# Status of HPC & Gateway

**R. Kamendje**

**EUROfusion**

**ENEA**
Italian National Agency for New Technologies, Energy and Sustainable Economic Development

**CINECA** **SCAI**
SuperComputing Applications and Innovation

# Marconi, Leonardo & Gateway

| | CPU | GPU | Gateway | | CPU | GPU | TOTAL |
|---|---|---|---|---|---|---|---|
| | # nodes | # nodes | # compute nodes | | HPL (Pflops) | HPL (Pflops) | HPL (Pflops) |
| | Marconi | Leonardo | Gateway | | Marconi | Leonardo | TOTAL |
| From August 3, 2023 until July 31 2024 | 2848 | 72 | 88 | | 5.96 | 5.00 | 10.96 |
| From August 1, 2024 until February 2025 | 1424 | 100 | 88 | | 2.98 | 6.94 | 9.92 |



Marconi - SKL
partition (A3) -
1424 nodes
3.0 Pflops (HPL)



Leonardo
GPU partition (C2)
100 nodes
6.9 Pflops (HPL)

# Flooding of the CINECA data centre in the evening of October 19

- Due to heavy rain in the Bologna area, the CINECA data centre in Casalecchio di Reno (where Marconi and the Gateway are located) was flooded
  - 170 cm of water in the basement where most of the main electrical panels, UPS, and batteries are located – electrical equipment were severely damaged and irreparable
  - The flooding induced a short circuit that triggered the gas-based fire extinguishing system of the data centre which generated a sound wave in the computer room that severely damaged the disks containing EUROfusion data – scratch and work fs (located close to the nozzles of the extinguishing system)
  - Leonardo was not affected since it is hosted in the CINECA Technopolo data centre (access to Leonardo was however disturbed for a few days)
- As far as we understand, the strategy of CINECA for restoring access to systems affected by the flooding has been the following
  - When one power supply becomes available for a given system, testing and restoring of the associated storage is started
  - When two power supplies become available for a given system, the system is opened to users as soon as storage testing/restoring is complete

# Flooding of the CINECA data centre – Current situation

- ## Marconi

  - Access to Marconi login nodes and home file system restored on December 5

  - For **scratch** and **work file systems**, restoration work is still on-going – at this time, it is likely that all the storage will be lost due to the large quantity of disks that were damaged

  - It was agreed **not to restart Marconi and to provide** (as compensation) **nodes of Leonardo-CPU partition ("Data Centric")**

- ## Gateway

  - Restored during the second half of December
    - Mid December for login nodes,
    - End December for compute nodes

  - Access to all storage restored except DRES/GATEWAYDB_FS (simdb and simdb1 cannot access the DRES/GATEWDB area)

|  | Total Dimension (TB) | Quota (GB) | Notes |
|---|---|---|---|
| $HOME | 200 | 50 | • permanent/backed up, user specific, local |
| $CINECA_SCRATCH | 2.500 | no quota | • temporary, user specific, local<br>• no backup<br>• automatic cleaning procedure of data older than 40 days (time interval can be reduced in case of critical usage ratio of the area. In this case, users will be notified via HPC-News) |
| $WORK | 7.100 | 1.000 | • permanent, project specific, local<br>• no backup<br>• extensions can be considered if needed (mailto: superc@cineca.it) |

**Data Centric Module**
BullSequana X2140 three-node CPU Blades.
 Each computing node is composed of
- 2x Intel Sapphire Rapids, 56 cores, 4.8 GHz
- 512 (16 x 32) GB RAM DDR5 4800 MHz
- 3xNvidia HDR cards 1x100Gb/s cards
- 8 TB NVM

# Main KPIs

≈ 3% due to nodes out of order

Assuming 100% availability before the flooding

| MARCONI-SKL (A3) | | July 2024 | Aug. 2024 | Sept. 2024 | Oct. 2024 | Nov. 2024 | Dec. 2024 |
|---|---|---|---|---|---|---|---|
| Availability | A3+ | 90.8% | 95.9% | 96.9% | < 61.0% | 0.0% | 0.0% |
| | login | none | none | none | none | none | none |
| Usage | A3 | 94.6% | 86.1% | 88.7% | na | 0.0% | 0.0% |
| | | 85.0% | 85.0% | 85.0% | 85.0% | 85.0% | 85.0% |

