# Guidelines for EuroHPC applications

relative to LUMI (CSC) and Leonardo (CINECA)

Additional details here

www.elliot-ai.eu

**MINERVA is an European Support Centre for Scalable AI Research and Deployment.**

- Establish and operate a Europe-wide Support Centre.
- Interact with AI communities to identify needs.
- Offer a rich service portfolio, aligned with the European needs.
- Ensure models are developed according to ethical and responsible AI regulations

# Agenda

- General guidelines + hands-on look at the forms (15 min)

- Technical specifications of LUMI and Leonardo (20 min)

- Q&A

# How to apply to EuroHPC Calls

**elliot**

| | Call name | GPU hours provided | cut-off | Time to provision* | Allocation time |
|---|---|---|---|---|---|
| **AIFactory** | Playground | 5.000 GPU hours | continuous | 2 working days | 1, 2, or 3 months |
| | Fast Lane | 10.000 - 50.000 GPU hours | continuous | 4 working days | max 3 months |
| | Large scale | > 50.000 GPU hours** | twice a month | 10 working days | 3, 6, or 12 months |
| **EuroHPC** | Science & EU projects | 20.000 - 90.000 GPU hours | CLOSED | 1 month | 6 months |
| | Benchmark Access | 3.500 | 1st day of each month | 2-3 weeks after the cut-off date | 3 months |
| | Development Access | 4.500 | 1st day of each month | 2-3 weeks after the cut-off date | 6 months - 1 year |
| | Regular Access | 20.000 - 150.000 | CLOSED | 4 months after the cut-off date | 1 year |
| | Extreme Scale Access | >= 245.000 | CLOSED | 6 months after the cut-off date | 1 year |
| | AI and Data Intensive | 50.000 | CLOSED | 1 month after the cut-off date | 1 year |

- Website: https://www.eurohpc-ju.europa.eu/index_en
- Here you can check the currently open calls: https://access.eurohpc-ju.europa.eu/
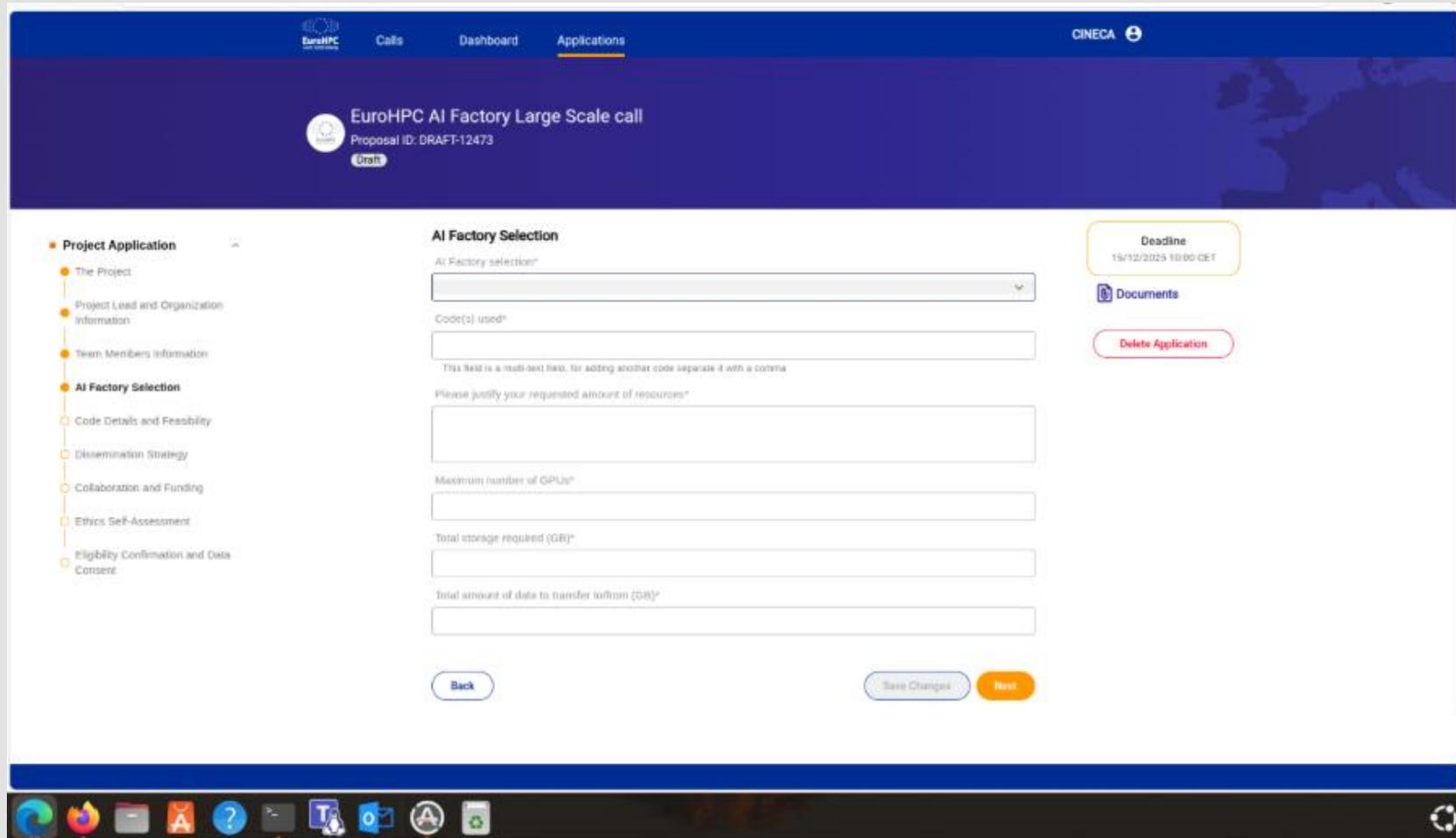
# General Guidelines - Before to apply

