# WHEN PERFORMANCE MATTERS

# HETEROGENOUS COMPUTING: FROM BENCHMARK TO RECONFIGURABLE INFRASTRUCURE

Giordano Mancini, CTO

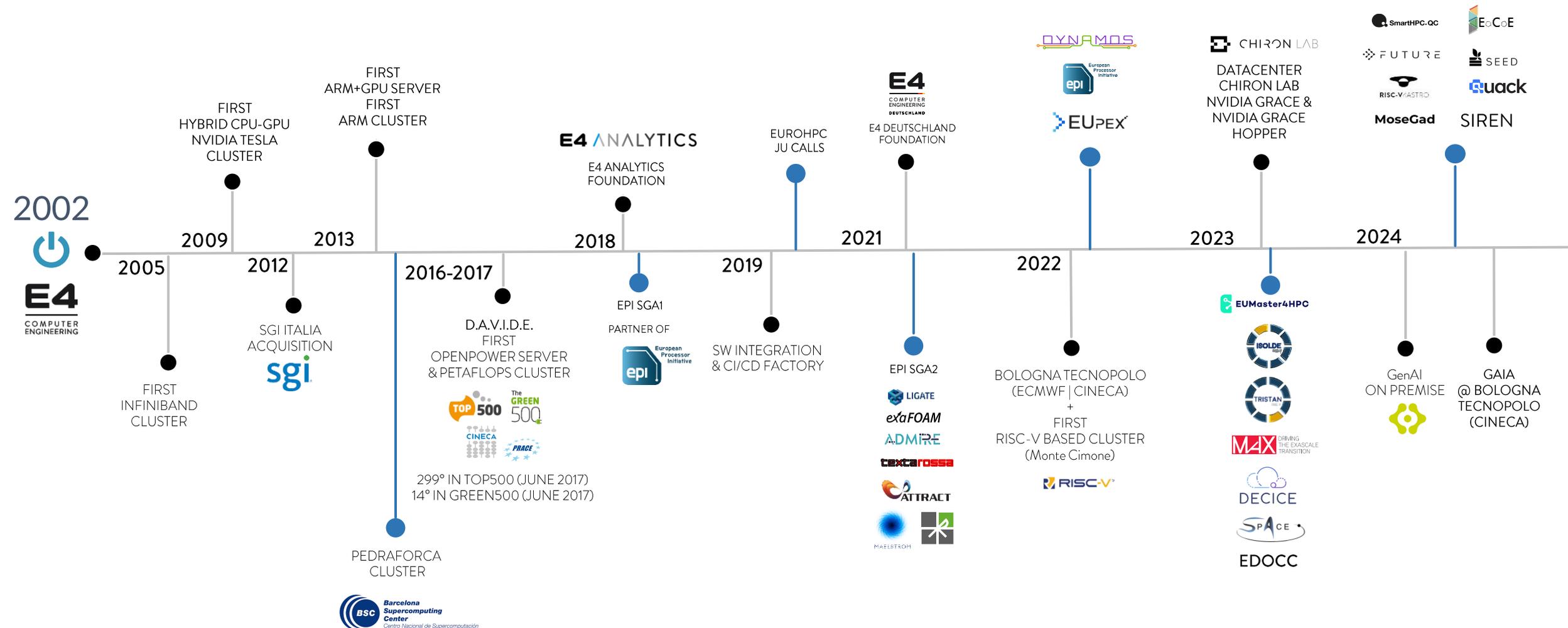Workshop sul calcolo INFN, La Biodola, 26 – 30 maggio 2025

www.e4company.com

# AGENDA

- Who we are

- Limiting factors of HPC in 2025

- Case study A: Benchmarks on heterogeneous architectures

- Case study B: The Čerenkov array pilot project

- Case study C: Genomics cluster @IIT

- Flagship project 2024-2025: GAIA

- Conclusions

# E4 TIMELINE

E4 COMPUTER ENGINEERING

**2002**

**2005**
FIRST INFINIBAND CLUSTER

**2009**
FIRST HYBRID CPU-GPU NVIDIA TESLA CLUSTER

**2012**
SGI ITALIA ACQUISITION

**2013**
FIRST ARM+GPU SERVER
FIRST ARM CLUSTER

**2016-2017**
D.A.V.I.D.E.
FIRST OPENPOWER SERVER & PETAFLOPS CLUSTER

PEDRAFORCA CLUSTER

299° IN TOP500 (JUNE 2017)
14° IN GREEN500 (JUNE 2017)

**2018**
E4 ANALYTICS
E4 ANALYTICS FOUNDATION

EPI SGA1
PARTNER OF

**2019**
EUROHPC JU CALLS

SW INTEGRATION & CI/CD FACTORY

**2021**
E4 COMPUTER ENGINEERING DEUTSCHLAND
E4 DEUTSCHLAND FOUNDATION

EPI SGA2
LIGATE
exaFOAM
ADMIRE
textarossa
ATTRACT
MAELSTROM

**2022**
EUPEX

BOLOGNA TECNOPOLO (ECMWF | CINECA)
+
FIRST RISC-V BASED CLUSTER (Monte Cimone)
RISC-V

DYNAMOS
European Processor Initiative
epi

**2023**
CHIRON LAB
DATACENTER CHIRON LAB NVIDIA GRACE & NVIDIA GRACE HOPPER

EUMaster4HPC
ISOLDE
TRISTAN
MAX DRIVING THE EXASCALE TRANSITION
DECICE
SPACE
EDOCC

**2024**
SmartHPC-QC
FUTURE
RISC-V.MAESTRO
MoseGad

EoCoE
SEED
Quack
SIREN

GenAI ON PREMISE

GAIA @ BOLOGNA TECNOPOLO (CINECA)

# E4 IN NUMBERS

**+287%**

Revenues growth from 2021

**>280**

Active customers (last 3 years)

**92**

Employees

**>10.000**
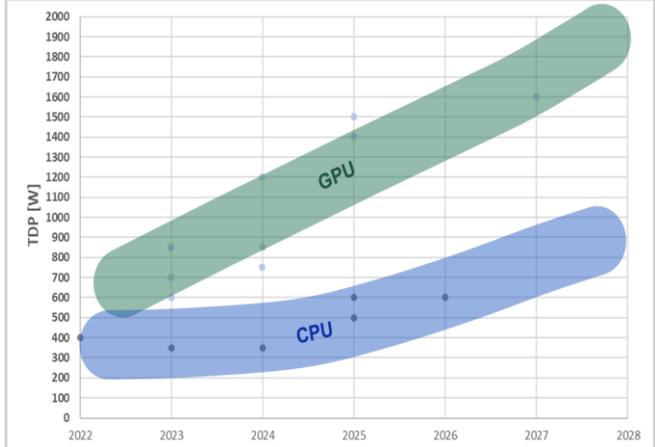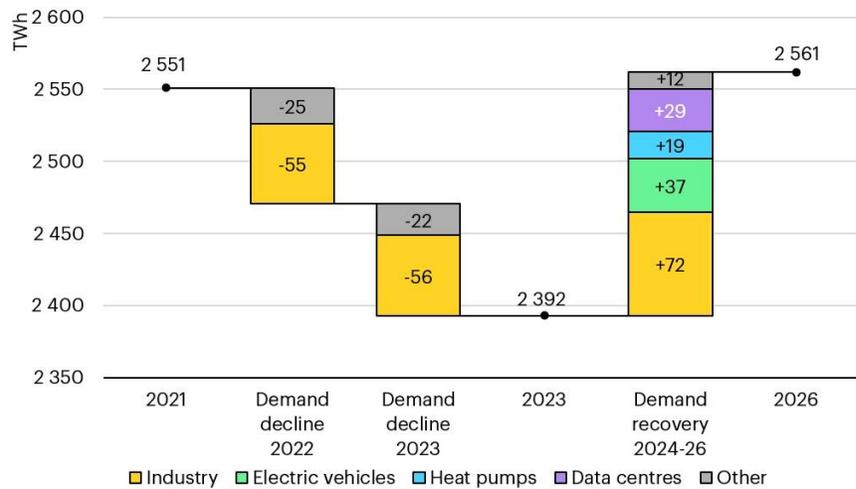
Systems produced in the last year