| LEONARDO (C2) | | July 2024 | Aug. 2024 | Sept. 2024 | Oct. 2024 | Nov. 2024 | Dec. 2024 |
|---|---|---|---|---|---|---|---|
| Availability | C2+ | 100.0% | 100.0% | 100.0% | 100.0% | 100.0% | 100.0% |
| | login | 97.0% | 97.0% | 97.0% | 97.0% | 97.0% | 97.0% |
| Usage | C2+ | 74.9% | 83.1% | 85.5% | 89.5% | 85.0% | 90.5% |
| | | none | none | none | none | none | none |

| | | July 2024 | Aug. 2024 | Sept. 2024 | Oct. 2024 | Nov. 2024 | Dec. 2024 |
|---|---|---|---|---|---|---|---|
| Major incidents (from mor | | 1 | 3 | 0 | 4 | 1 | 0 |

| | | July 2024 | Aug. 2024 | Sept. 2024 | Oct. 2024 | Nov. 2024 | Dec. 2024 |
|---|---|---|---|---|---|---|---|
| Maintenance | A3 | 4.7 | 8.5 | 3.3 | 0.0 | 0.0 | 0.0 |
| | C2 | 17 | 22.8 | 14.3 | 0.0 | 42.5 | 10.0 |
| | A3 | 167.7 | 167.2 | 167.2 | 167.2 | 167.2 | 167.2 |
| | C2 | 76 | 61.2 | 54.9 | 54.9 | 20.4 | 18.4 |

D_2024.09                    D_2024.12

## Availability

- Very good availability of the C2 partition (Leonardo)

- A3 partition (Marconi) fully unavailable since the flood on October 19 - No statistics on the availability the partition before the flood
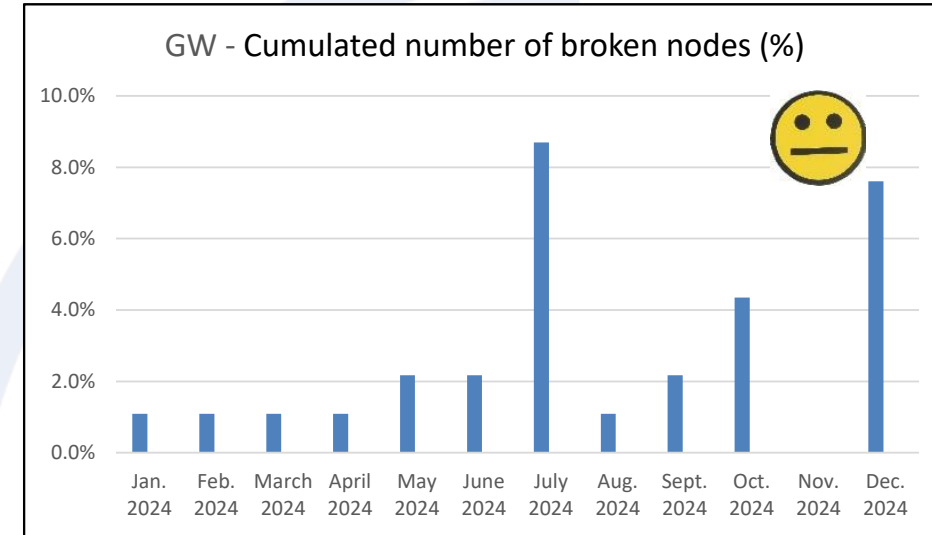
## Usage

- Very high for Leonardo leading to very long waiting times
- The Operations Committee decided to limit to 4 the number of job one user can execute at the same time

# Gateway Operation, Availability

- Computing resources
  - **92 SKL nodes (incl. 4 login nodes) as independent HPC system based on**
    - 2 x CPU Intel Xeon Platinum 8160c (SKL) with 24 core at 2.1 GHz
    - 1 x 240 Gib SSD hard disk
    - 1 x link OPA (OmniPath) 100Gbps, 1 Link 1GbE, 2 link 10GbE
  - 12 management nodes
  - Private storage (3.1 PB)

- Availability affected by the flooding



GW - Cumulated number of broken nodes (%)

