# Cloud privato per l'AI

**Un caso d'uso GARR**

Alex.Barchiesi @ garr.it

NVIDIA Corporation (NVDA) ☆
NasdaqGS - NasdaqGS Real Time Price. Currency in USD

135.37 +2.61 (+1.97%)
At close: November 1 03:00PM EST

139.32 +3.92 (+2.90%)
After hours: Nov 1, 06:59PM EST

✦ Comparison    ≅ Indicators    📢 Corporate Events    📈 Mountain ⌄    ✎

O 143.00  H 143.14  L 132.11  C 135.40  Vol 983m

yahoo!finance

^IXIC 18,239.92
BTC 6.14

1,905.41%
1,654.51%
112.90%

2,000.00%
1,750.00%
1,500.00%
1,250.00%
1,000.00%
750.00%
500.00%
250.00%
0.00%

Apr  Jul  Oct  2021  Apr  Jul  Oct  2022  Apr  Jul  Oct  2023  Apr  Jul  Oct  2024  Apr  Jul  Oct

1D  5D  1M  3M  6M  YTD  1Y  2Y  5Y  Max    📅 Date Range    Interval: 1 week  ⌄

WORK SHOP GARR 2024    NET MAKERS

1905%

NVIDIA Corporation (NVDA)
NasdaqGS - NasdaqGS Real Time Price. Currency in USD
135.37 +2.61 (+1.97%)
At close: November 1 03:00PM EST
139.32 +3.92 (+2.90%)
After hours: Nov 1, 06:59PM EST

Comparison    Indicators    Corporate Events    Mountain

O 143.00  H 143.14  L 132.11  C 135.40  Vol 983m
^IXIC 18,239.92
BTC 6.14

**Market cap : 360 B$  → 33000 B$**

1,905.41%
1,654.51%
112.90%

1D  5D  1M  3M  6M  YTD  1Y  2Y  5Y  Max    Date Range    Interval: 1 week

yahoo!finance

WORK SHOP GARR 2024    NET MAKERS

NVIDIA Corporation (NVDA) ☆
NasdaqGS - NasdaqGS Real Time Price. Currency in USD

135.37 +2.61 (+1.97%)
At close: November 1 03:00PM EST

139.32 +3.92 (+2.90%)
After hours: Nov 1, 06:59PM EST

**1905%**

⤴ Comparison   ≋ Indicators   📢 Corporate Events      📈 Mountain ⌄   ✏

O 143.00  H 143.14  L 132.11  C 135.40  Vol 983m

^IXIC 18,239.92

BTC 6.14

yahoo!finance

2,000.00%
1,905.41%
1,750.00%
1,654.51%
1,500.00%
1,250.00%
1,000.00%
750.00%
500.00%
250.00%
112.90%
0.00%

## GARR cloud ~ dal 2020:

- 52 (13 x 4) GPU (A100 - A30)
- 333TFLOP

Apr   Jul   Oct   2021   Apr   Jul   Oct   2022   Apr   Jul   Oct   2023   Apr   Jul   Oct   2024   Apr   Jul   Oct

1D  5D  1M  3M  6M  YTD  1Y  2Y  5Y  Max    📅 Date Range    Interval: 1 week  ⌄

WORK SHOP GARR 2024   NET MAKERS

# Outline: cloud privato per l'AI

- Perchè - Why?

- Cosa - What?

- Come - How?

- Sviluppi successivi - Then?

# Why: GPU – ML – not only a buzzword

- **Network optimization:** ML can **analyze** historical network traffic data to predict demand peaks and allow to adequately size network resources, avoiding congestion and ensuring efficient service. It can also **detect** anomalies in the network, such as cyberattacks or hardware malfunctions, enabling a rapid and effective response. Additionally, ML can be used to **find** the most efficient routes for network traffic, reducing latency and increasing network capacity.

- **Data management:** ML can be used to **classify and group** large amounts of data, facilitating search and analysis. It can also reduce data dimensionality, making it easier to visualize and analyze. ML can be used to identify **patterns** in data, which can be used to discover new knowledge and make informed decisions.

- **Research:** ML can be used to analyze large amounts of **scientific data**, discovering new knowledge and accelerating the research process. It can also be used to develop new mathematical models and algorithms, which can be used to solve complex problems in various fields.