**>85%**

Employees in Tech/R&D

**26**

E4 contribution in European Projects

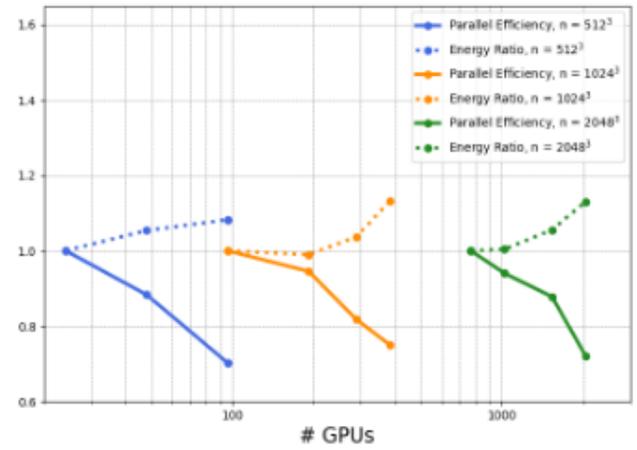# ENERGY CAP

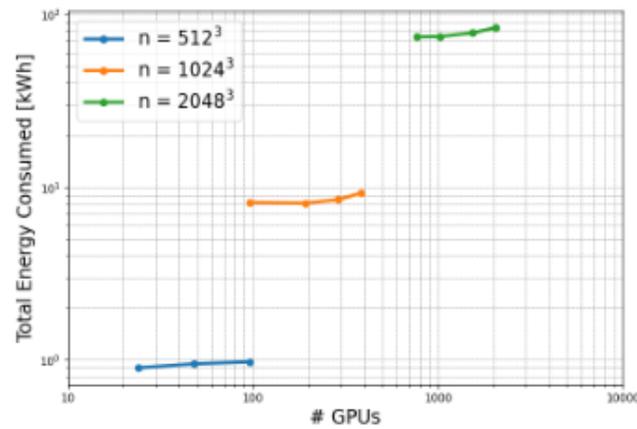## From "time to solution" to "energy to solution"

The thermal powers of GPUs will tend to 1,600 W by 2027



### IEA report 2024 for the EU

# DATA PRODUCTION



RADIO ASTRONOMY DATA VOLUMES (TB)



SKA Science Archive

Different technologies can produce short or long reads of ~ 10 Tb in 24- 72 hours



Sequence Read Archive (SRA) growth

NIH short read archive database growth

# SIESTA:
# A CODE FOR AB INITIO MOLECULAR DYNAMICS

Based on density-functional theory (DFT)

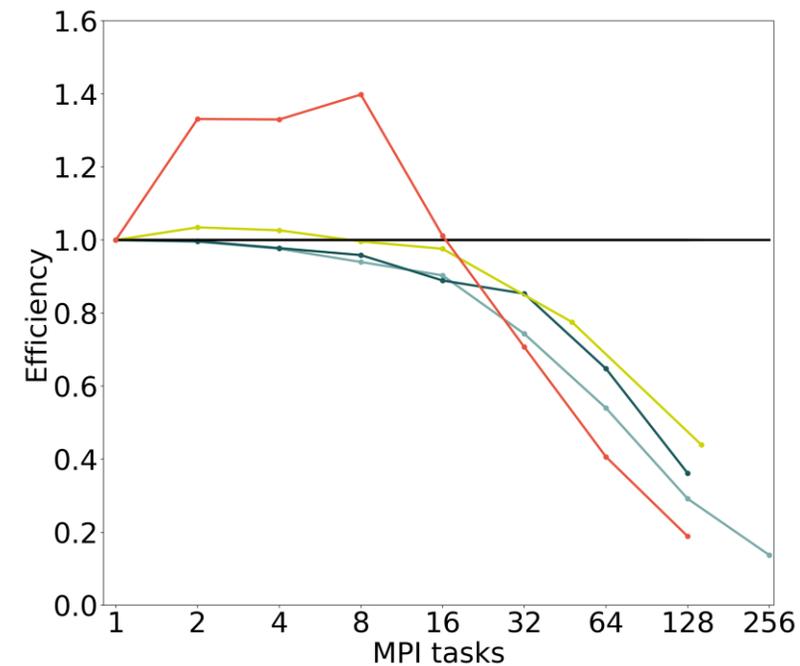Written in Fortran
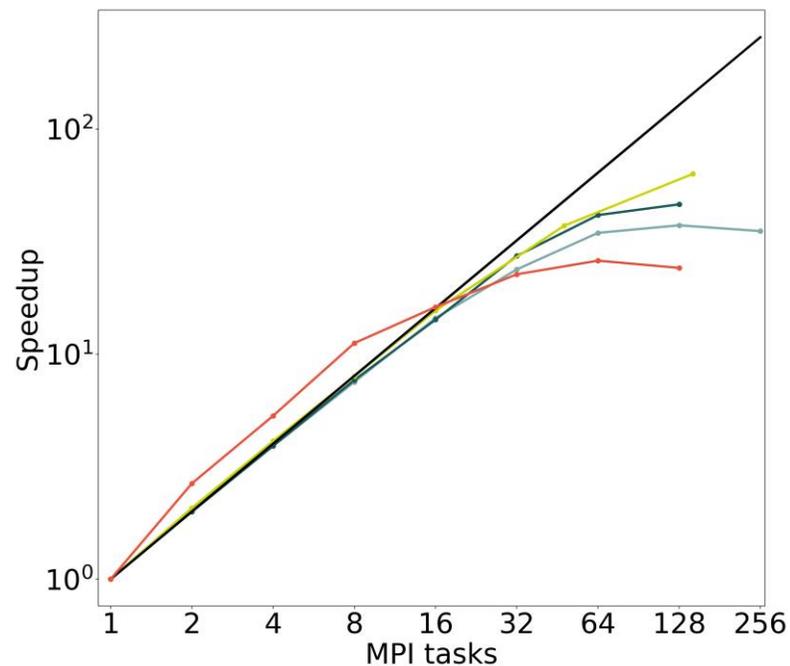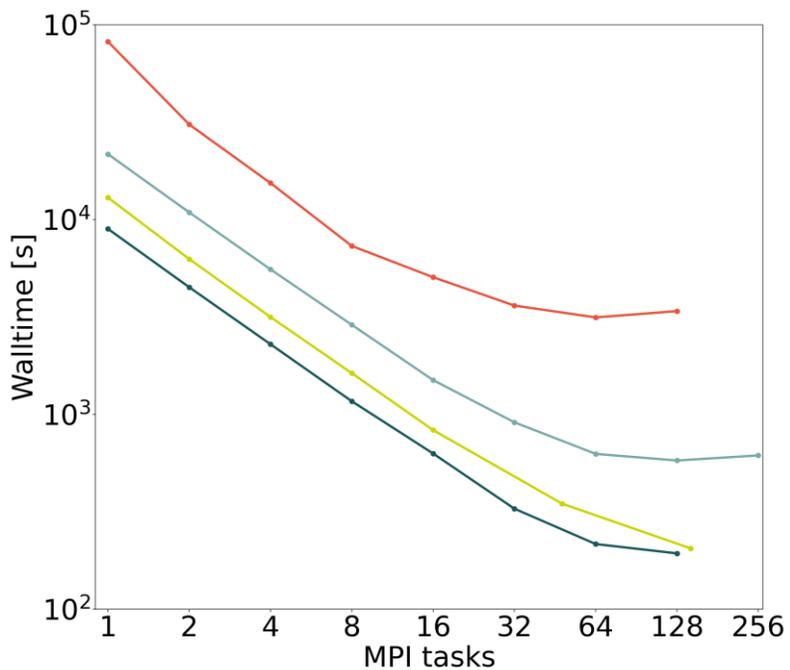
Parallelised with MPI and OpenMP
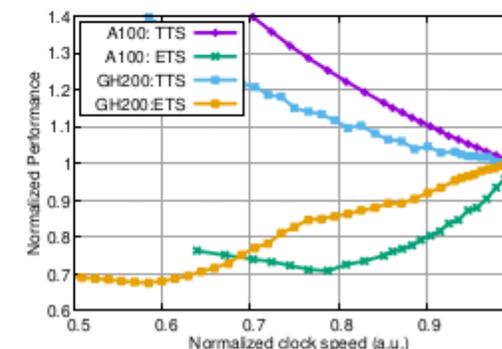
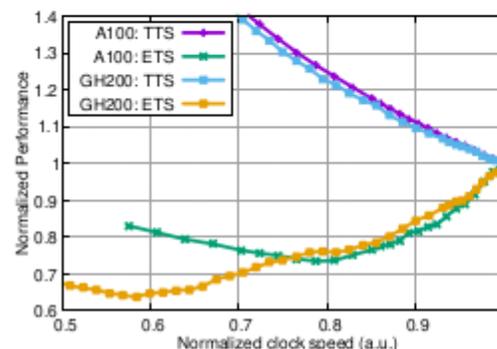Uses BLAS, LaPACK, ScaLAPACK, and ELSI