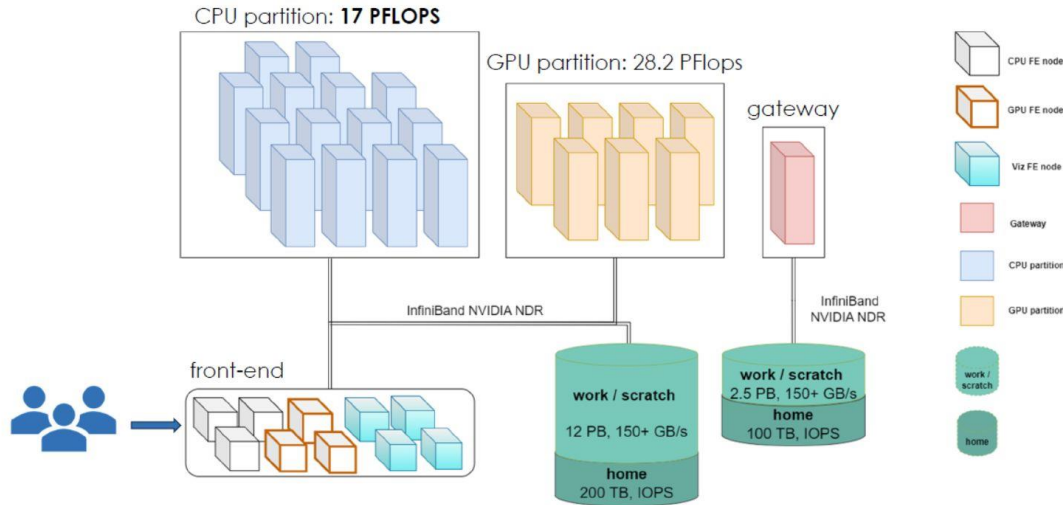| GW | | July 2024 | Aug. 2024 | Sept. 2024 | Oct. 2024 | Nov. 2024 | Dec. 2024 |
|---|---|---|---|---|---|---|---|
| **Availability** | GW | 92.6% | 98.8% | 98.6% | 60.3% | 0.0% | 13.1% |
| | | 97.0% | 97.0% | 97.0% | 97.0% | 97.0% | 97.0% |
| **Usage** | GW | 77.6% | 46.4% | 49.0% | 67.6% | na | 13.8% |
| | | none | none | none | none | none | none |
| **Incidents** | GW | 0 | 2 | 0 | 1 | 1 | |
| **Maintenance** | GW | 0 | 7.5 | 0 | 0 | 0 | 0 |

# Situation with Pitagora

- The flooding did only marginally affect the infrastructure equipment prepared for Pitagora

- Lenovo had to solve one hardware issue affecting AMD Turin nodes that further delayed the delivery of CPU nodes (CPU partition and CPU nodes of the Gateway)
  - Pitagora is the first machine worldwide installed by Lenovo based on AMD Turin processors





