# Virtual Data Centers: Fueling Data Science with OpenStack

Stefano Cacciaguerra
Stefano Chiappini
INGV

# NE.RE.I.D.E.

INGV implemented the **NEw REsearch Infrastructure Datacenter for EMSO** in the **Western Ionian Sea EMSO site**, at Portopalo di Capopassero (SR), Italy:

➔ an European **research infrastructure** of **EMSO-ERIC** funded by **PON InSEA**

➔ a **ICT infrastructure** for archiving, processing, and sharing **scientific data** from marine observatories and for developing advanced services

➔ it promotes multidisciplinary **scientific/technological research** to understand anthropogenic phenomena in the deep marine environment

➔ **Data science** supports the understanding of complex marine phenomena through **advanced processing of the collected data**

# What do Data Scientists want in NEREIDE ?

**Data Analysis**: Using **scalable computational resources** to analyze large volumes of data

**Model Development**: Exploiting **computing power** to perform complex simulations

**Data Management**: Organize and manage complex datasets, using distributed storage features to ensure **data integrity** and security

**Data Visualization**: Making **complex data clear** with visual tools like maps and charts

**Collaboration**: Working with other scientists to promote **interdisciplinary research**

**Workflow Automation**: Using management tools to **automate** processes and running periodic analyses

**Secure Remote Access**: Accessing to infrastructure from **anywhere**, ensuring real-time research continuity

# How NEREIDE supports Data Science

**Virtual Data Centers** based on Openstack provide a custom and adaptable virtual environment, enabling precise control over applications and data management

➔ **Technologies** like **JupyterHub**, ERDDAP and ElasticSearch Cluster power data analysis and visualization

➔ Openstack **scalable nature** allows data scientists to adjust processing and storage resources for handling large datasets or complex simulations

➔ Enhanced **interdisciplinary** collaboration through **IDEM** and **GARR Cloud Federation**

➔ **Automation tools**, like **MaaS/JuJu**, simplify data science workflows, from data analyses to results sharing

# Openstack → Users, Projects and Tenants

Openstack is a free open standard cloud computing platform deployed as **Infrastructure-as-a-Service** where **cloud resources are made available to users**

**Users** can manage cloud resources through a **web-based dashboard**, **command-line tools**, or **RESTful web services**

**Project** is the base unit of **ownership** in OpenStack (all resources must be owned by a specific project). In OpenStack Identity, a project must be owned by a specific domain

**Tenant** is a group of users in charge of a **logical grouping** of cloud resources

# Tenant & Virtual Data Center

In our solution, **Tenants** are in charge of cloud resources where users could install, configure and manage virtual machines behind a **Gateway** owning a **public IP address** realizing a own **Virtual Data Center** (**VDC**).

➔ **cloud admins** make **cloud infrastructure** and **real Data Center** works

➔ **tenant users** create **Service** in their **VMs** on **VDC**

*"Admins are owners of a mall. Admins entrust a tenant with the management of a shop. Admins are in charge of managing the whole infrastructure, managers of their own shop"*

# Networks & Tenants

There are different types of networks:

**provider** → mapped directly to an **existing physical network** You can use flat provider networks to connect instances directly to the **external network**. (managed by *cloud admins*)

**project** → multiple **private networks** are fully isolated by default and are not shared with other projects. (managed by *tenant users*)

**shared** → networks shared among **all tenants**! (managed by *cloud admins*)

**shared (RBAC)** → Role-Based Access Control (RBAC) networks shared among **specific tenants**! (managed by cloud admins)

# Tenants on Openstack Networking

# VDCs

# Optimization of Tenants

# Tenant's Gateway

In order to share services on Internet, Tenant must use a Gateway owning a **specific Public IP address** and it must work as:

➜  a **Border Router**

➜  a **Firewall**

➜  an **OpenVPN** server

The Gateway is a VM inside the **T0 tenant** (cloud admins) with **two interfaces**:

➜  one on the **Provider** Network

➜  one on the **RBAC Shared** Network

A simple linux VM or something customed like **Endian Firewall**, IPFire or OpenWrt

# Endian Firewall as Gateway

It is an open-source **router**, **firewall** and **gateway** security Linux distribution

Credentials of Endian Firewall:

➔ **root** of SO → Cloud Admins

➔ **admin** of Web Dashboard → Cloud Admins

➔ **admin user** of Web Dashboard → Tenant users

# Endian Firewall Services

## Port Forwarding / DNAT

to make **accessible services** from Internet

| # | Incoming IP | Service | Policy | Translate to | Remark | Actions |
|---|---|---|---|---|---|---|
| 1 | 10.240.3.122 (Uplink main) | TCP/8888 | 🔍➡ | 172.16.0.20 : 8888 | | ⬇ ☑ ⊕ ✏ 🗑 |
| | ALLOW with IPS from: | | | <ANY> | | ✏ 🗑 |
| 2 | 10.240.3.122 (Uplink main) | TCP/8000 | 🔍➡ | 172.16.0.20 : 80 | | ⬆ ☑ ⊕ ✏ 🗑 |
| | ALLOW with IPS from: | | | <ANY> | | ✏ 🗑 |

## Masquerading / SNAT

to allow VMs to **access Internet**

| # | Source | Destination | Service | NAT to | Remark | Actions |
|---|---|---|---|---|---|---|
| 1 | <ANY> | Uplink ANY | <ANY> | Auto | standard uplink SNAT | |

## OpenVPN

to allow users to

**access VMs** (ssh,https,etc)

OpenVPN settings

Authentication type
X.509 certificate

Server certificate

Certificate configuration *          10.242.3.122
Use selected certificate             ⓘ View details

Certificate Authority
ca
📄 Download certificate

▸ Advanced options

Save                                 * This Field is required.

OpenVPN server configuration

Bind only to                         Port *
10.240.3.122                         1194

Network options

Device type                          Protocol
TUN                                  UDP

Bridged                              VPN Subnet
☐                                    192.168.16.0/24

▸ Advanced options

Save  or Cancel                      * This Field is required.

# Jupyter Ecosystem

**JupyterHub** (**JH**) is an open-source web application that allows **multiple users** to interact with **Jupyter Notebooks** (**JN**) on a shared server

With JH, users can log in to a **central server** using their **own credentials** and access their **own JNs**, which are hosted on this server

JNs are **interactive documents** containing executable code (like **Python**, **R**, **Julia**), visualization and text editing capabilities, it is a useful tool for **data science**

JNs can be used for **data cleaning** and **transformation**, **numerical simulation**, **statistical modeling**, **machine learning**, …
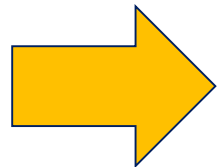
# Jupyter As a Service of VDC

In order to implement  JupyterHub **as a service**, it is necessary:

1. create your VM via **Horizon**

2. connect the **openVPN** server (the Tenant's Gateway)

**3.** **ssh** with **key pair** to your VM

4. install the **Littlest JupyterHub** (**TLJH**)

> **TLJH** places 2 systemd units on your **VM**
>
> ➔ **jupyterhub.service** - starts the JupyterHub service
>
> ➔ **traefik.service** - starts proxy HTTPS