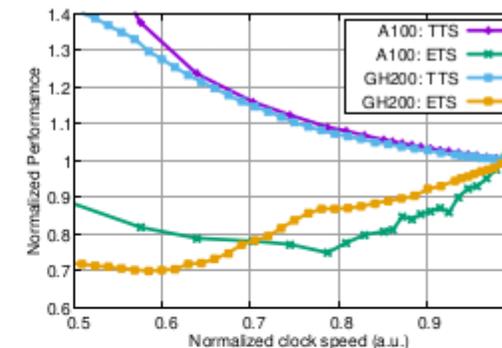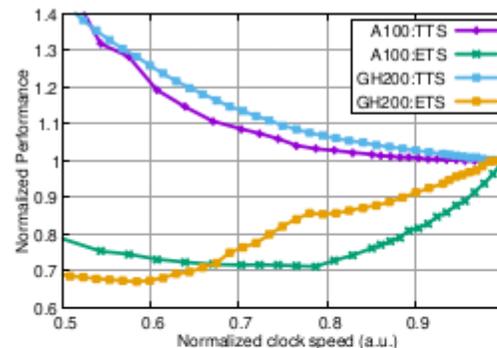P. Ordejon, SIESTA (nano) TUTORIAL (2010).

Quantum dot
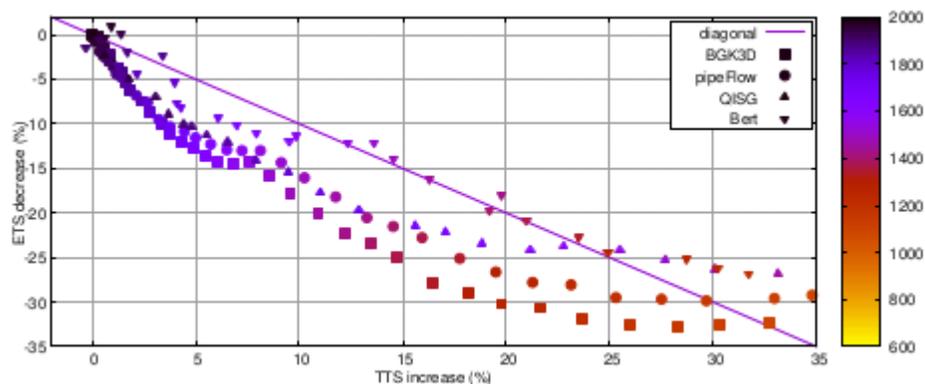gcc 8.5.0 – 13.2.1 + OpenMPI 4.1.4



-- SOPHON SG2042, -- AMD GENOA, -- AMPERE ALTRA MAX, -- NVIDIA GRACE

Energy-to-Solution full node estimate: AMD GENOA: 139.5 kJ, NVIDIA GRACE: 102.7 kJ

- **BGK3D**: CFD / LB; F90 / OpenACC / MPI; memory bandwidth-limited code. Test: 3D lid-driven cavity flow, 512 x 512 x 512 grid, *Re* = 256
- **PipeFlow**: CFD code; CUDA Fortran / cuTensor / MPI; memory bandwidth-limited. Test: incompressible fluid flow in a circular pipe, *Re* = 2.85 × 10$^5$ .
- **QISG**: MC of 2D ising quantum spin glass; C / CUDA; memory bandwidth-limited
- **BERT**: transfomer; Python / TensorFlow; compute bound. Test: 25 epochs  bert-large-un- cased variant (~340M parameters), batch size of 256.
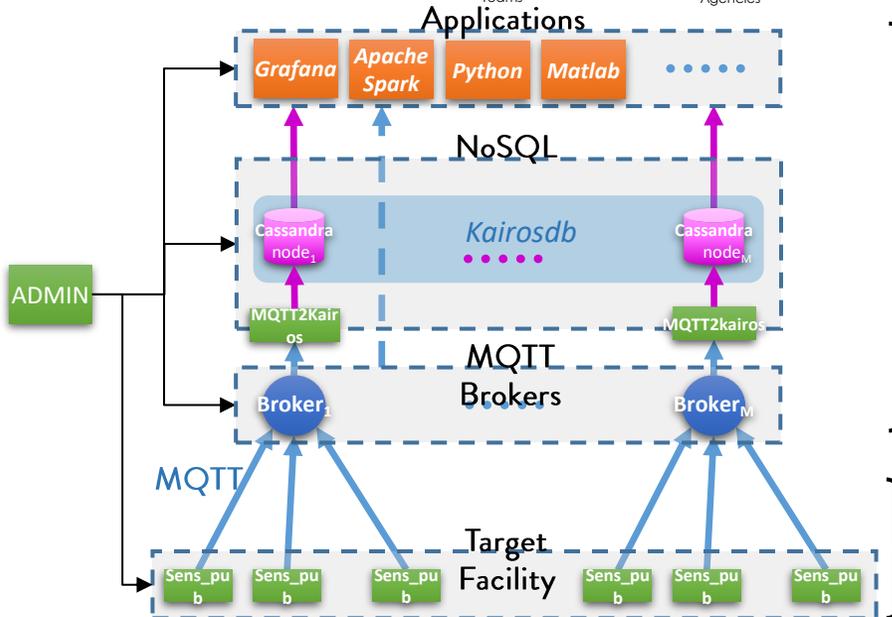


ETS and TTS versus clock speed

GH200 ETS decrease vs TTS increase. Colors indicates the clock speed in MHz. Markers refers to the four codes under test.

