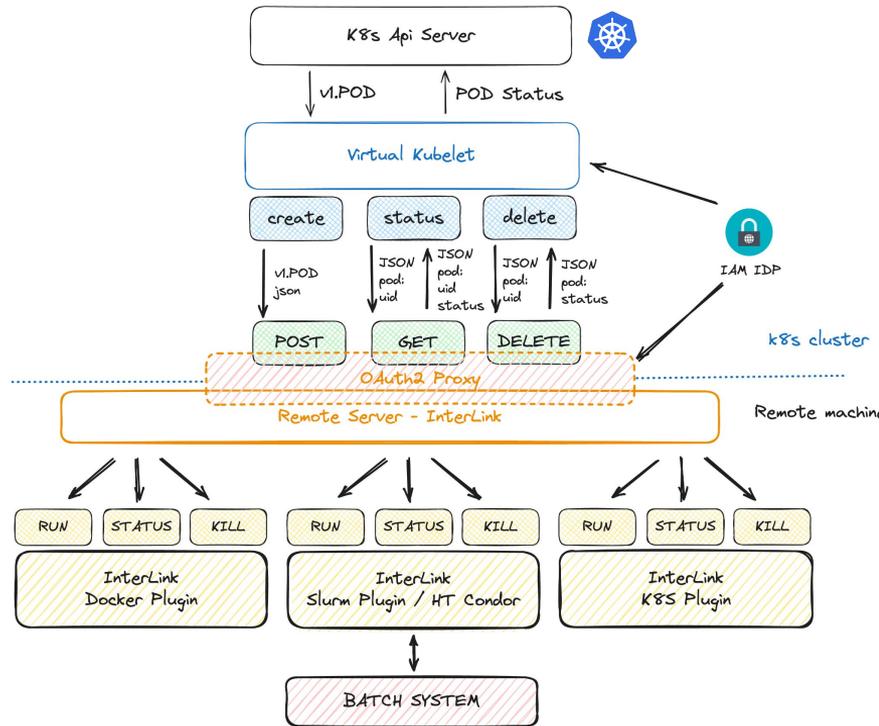
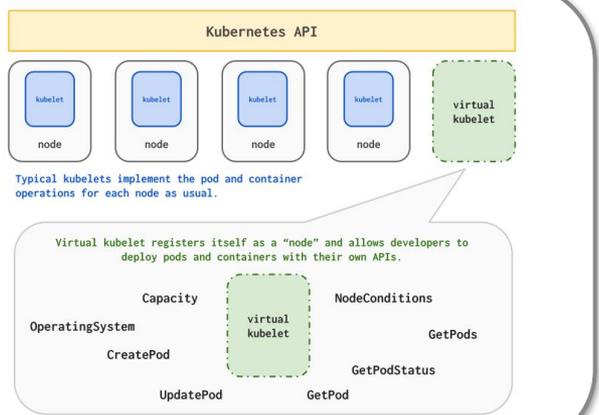
**InterLink Server**

Si interpone tra il VK e il sidecar. Gestisce le richieste provenienti dal VK e le inoltra al sidecar

**Sidecar**

Esegue i container sull'infrastruttura e ritorna il risultato. Comunica con il server InterLink.

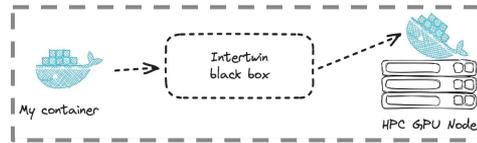
interLink



## Cosa vogliamo abilitare

### Esecuzione di un POD semplice

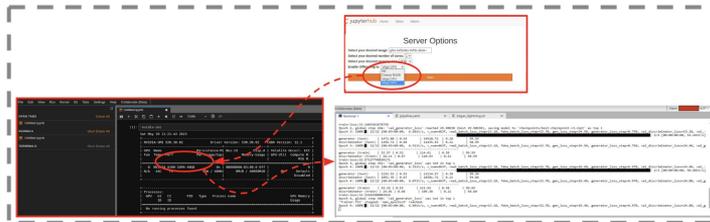
Per creare un semplice container e farlo eseguire da un batch slurm remoto in un HPC



### Sessioni interattive

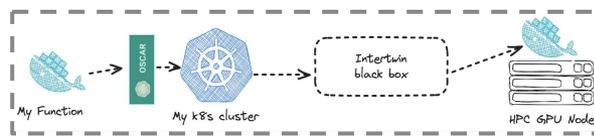
Eeguire istanze JupyterLab su richiesta sia in ambiente HPC

sia in ambienti più tipicamente cloud, tramite Kubernetes.



### Scale out workload

Per lanciare un payload al verificarsi di un trigger esterno



### Gestire backend / modelli di provisioning

- Un insieme unificato di API per integrare risorse eterogenee fornite da provider con tecnologie e architetture differenti



### Spostare i payload di un servizio cloud in base alle specifiche esigenze

- payload compute-intensive, memory-intensive, gpu-intensive, ...



### "nascondere" all'utente l'eterogeneità

- Per l'utente finale l'intero processo è trasparente: il sistema di offloading si occupa di orchestrare l'esecuzione dei workload, decidendo automaticamente su quale backend (Slurm, HTCondor, Kubernetes, ecc.) eseguirli.



### Usare un sistema lightweight e facilmente mantenibile

- Architettura modulare che consente l'integrazione di diversi backend provider attraverso un sistema a plugin