- GPU-hours = Number of GPUs x Duration (hours)
- Node-hours = GPU-hours / GPUs per node

Before applying, users should have an idea of:

- **Workflow Type:** The specific pipeline or process to be executed.

- **Data & Models:** Total volume of training+output data and the size of the AI models.

- **Time Estimates:** Estimated running time for a (non-optimized) version of the code.
- **Machine Architecture:** Hardware specifications of the chosen machine.
                                         (GPU model, GPUs per node, interconnects, storage, …)

# General Guidelines - How to apply

You will need to complete a form similar to this one. Let's scroll through it together.

https://access.eurohpc-ju.europa.eu/
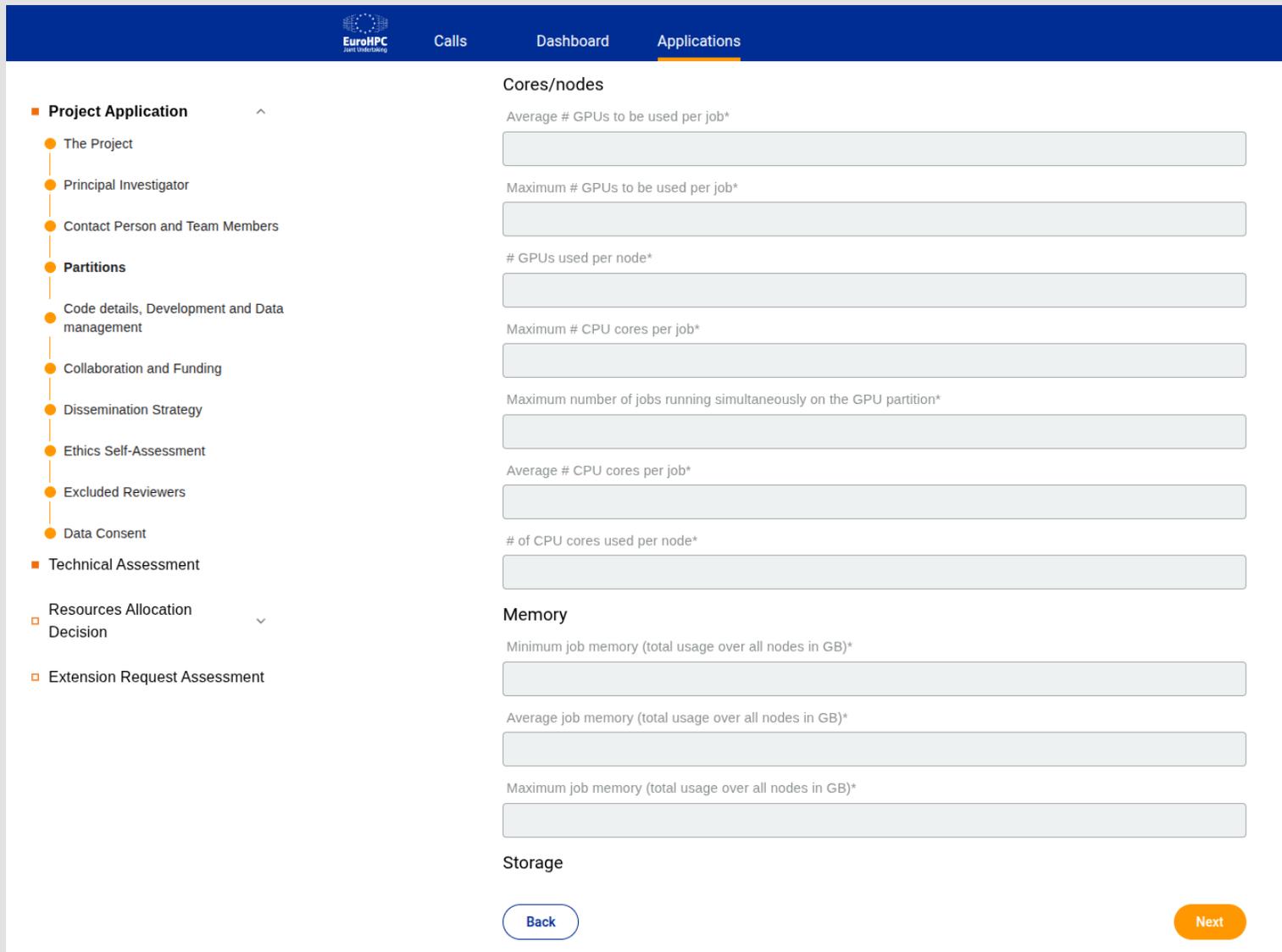
# Example from AIF Science and Collaborative EU Projects



Leonardo, LUMI, Marenostrum, Jupiter,....

Node-hours = Number of nodes x Duration (h)
Custom or fixed, depending from the call type

Recommended for long-running jobs.

# Example from AIF Science and Collaborative EU Projects



Refer to the chosen machine's technical specification!

More info in the next slides.

# Example from AIF Science and Collaborative EU Projects



e. elliot

**Storage**

Maximum amount of SCRATCH needed at a time (TB)*

Maximum amount of WORK needed at a time (TB)*

Maximum amount of HOME needed at a time (TB)*

Maximum amount of ARCHIVE needed at a time (TB)*

Maximum # files to be stored on SCRATCH (thousands)*

Maximum # files to be stored on WORK (thousands)*