*"Experience on clock rate adjustment for energy-efficient GPU-accelerated real-world codes"* — Giorgio Amati, Matteo Turisini, Andrea Monterubbiano, Mattia Paladino, Elisabetta Boella, Daniele Gregori, and Danilo Croce

Community Initiative: https://examonhpc.github.io/examon

Current Installations:
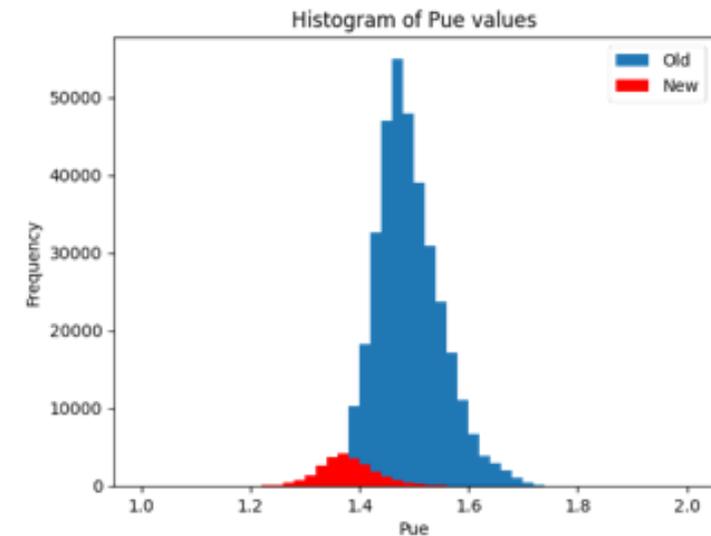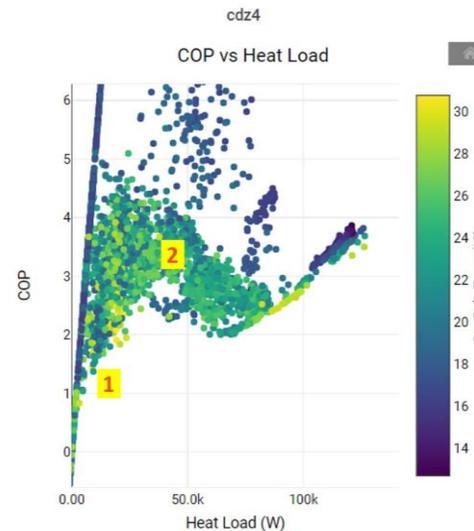- ENI: HPC4 & HPC5 Clusters, DaVinci 1



**Front-end**
- **MQTT** Brokers
- Data **Visualization**
- **NoSQL** Storage
- **Big Data** Analytics

**Back-end**
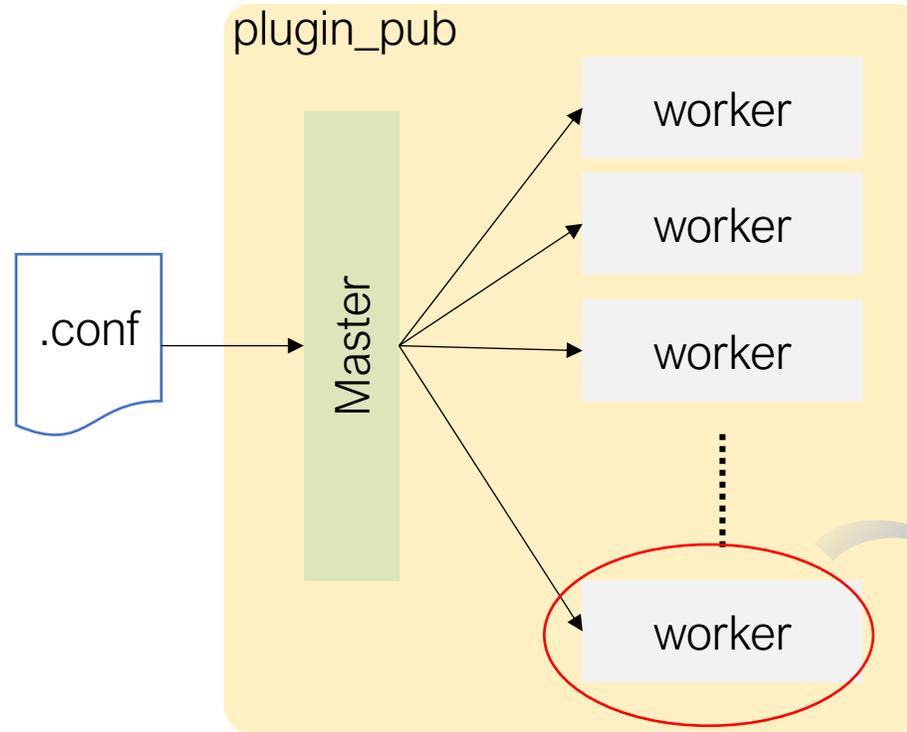- **MQTT**-enabled sensor collectors

# ExaMon Plugins: Internals

- Features:
  - Multithread/ multiprocessing
  - Internal scheduler / retry logic
  - logging

.conf

**plugin_pub**

Master

worker

worker

worker

worker

- Sensor API
  - Slurm, IPMI,...
- Publisher API
  - MQTT
  - HTTP
  - ....

Data Source

Sensor API

Examon (Publisher API)

HTTP/MQTT stream

Da Vinci POC: overview

# EnelX – Power/Energy



Ts=15min

# EnelX – Carbon Emissions

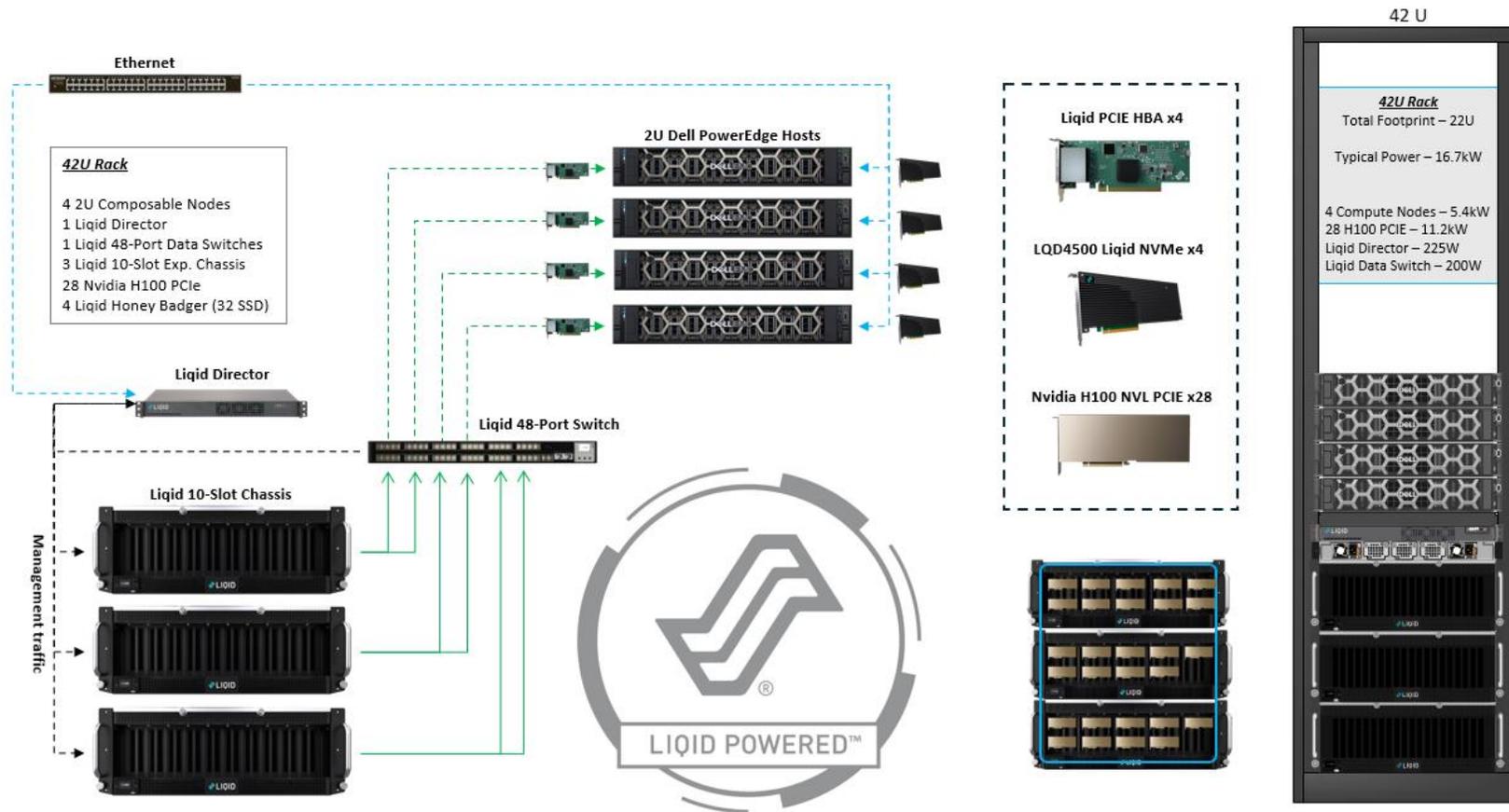Ts=1day

Liqid system installed in E4 Data Center and on customer's infrastructure
Different GPUs mounted (various models of NVIDIA and AMD
Connetected to ExaMon
Planned activity:
- Standard benchmark (e. g. HPL) and dedicated workloads (e. g. GROMACS)
- Check flop per watt tradeoff