# Dashboard of Jupyter

# What is the main result?

From "**Tenant Users**" side → **Data Scientist**

⇒ lets tenant users operate complex infrastructure (**ready-to-use**) without the task of setting up and managing a **Real Data Center**
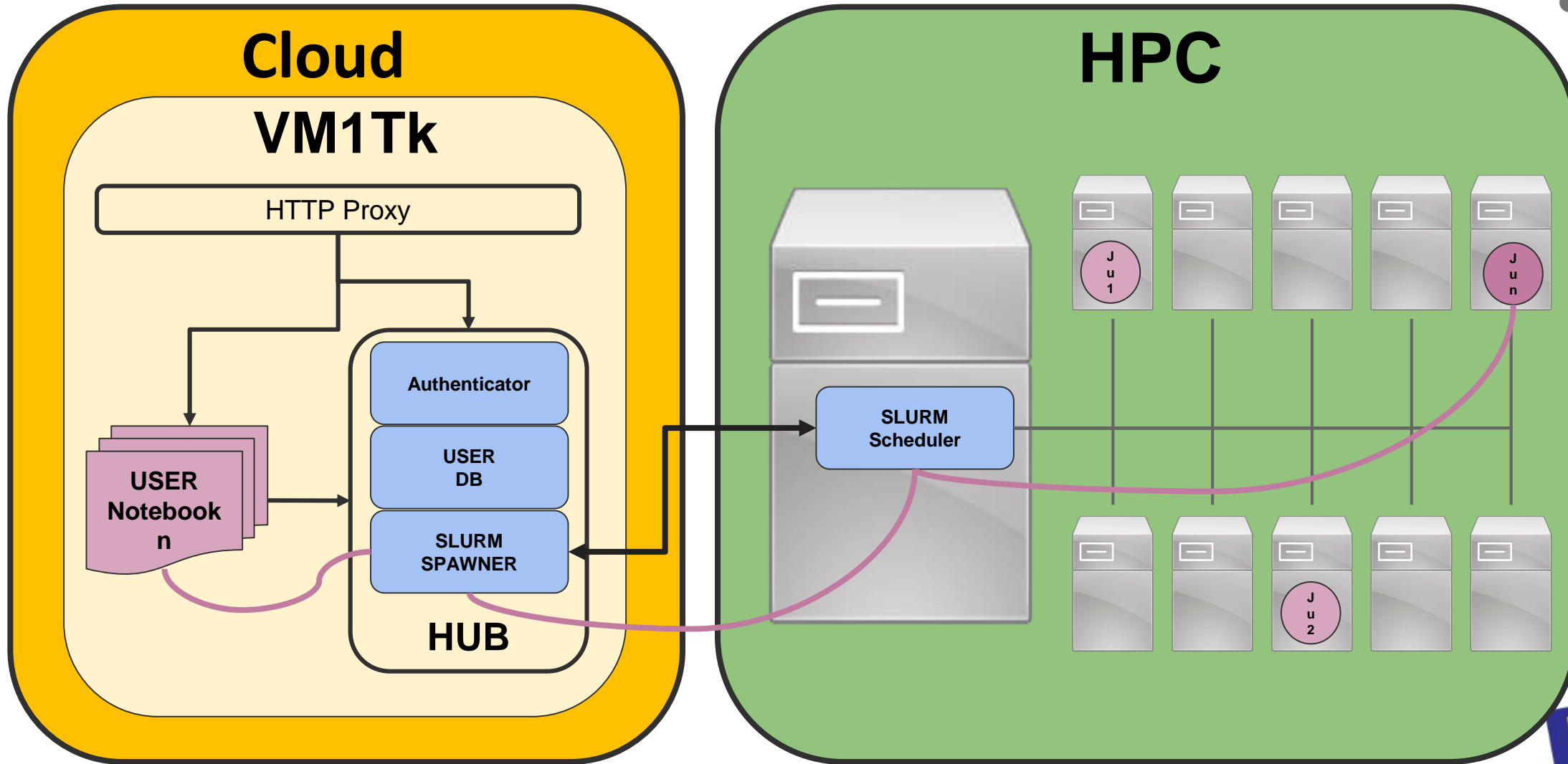
From "**Cloud Admins**" side → **Data Engineer**

⇒ **Centralizing** e **Optimizing** the resources

⇒ **Virtualizing** all the resources

⇒ **Rationalizing** the infrastructure investments

# Jupyter Slurm Spawner from Cloud to HPC

# VDC Migration

Is it possible to Migrate a VDC from a source cloud openstack to a target one?

# Acknowledgement

We would like to express our gratitude for their collaboration to:

⇒ **Alex Barchiesi**

⇒ **Alberto Colla**  } **GARR**
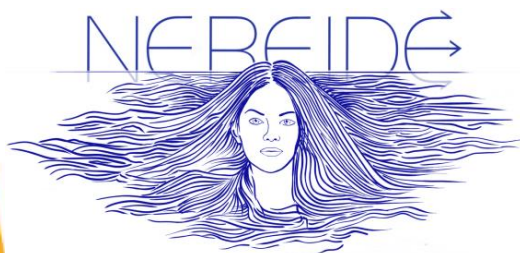
⇒ **Claudio Pisa**

and we would like to mention:

**NereideBO tenant** created in Bologna on **Cloud SUPER** (**POR-FESR** - **Supercomputing Unified Platform - Emilia-Romagna**) is used for part of the development and experiments on VDC