Maximum # files to be stored on HOME (thousands)*

Maximum # files to be stored on ARCHIVE (thousands)*

Total amount of data to transfer to/from (TB)*  →  data to load at the beginning (to) **+** data do download at the end (from)

Justification of data transfer*  →  Explain the section above.

Back    Next

**Project Application**
- The Project
- Principal Investigator
- Contact Person and Team Members
- **Partitions**
- Code details, Development and Data management
- Collaboration and Funding
- Dissemination Strategy
- Ethics Self-Assessment
- Excluded Reviewers
- Data Consent

Technical Assessment

Resources Allocation Decision

Extension Request Assessment

EuroHPC Joint Undertaking | Calls | Dashboard | Applications

# Example from AIF Science and Collaborative EU Projects



report main software / libraries specifying the version if it is a strict requirement.

# Example from AIF Science and Collaborative EU Projects



If you have not scalability results can be cited the scalability of the distributed framework used (Megatron,FSDP,...)

If you will need technical help on HPC (code porting, multi-gpu/node running,...)

# Leonardo specs

# Leonardo ( hosted by CINECA, Italy )

- General info: https://www.hpc.cineca.it/systems/hardware/leonardo/

- 2 partitions:

    - DCGP   (CPU only)
    - BOOSTER (4 A100 - 64GB GPUs per node, 32 process/threads per node)

# "Partitions" section



max 256 concurrent nodes

24h single job wall time

~ 6 hours.