**Ethernet**

**2U Dell PowerEdge Hosts**

*42U Rack*

4 2U Composable Nodes
1 Liqid Director
1 Liqid 48-Port Data Switches
3 Liqid 10-Slot Exp. Chassis
28 Nvidia H100 PCIe
4 Liqid Honey Badger (32 SSD)

**Liqid Director**

**Liqid 48-Port Switch**

**Liqid 10-Slot Chassis**

Management traffic

**LIQID POWERED™**

Liqid PCIE HBA x4

LQD4500 Liqid NVMe x4

Nvidia H100 NVL PCIE x28

42 U

*42U Rack*
Total Footprint – 22U

Typical Power – 16.7kW

4 Compute Nodes – 5.4kW
28 H100 PCIE – 11.2kW
Liqid Director – 225W
Liqid Data Switch – 200W

*Power Consumption estimates based on typical power usage*

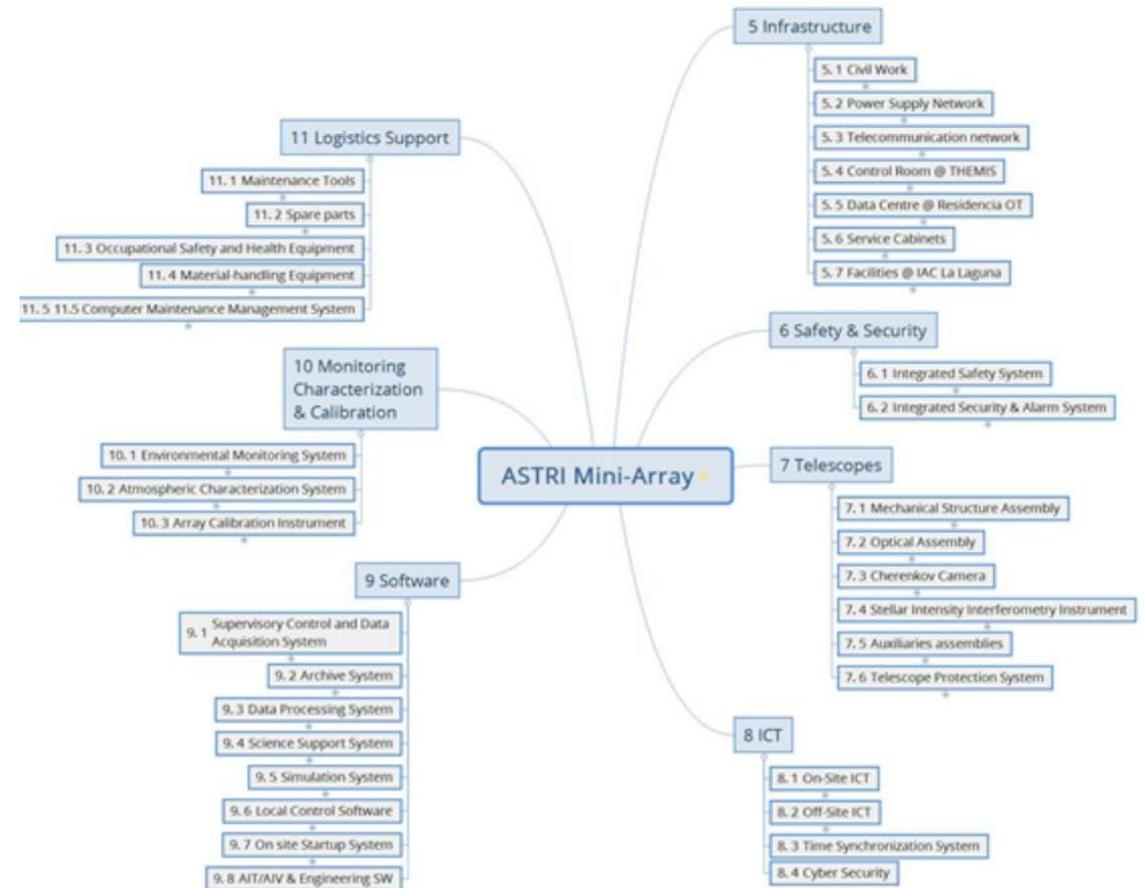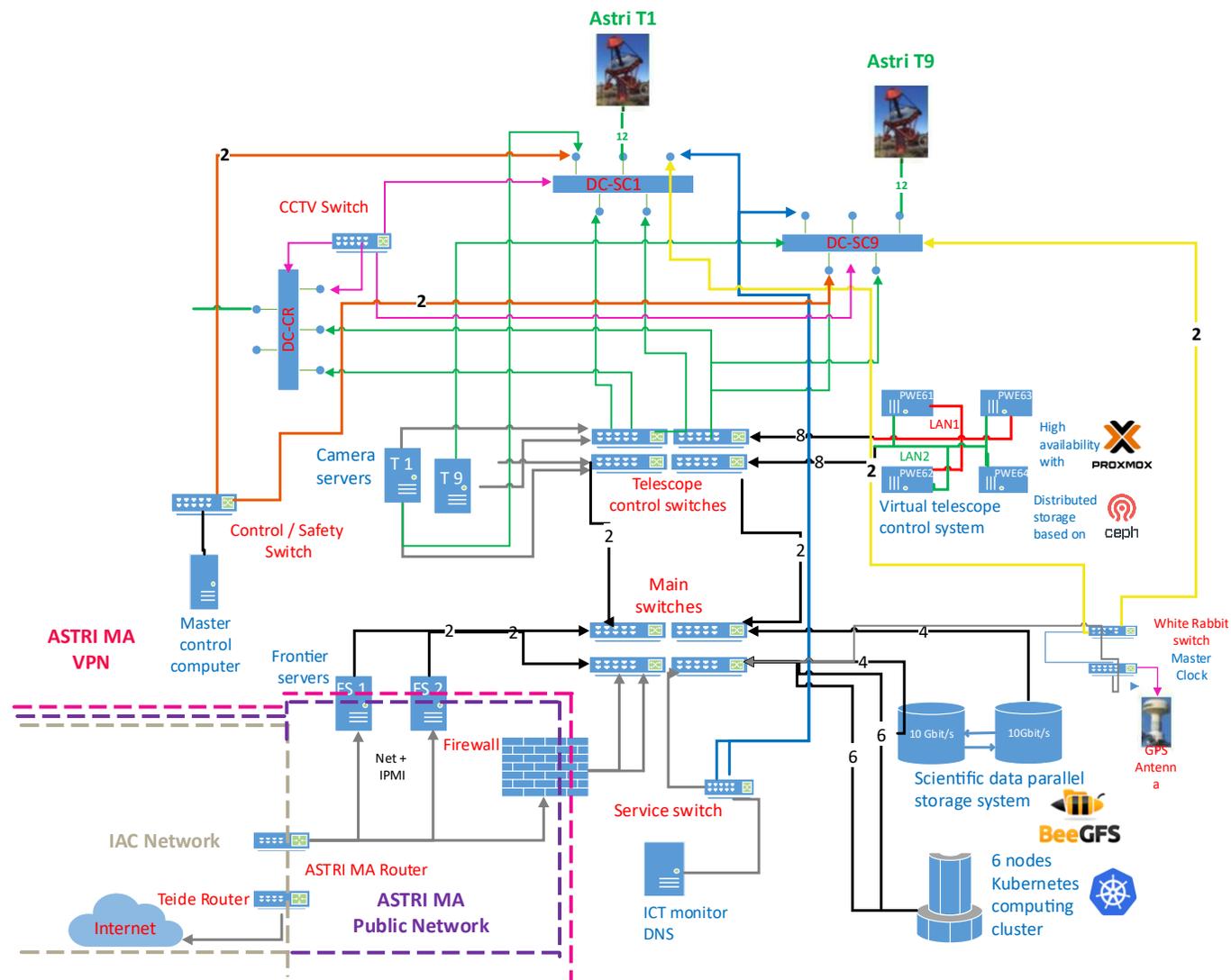| Year | 2024-2025 (active) |
|---|---|
| Customer | INAF @Teide observatory |
| Solution | Design, Installation |
| | The **ASTRI Mini-Array** is an **INAF**-led project to build nine double-mirror Čerenkov telescopes to study astrophysical sources emitting high-energy gamma photons with the goal of **measuring the energy, direction, and arrival time of gamma photons arriving on Earth from astrophysical sources**, by exploiting the Čerenkov atmospheric imaging technique |
| | E4 has designed and installed an infrastructure to manage and monitor an array of nine Čerenkov telescopes and process the data collected in real time |