**4 racks of AMD Turin CPU nodes already on site**



**All racks of GPU nodes already on site**
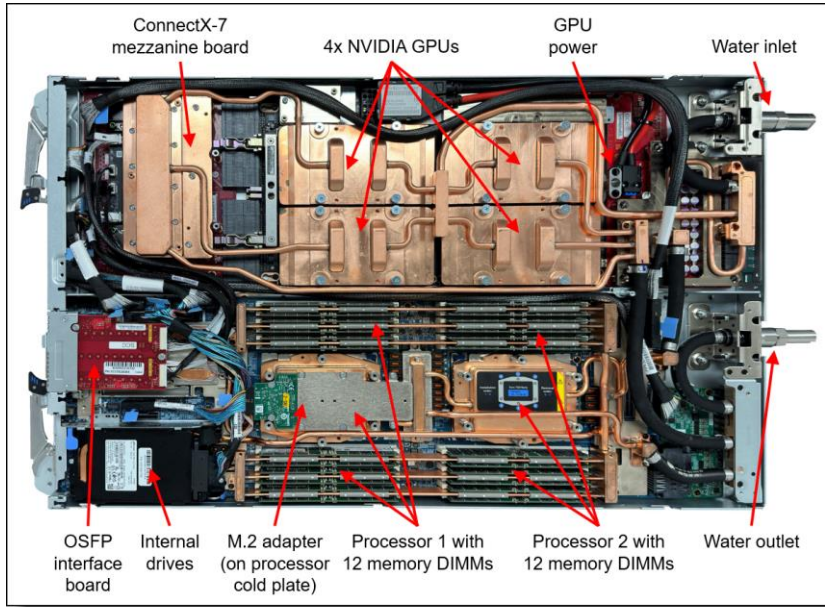
## CPU partition

## Gateway

- 14 racks
- **17 PFlops (Rmax)**
- 1008 Compute nodes
  - 72 CN x rack (6x6 tray enclosures)
  - 2x AMD Turin 128c (Zen5 microarch) 2.3 GHz
  - **1.7X performance** over Genoa
  - 768 GB DDR5 6400 MT/s
  - 3GB DDR5 per core, **1.33x higher bdw**
  - Integrator w/ earliest availability from AMD Turin
  - 1x NDR200 adapter
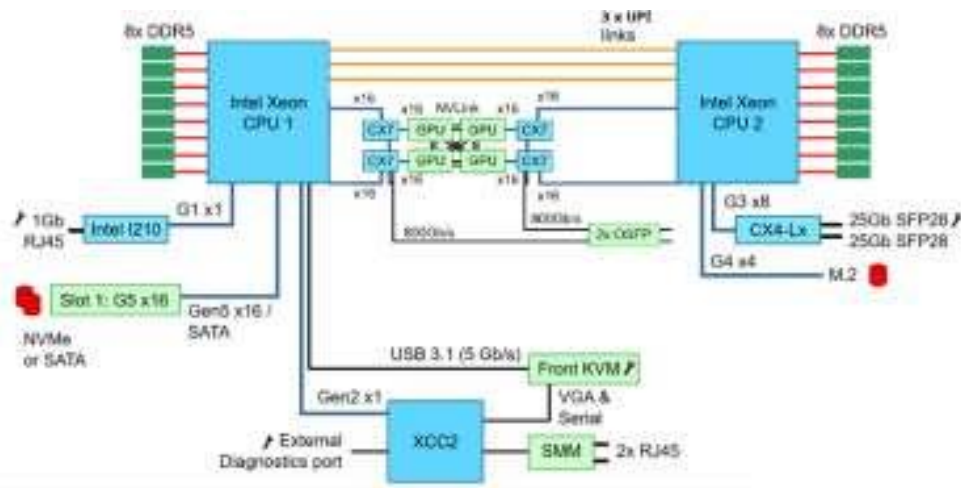- Full DLC (**97%** heat removal)

- The proposed GW system is a **fully separated** system
  - **Separated fabric** from main SC
- CPU nodes (14x **identical** to main SC)
  - 4 nodes 24*48 = **1TB DDR5**
- 1x GPU nodes (GPUs **identical** to main SC)
  - CPU AMD EPYC 9354 32C 280W 3.25GHz + 24 x 4800Mhz DIMMs
- VIZ nodes 4x (**identical** to main SC)
  - 2x NVIDIA L40s GPU
- Service nodes (3 masters + 3 workers) + service storage
- **Dedicated** networks, both IB (2 x QM9790) and Ethernet
- **Dedicated** storage
- **Dedicated** management stack
  - identical to the main SC

# GPU Configuration: Pitagora





- 7 racks

- **28.2 PFlops (Rmax) -> GPU at 700 W**

- **700 W -> + 7% HPL (and nothing else)**

- 168 Compute nodes
  - **24 CN** x rack (6x4 tray enclosures)
  - 2x Intel Emerald Rapids 32c
  - 512 GB DDR5 6400 MT/s
  - **4x NVIDIA H100 SXM 94GB HBM2e**
  - 2.3x performance over A100
  - Expected to run at 600 W (w.r.t max 700 W)
  - 1TB NVMe
  - 4x NDR200 adapters

# Pitagora – Current situation

- **GPU partition**
  - All the nodes are on site since summer 2024
  - Rank #44 in the Top500 of November 2024
    - Even if EUROfusion indicated to CINECA that the Top500 was not the priority

| 44 | **PITAGORA** - ThinkSystem SD650-N V3, Xeon Gold 6548Y+ 32C 2.5GHz, NVIDIA H100, Infiniband NDR, Linux, Lenovo CINECA Italy | 55,104 | 27.27 |

- **CPU partition**
  - 4 racks delivered in the second week of January – all racks supposed to be delivered by mid-February

- **Gateway**
  - Delivered to CINECA in December – CPU nodes are missing and are expected in the next two weeks

# Pitagora – Plans for deployment

- Opening a new partition in production involves several steps:

  - Delivery of the hardware

  - Installation/testing of the hardware and configuration for acceptance by Lenovo

  - Acceptance/handover to CINECA

  - Personalization by the CINECA system team and the CINECA user support team – 2 or 3 weeks

  - Testing by EUROfusion ACHs (MPG and EPFL) for checking the partition can be open to production including validation by the Operations Committee – 2 weeks

  - Documentation and webinar for user access by CINECA

  - Opening to production

  - Webinar for best practices and optimization by ACHs

# Plans for the deployment of Pitagora

*Information communicated by CINECA (D. Galetti on January 20/21)*

## PITAGORA

**Milestones**

**GPU partition**

**31st of January:** handover Lenovo - Cineca

**20th of February:** pre-production start

**Mid March:** production start (TBC by Eurofusion)

**CPU Turin partition**

**17th of February:** delivery completed

**28th of Februrary:** physical installation completed

**31st of March:** handover Lenovo – Cineca

**15st of April:** pre-production start

**Mid  May**  production start (TBC by Eurofusion)

## Proposal by the Operations Committee

### GPU
20th of February: start of pre-production (testing by ACH)
Webinar by CINECA
6th of March: end of testing by ACH
Validation by Operations Committee
Mid-March: Production starts (tentative date to be communicated to users)