# "Partitions" section



max 4
(4xA100 per node on Leo. Booster)

max 256 concurrent nodes

max 32
(32 processes per node on Leo. Booster)

Memory per node:
512GB RAM + 64GB x 4 VRAM

# "Partitions" section

EuroHPC
Joint Undertaking

Calls | Dashboard | Applications

**Storage**

Maximum amount of SCRATCH needed at a time (TB)*

the limits depend on the call type

Maximum amount of WORK needed at a time (TB)*

MAX 100TB

Maximum amount of HOME needed at a time (TB)*

MAX 50GB

Maximum amount of ARCHIVE needed at a time (TB)*

Maximum # files to be stored on SCRATCH (thousands)*

Maximum # files to be stored on WORK (thousands)*

WORK and ARCHIVE remain active for 6 months after the project ends.

Maximum # files to be stored on HOME (thousands)*

Maximum # files to be stored on ARCHIVE (thousands)*

Total amount of data to transfer to/from (TB)*

Data transfer:
Consider max ~ 2TB per day.
Do the math considering also the project duration.

Justification of data transfer*

**Project Application** ⌃
- The Project
- Principal Investigator
- Contact Person and Team Members
- **Partitions**
- Code details, Development and Data management
- Collaboration and Funding
- Dissemination Strategy
- Ethics Self-Assessment
- Excluded Reviewers
- Data Consent
- Technical Assessment
- Resources Allocation Decision ⌄
- Extension Request Assessment

Back | Next

# "Code" section



report main software / libraries specifying the version if it is a strict requirement.

Environment:

you can use either:
- python env
- Leonardo modules (cineca-ai)
- **Singularity** containers (not Docker).

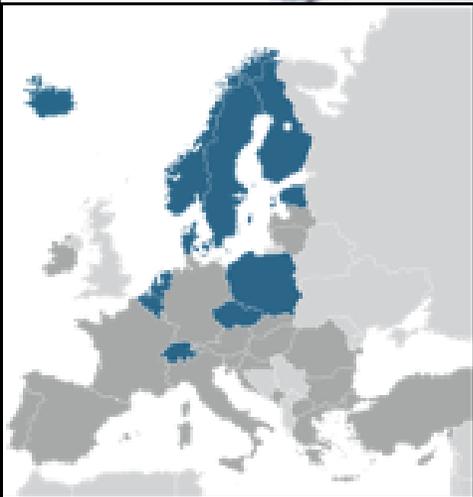# LUMI specs

# LUMI is an HPE Cray EX Supercomputer

- https://lumi-supercomputer.eu/

- ~12,000 AMD Instinct MI250X GPUs

- Hardware overview: https://docs.lumi-supercomputer.eu/hardware/

# "Partitions" section

- Partitions: LUMI-C and **LUMI-G**
- 4 GPUs and 8 GCDs per node
  - Node hours: 4 * GPU hours
  - 512 GB host mem + 8 * 64 GB device mem



**2978x** compute nodes

**1x** 64 cores AMD EPYC 7A53

**8x** 64 GB DDR4 memory
**512 GB total**

**4x** AMD MI250x GPU modules per nodes

**2x** Graphics Compute Die (GCD) per module
**8x per compute node**

**4x** 16GB HBM2e stacks
**64 GB per GCD**
**128 GB per module**

**1x** 200 GB/s HPE Cray Slingshot-11 interconnect
**4x per compute node**

# "Code" section

- Main issues: scalability and AMD compatibility

- AMD GPUs: ~~CUDA~~ → ROCm

- ROCm versions of Pytorch, TensorFlow, JAX preinstalled as containers

  - PyTorch pretty well supported out-of-the-box

  - Although we do occasionally stumble on some edge cases…

- Most of common libraries have been ported and optimized for AMD

- CUDA kernels need to be converted

  - E.g. with the HIPIFY tool

# EuroHPC JU access modes and LUMI

- **LUMI is open for all calls**

- Recommended path is to begin with Benchmark or Playground

- Software compatibility and scalability **important** in particular for large access modes

- Humans evaluate the proposals
  - all useful technical information is welcome
  - the reasoning behind the provided values can be very useful

- LUMI AI Factory can help, contact us at lumi-ai-factory.eu



LUMI AI Factory