**Hardware and operational requirements:**

- Install a cluster to control the telescopes managed by a central control plane
- Install a data acquisition and preliminary analysis platform with the purpose of running tests on acquired data and transmit only relevant data
- Within the budget provide both computational power and resiliency to operate in a "semi unsupervised mode"
- Roll out the existing preliminary infrastructure while installing the new one
- Maximise ease of transistion by using as much possible software tools known to INAF staff
- *Assist the customer through the whole process up to the final run of the test bed in Tenerife*
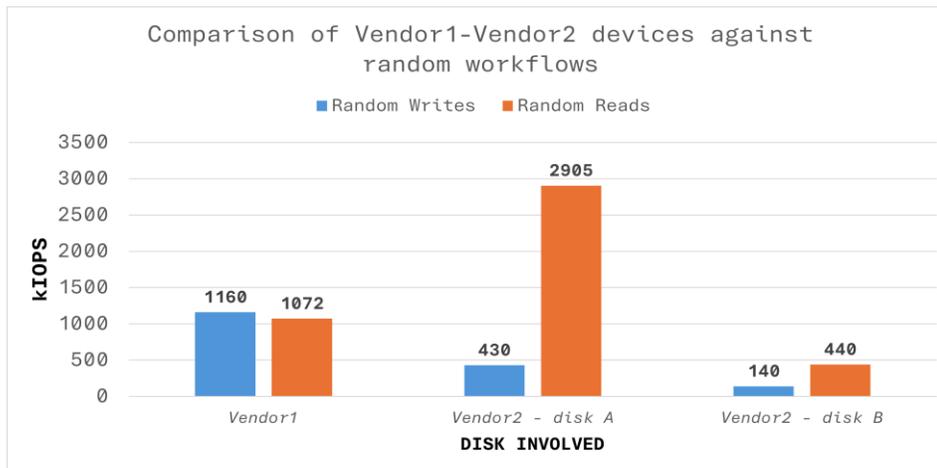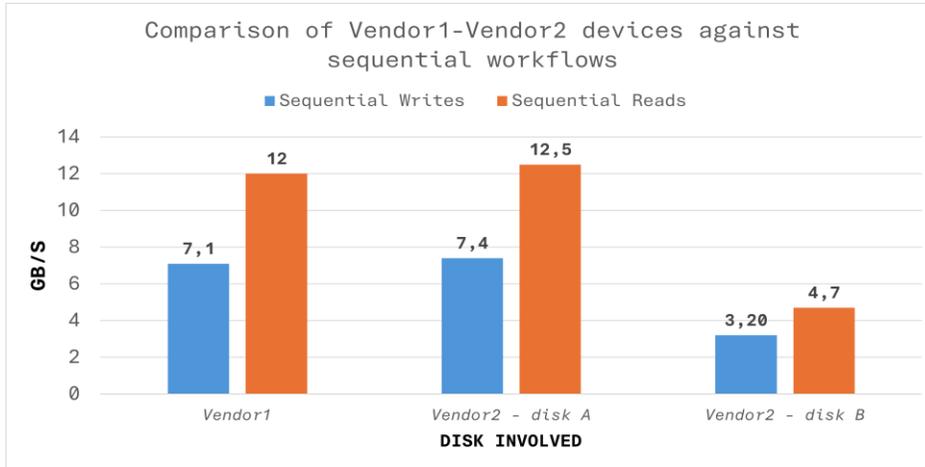
**Key points**

- Create a control plane based on ProxMox for simplicity but using OpenStack components (Ironic, Bifrost) for the lifecycle management of the infrastructure
- Use Medooza components for setup
- Create a Kubernets / Rancher cluster to analyze data from telescopes
- Use BeeGFS to provide the storage; to contain costs and complication consolidates Storage and Metadata Services into two storage systems, while Management is virtualized. Redundancy is ensured through Buddy Mirroring for both data and metadata servers.
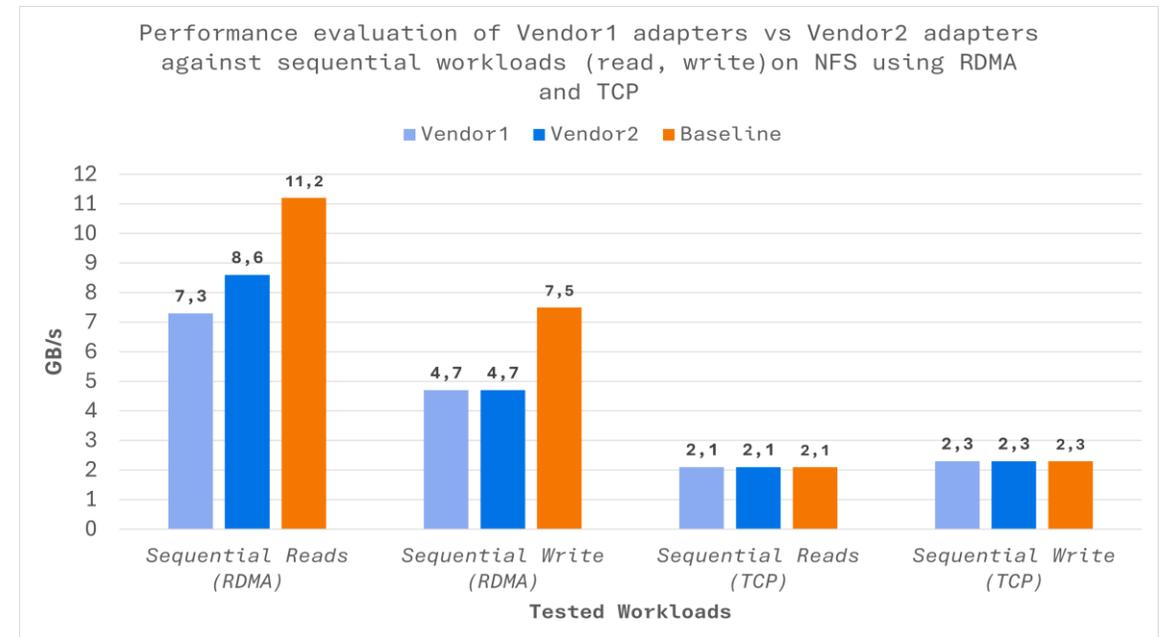
# Evaluating NVME and NICs performance under different conditions

Test platform:

| | |
|---|---|
| OS | AlmaLinux 9.5 (Teal Serval) |
| Kernel_version | 5.14.0-503.31.1.el9_5.x86_64 |
| FIO_version | fio-3.35 |
| mdRAID_version | v4.3 |
| NVME-CLI_version | 2.9.1 |

NFS preliminary results carried out from VMs and baremetal nodes

## Comparison of Vendor1-Vendor2 devices against sequential workflows

■ Sequential Writes  ■ Sequential Reads



## Comparison of Vendor1-Vendor2 devices against random workflows

■ Random Writes  ■ Random Reads



## Performance evaluation of Vendor1 adapters vs Vendor2 adapters against sequential workloads (read, write) on NFS using RDMA and TCP

■ Vendor1  ■ Vendor2  ■ Baseline

# STORAGE SOLUTIONS

Among the first to Implement NVMe RAID Solutions for the Beegfs Parallel Filesystem!