- **Scientific collaboration:** ML can be used to develop tools that facilitate the sharing and analysis of scientific data among researchers. It can also be used to identify potential collaborators for research projects, based on their interests and expertise.

**In summary**, Machine Learning can offers enormous potential to improve the efficiency, reliability, and capacity of the network, as well as to develop new services and support scientific research.

WORK SHOP GARR 2024

NET MAKERS

# Challenges and Considerations on private infrastructure

- High initial **investment** in hardware and infrastructure.

- Requires specialized **expertise** in Kubernetes, Kubeflow, and OpenStack.

- Maintenance and management of the **infrastructure**.

- *Potential* **complexity** in integrating different components.

- Careful planning and **resource allocation** are crucial for optimal performance.

# What: Kubeflow the AI Orchestration Engine

- **simplifies** the deployment and management of machine learning (ML) workflows on Kubernetes.

- It provides a user-friendly **interface**

- provides a comprehensive **platform** for building, training, and deploying AI models.

- supports various ML **frameworks** like TensorFlow, PyTorch, and scikit-learn.

- offers **tools** for experiment tracking, model versioning, and pipeline management.

- allow researchers to focus on rapid experimentation with **shared** notebooks and data without worrying about the

  underlying infrastructure.

# Kubeflow vs Jupyter hub

- Jupyter hub is a tool for **individual** data scientists to explore and experiment with data

- Kubeflow is a platform that enables **teams** to build, deploy, and manage large-scale ML pipelines.

# How: Architecture Overview

Layered architecture ensures scalability, resilience, efficient resource utilization.

**OpenStack** provides the underlying IaaS.

**Terraform** manages OpenStack resources (VMs, networks, etc.).

RKE2 forms the **Kubernetes** cluster.

**Ansible** automates Kubeflow installation and configuration.

**Kubeflow** components are deployed on the RKE2 cluster.



**Pipelines**

Kubeflow Pipelines (KFP) is a platform for building then deploying portable and scalable machine learning workflows using Kubernetes.

**Notebooks**

Kubeflow Notebooks lets you run web-based development environments on your Kubernetes cluster by running them inside Pods.

**Dashboard**

Kubeflow Central Dashboard is our hub which connects the authenticated web interfaces of Kubeflow and other ecosystem components.

# GPU timesharing and/or Multi Instance (?)

- **Multi-Instance GPU (MIG):** enables partitioning a single physical GPU into multiple isolated vGPU instances, providing granular control over GPU resources and allowing multiple users or applications to share a single GPU.



- **GPU Time-Sharing:** multiple applications or workloads can share a single physical/virtual GPU.

  **How it Works:**

  1. **Time Slicing:** The GPU's processing time is divided into smaller time slices.
  2. **Task Switching:** The GPU switches between different tasks, allocating a portion of its resources to each task.
  3. **Resource Sharing:** Multiple workloads can share the GPU's memory and compute resources.

# Kubeflow Notebooks

Interactive development environment for exploratory analysis, prototyping…

Some key features include:

- Native support for JupyterLab, RStudio, and Visual Studio Code (code-server).

- Users can create notebook **containers directly in the cluster**, rather than locally on their workstations.

- Admins can provide standard **notebook images** for their organization with required packages pre-installed.

- **Access control** is managed by Kubeflow's RBAC, enabling easier notebook sharing across the organization.

# Kubeflow Pipelines

Pipelines are for automating and managing entire ML workflows, from data ingestion to model deployment.

A pipeline is a definition of a more complex workflow that composes one or more components together. At runtime, each component execution corresponds to a single container execution, which may create ML artifacts. Pipelines may also feature control flow.