### CPU
15th of April: start of pre-production (testing by ACH)
Webinar by CINECA
30th of April: end of testing by ACH
Validation by Operations Committee
Mid-May: Production starts (tentative date to be communicated to users)

# Plans for the deployment of the new GATEWAY

*Information communicated by CINECA (D. Galetti on January 20/21)*



## new GATEWAY

### Milestones

**Gateway partition**

29th of January: Turin CPU delivery completed

           (some of them are expected during this week)

10th of February: physical installation completed (cabling part to do onsite)

Pre-production and production start dates TBD (plan details to be completed soon)

Cineca is pressing Lenovo to move up on gateway activities.

## Note by the Operations Committee

To be clarified as soon as possible

Note: The deployment of the Gateway requires the installation of a large number of software and services. Therefore, it needs to be carefully planned and started early
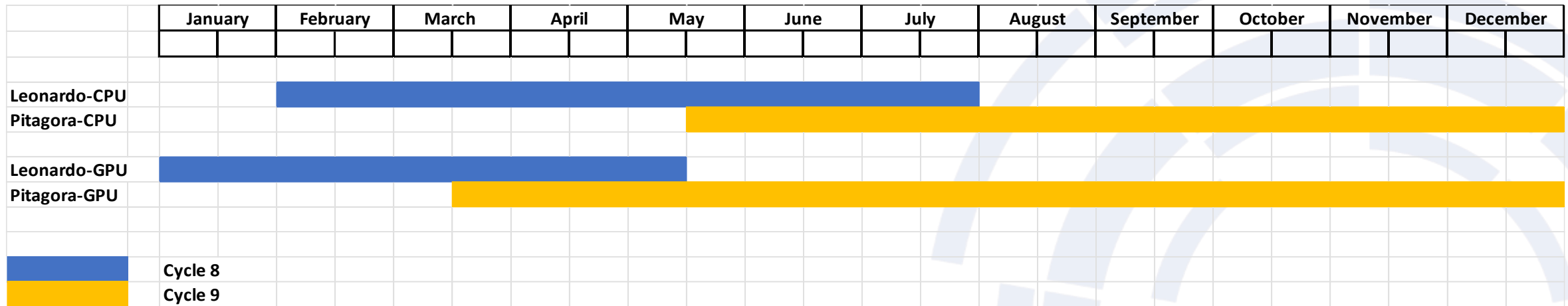
- ## CPU resources

  - **260 Leonardo-CPU nodes** are being made available to EUROfusion

  - The **performance ratio** (for the EUROfusion workload) between one node of Leonardo-CPU and one node of Marconi is difficult to figure out.

  - First results obtained by ACH-MPG as well as indications given by CINECA tend to converge to a figure between **2.2 and 3**.

  - ACH-MPG noted also that porting codes from Marconi was easy but that the stability of the Leonardo-CPU (in terms of operation and performances) needs to be improved.

  - ACH-EPFL performed tests on Leonardo-CPU and noted that "everything looks fine in terms of environment and performance".

  - If **2.5** is used for the conversion factor between CPU nodes of Leonardo-CPU and Marconi, then 260 nodes of Leonardo-CPU are equivalent to **650** nodes of Marconi (half the size of Marconi as it was in mid-2024).

# Next steps – Oganization of Cycles 8 & 9 of Production Runs

*Assumption: 2 months overlap between old and new partitions*

| | January | February | March | April | May | June | July | August | September | October | November | December |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | |

**Leonardo-CPU**

**Pitagora-CPU**

**Leonardo-GPU**

**Pitagora-GPU**

Cycle 8
Cycle 9

- ## Allocation of CPU resources

  - **Migration** of Cycle 8 projects from Marconi to Leonardo-CPU as of **28 January 2025**

  - **Extension** of Cycle 8 on Leonardo-CPU until end of July (providing compensation for the time lost by projects because of the unavailability of Marconi)

  - **Start** of Cycle 9 on Pitagora-CPU in **mid-May 2025**

- ## Allocation of GPU resources

  - **End** of Cycle 8 on Leonardo-GPU by mid-May 2025 (allowing for a 2 month overlap)

  - **Start** of Cycle 9 on Pitagora-GPU in mid-March 2025

**END**