**Automating:** Public Release of the Ansible Collection for Beegfs!

Fair CPU usage with DRBD CPU masking

# ARGUS
## STORAGE MONITORING & ALERTING SOLUTION

The upcoming release of Argus will allow observing the relative consumption of flash cells present in NVMe drives and thus predicting the lifespan of the drive and related carbon emissions: "Operational cost" and "Use phase", in line with the OCP sustainability initiative.



**Key features:**
It will extend to monitoring the entire E4's Medooza HPC solution & will integrate with Examon for metrics ingestion.

Installed 2021; starting situation - 2023

- CentOS 7.9
- 60 GPU nodes
- Intel based nodes
- 2 CPU + visualization nodes
- Two BeeGFS nodes
- xCAT based configuration control plane
- OpenPBS in HA
- Infiniband and 10 Gbit/s network
- Zabbix monitor
- Several software and software versions installed over the years

# FRANKLIN CLUSTER @IIT (2)

MEDOOZA  BeeGFS®  OpenPBS

## RESOURCES SELECTION



### After upgrade

- Expansion of networking
- OpenStack Hyperconverged control plane on 3 nodes
- RockyLinux 8.9
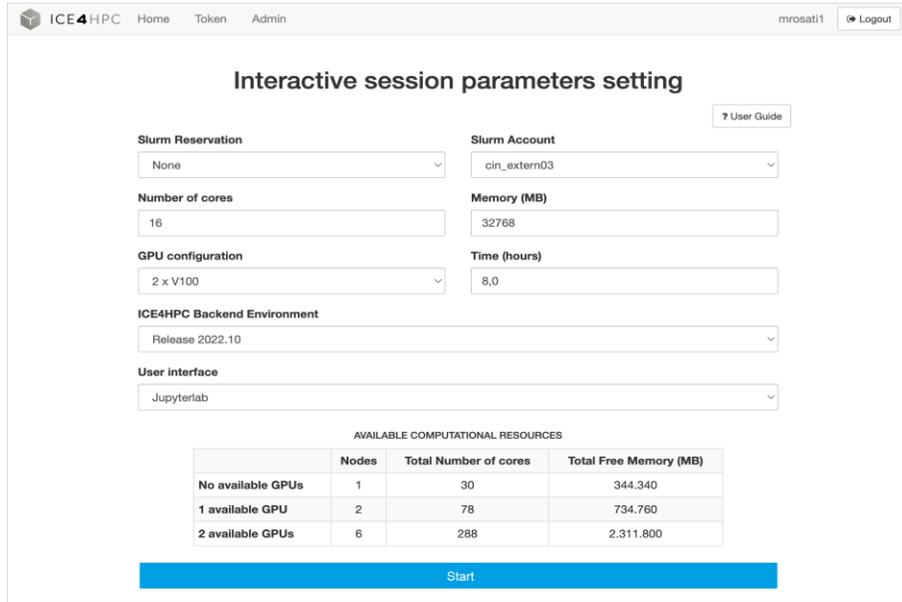- ~ 100 CPU and GPU nodes
- Intel and AMD mixed
- Interactive computing environment
- All software modules reinstalled with new stack and on different architectures

- Latest nodes with A100 GPUs needed by Genomics and Molecular Dynamics groups
- Baremetal & scheduled jobs vs container based workloads
- Strict isolation policies asked by DPO for Genomics data

- 100 nodes of mixed generation
- BeeGFS Storage
- Three node Open Stack control plane

- Cloud infrastructure allows to **reinstall with a different setup** nodes on demand
- Assign them to different networks
- Enforce firewall or other security measures

# MEDOOZA
## CURRENT REFERENCE ARCHITECTURE

**E4** COMPUTER ENGINEERING

MEDOOZA = OpenStack + OpenTofu / OpenBao + TALOS + A + ARGUS + EXAMON + + ICE4HPC

ACCESS NETWORK
NETWORKS
ACADEMIC / BUSINESS

**HPC NODES**
BAREMETALS AS A SERVICE

FAT NODES

openstack

CONTROL NODE
CONTROL NODE
CONTROL NODE

OPENSTACK PRIVATE API

CPU NODES

VISUALIZATION | INTERACTIVE | GPU NODES
ICE4HPC

ODS CEPH NODES

CEPH NODE

LOW LATENCY NETWORK
ACCESS NETWORK

COMPUTE NODES

COMPUTE NODE

PARALLEL FILE SYSTEM | ARCHIVE STORAGE

DATA | DATA

MANAGEMENT

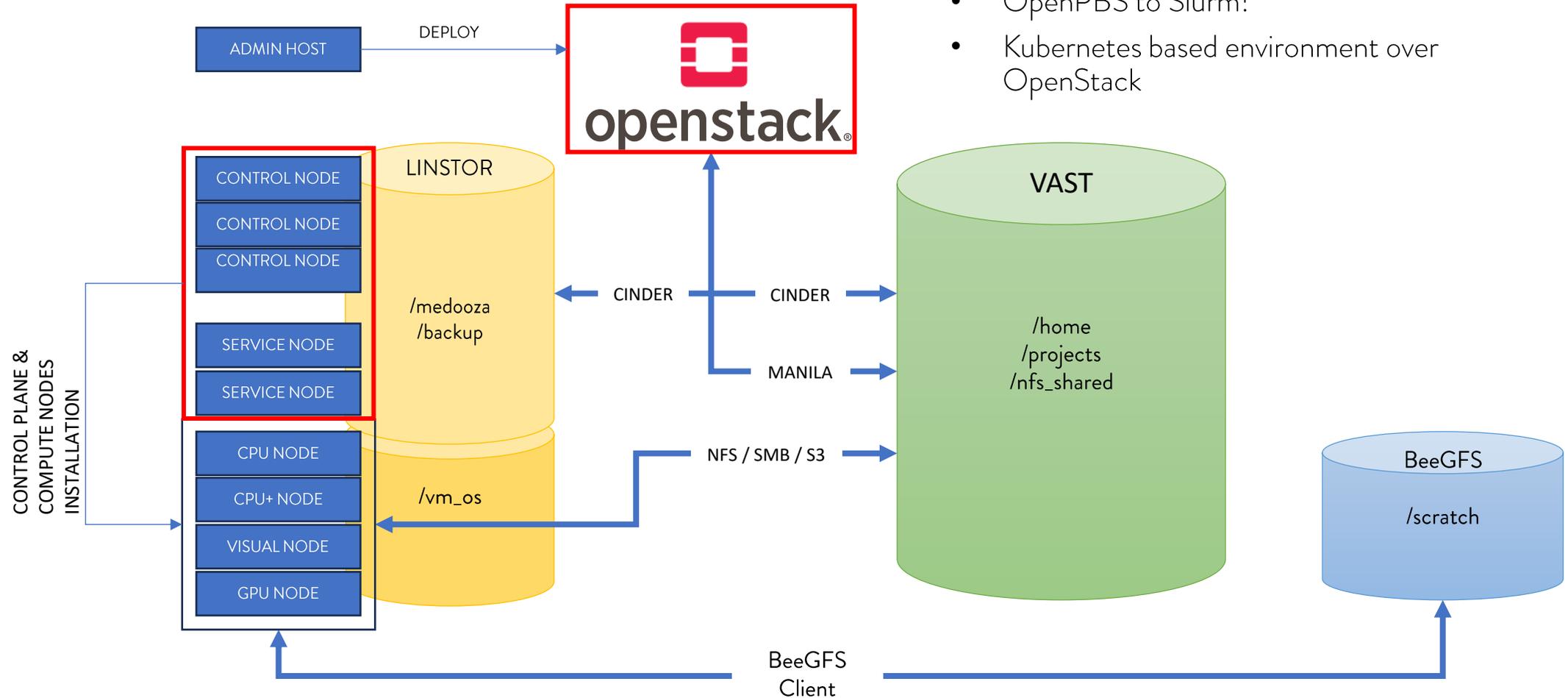**Medooza Infra & Medooza Ansible deployment solutions**