With KFP you can author components and pipelines, compile pipelines to an intermediate representation YAML, and submit the pipeline to run on a KFP-conformant backend.
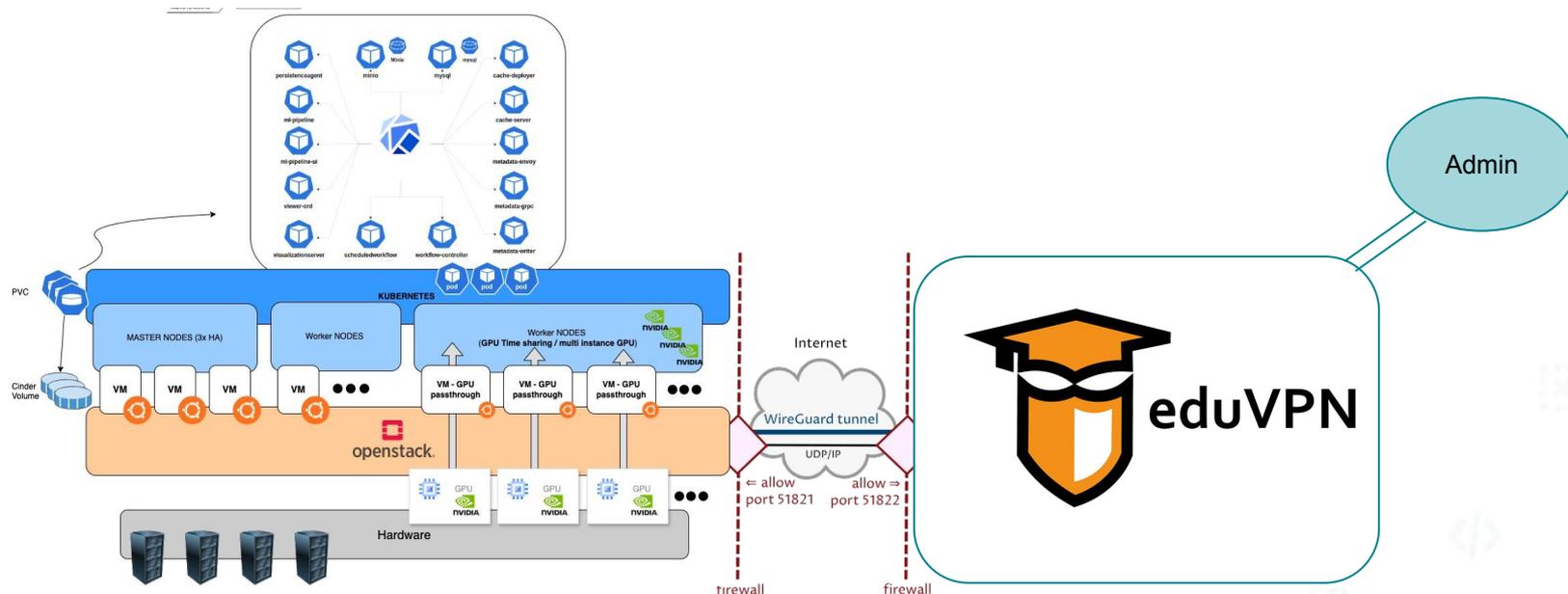
# Details

# Security Considerations

By design: secure communication between components through TLS encryption.

Additional security: tunnel WG - eduVPN to access VM private network

TECH NOTE: OpenStack side you must take care of ALL the security and **disable** Neutron VM security management (port security)

# Then: Future development

- Scalability

- Administration

- Governance

  - Policy to maximize availability of resources (reduce inactivity)

  - Scheduling priority policy

# Thanks

Special thanks and good luck to Alfredo.Funicello ~~@garr.it~~

# Kubeflow Interfaces

LAN Sede Utente
(Direzione GARR)

DNS
che conosce
Jupyter

Host fissi

Istanza
EduVPN

e.g.
EduVPN-test
GARR

Client
remoti

Client
remoti

Client
remoti

Config extra
Wireguard

(7) Ricetta deployment
e configurazione EduVPN
già in uso per eduvpn.garr.it
(da generalizzare)

wg-cloud

(8) configurazione endpoint
Wireguard site-to-site
per terminare "la cloud"
su eduVPN

Tunnel cifrato

Rete GARR

Scenario avanzato ipotetico
(post Conf GARR 2024)

floatingIP
vRouter

floatingIP
Proxy

(5) Provisioning VM
con Terraform

(6) Deployment
e config con
Ansible template-base

cluster
k8S terzo
(on-premise,
su cloud pubblico,
su DC-INFRA,...)

vRouter
Neutron

Proxy Wireguard
config site-to-site

Liqo

Risorse
da
raggiungere
dalla cloud

Liqo usa
Wireguard.
Possibili offload
verso il proxy o
direttamente
con il cluster K8S
con le GPU

) tenant con vDC
u OpenStack
ostruito con
erraform
a template-base

net/subnet
privata

(2) pool di VM (3+)
dedicate a
K8s x Jupyter
(flavor specifico)
con Terraform in

Driver
Container

Driver
Container

Driver
Container

Scenario avanzato
DA NON

NET MAKERS