- Create management VMs on the OpenStack control plane ( e. g. Scheduler, BeeGFS)

- User's secrets management during deployment

- Bootstrap Ansible installation

- Additional OS packages, basic settings

- Drivers

- Configuration of services on VMs

- Base stack for including compilers, MPI and mathematical libraries, container stack and Python/R environments. Each release is tested against a set of versions and options. It's basic environment the admins find on login.

- Compilers and mathematical libraries from scratch

- Release tests for the HPC Cluster from services to job submission and benchmarks

- Rocky 8.x, 9.x, Ubuntu Jammy and Noble

- Update plan: 1 release per year, ships a validated set of components and versions, support up to 3/5 years
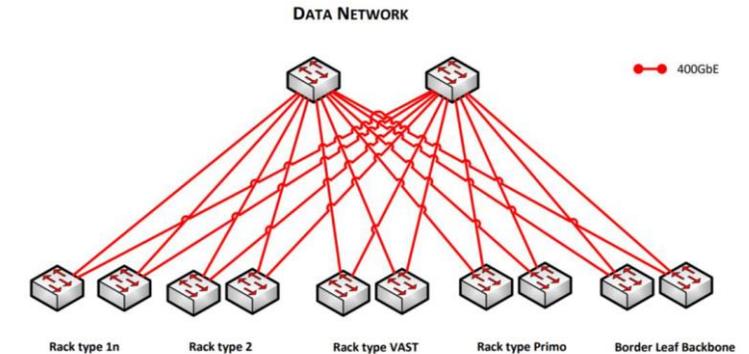
# NEW CLUSTER 2025

- Completely new cluster
- Standard Open Stack deployment with Kayobe
- Upgrade of Franklin after installation
- OpenPBS to Slurm?
- Kubernetes based environment over OpenStack

OpenStack Cluster:

- **420 CPU node DELL** PowerEdge **R670** Server CSP Edition with 2x CPU **Intel Xeon Sierra Forest 144 core**
- **80 GPU node** with NVIDIA GPU **112x L40S + 80x A30 + 8x H100**
- **4 GPU Server with 8x NVIDIA Hopper H100**
- **62 node** Ceph / Management / Service
- **400 GbE spine/leaf network** with DELL Powerswitch Z9432F-ON + Z9664
- **50 PB full flash storage Datalake**
- **4.35PB Disk-to-Disk Long Term Storage**

- **Examon HPC Monitoring (Exascale Monitoring) for** Data Collection and Analysis.
- E4 Medooza Cluster Management Suite on the c. p.
- Complete system installed with OpenStack Kayobe
- E4 on-site permanent support with 60 minutes intervention for critical incident
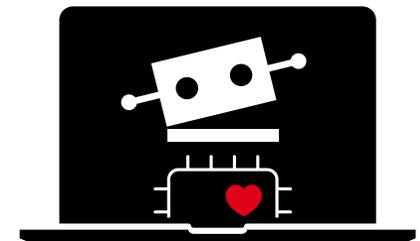
**TOTAL:** 131.200 physical cores

# SUMMARY & PERSPECTIVES

- Cloud based HPC scales from small to big infrasctures
- «One size fits all» is impossible but is important to focus on a few but flexible tools
- Small penalties in performance lead up to 20% energy savings

- Support for mixed architectures in Medooza (x86_64, ARM, RiscV)
- Deployment of MLOps/AIOps solutions in production in ExaMon

# THANK YOU FOR YOU ATTENTION!

This slide left blank

# EnelX – Room Meters
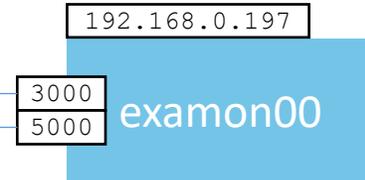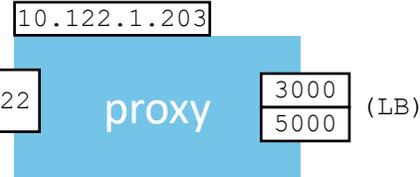
**E4** COMPUTER ENGINEERING

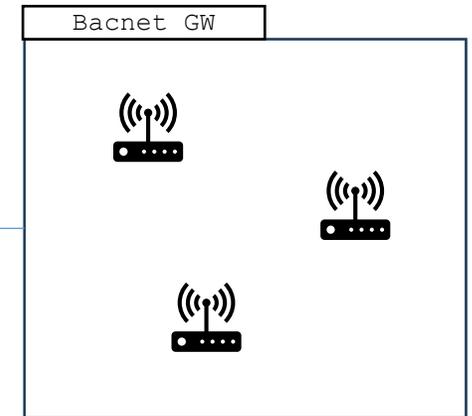**Citrix**

**Openstack Leonardo**

**EnelX**

VDI-HPC-DV1

53000
55000

localhost
Grafana:    53000
Examon API: 55000

10.122.1.203

22

proxy

3000
5000

(LB)

192.168.0.197

3000
5000

examon00

Grafana:    3000
Examon API: 5000
MQTT:       1883

Enelx Web API

enel x

Utente(email)

Password

Password dimenticata?          Italiano

Login d'accesso

Bacnet GW

## Goals

- Find anomalies, solve problems (the current project)
- Improve efficency, investigate job properties from metadata (possible?)
- Extended to new HW in a standard way

## Data flow (now / tomorrow):

- metrics ingested into Cassandra / IoTDB
- Streamed to Kafka for real-time processing
- written to S3 with iceberg
- Queried with Trino & Spark
- Visualized with Grafana
- Basic ML models trained on historical data for forecasting
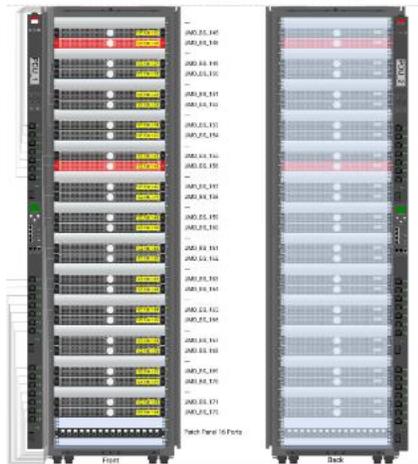- Integrated with foundation models & used with a conversational interface

## Tool

- Reference architecture for stand alone analytics platform: «ExaMon Box»

- Track and visualize rack-mounted equipment and available power, networks, and physical space
- Decommission equipment and connections fully with user-friendly DCIM software
- Quickly reconfigure devices and reserve equipment space via drag-and-drop
- Monitor capacity and changes at-a-glance with easy-to-understand color coding
- View floor plans by rack units or weight, and oversubscribe racks based on actual power

- Comprehensive reporting, including data center capacity, asset usage, and connectivity reports
- Manage work orders and tasks with details on due dates, owners, and statuses
- Change management
- Easily track rack utilization, weight, power, device types, etc.
- Stay organized with DCIM's endpoint connectors, link types, and cable usage information