

# NLP 2025

## CINECA HPC

---

Slides provided by:

- Luca Moroni
- Roberto Navigli

SAPIENZA  
NLP



# Cineca

Non profit Consortium, made up of 119 members:

- 2 Ministries
- 70 Italian Universities
- 48 Italian National Institutions and Agencies

# CINECA

# Cineca: CLUSTER

Cineca is one of the Large Scale Facilities in Europe

- **LEONARDO:** pre-exascale Tier-0 EuroHPC supercomputer. LEONARDO is classified in the 6° position in the Top500 list

# Cineca: CLUSTER

Cineca is one of the Large Scale Facilities in Europe

- **LEONARDO:** pre-exascale Tier-0 EuroHPC supercomputer. LEONARDO is classified in the 6° position in the Top500 list
- **MARCONI:** Tier-0 EuroHPC supercomputer less powerful than LEONARDO
- **GALILEO100:** Tier-1 EuroHPC

# CINECA DOCUMENTATION

The CINECA HPC is well documented, and the **helpdesk** is always active to help you!

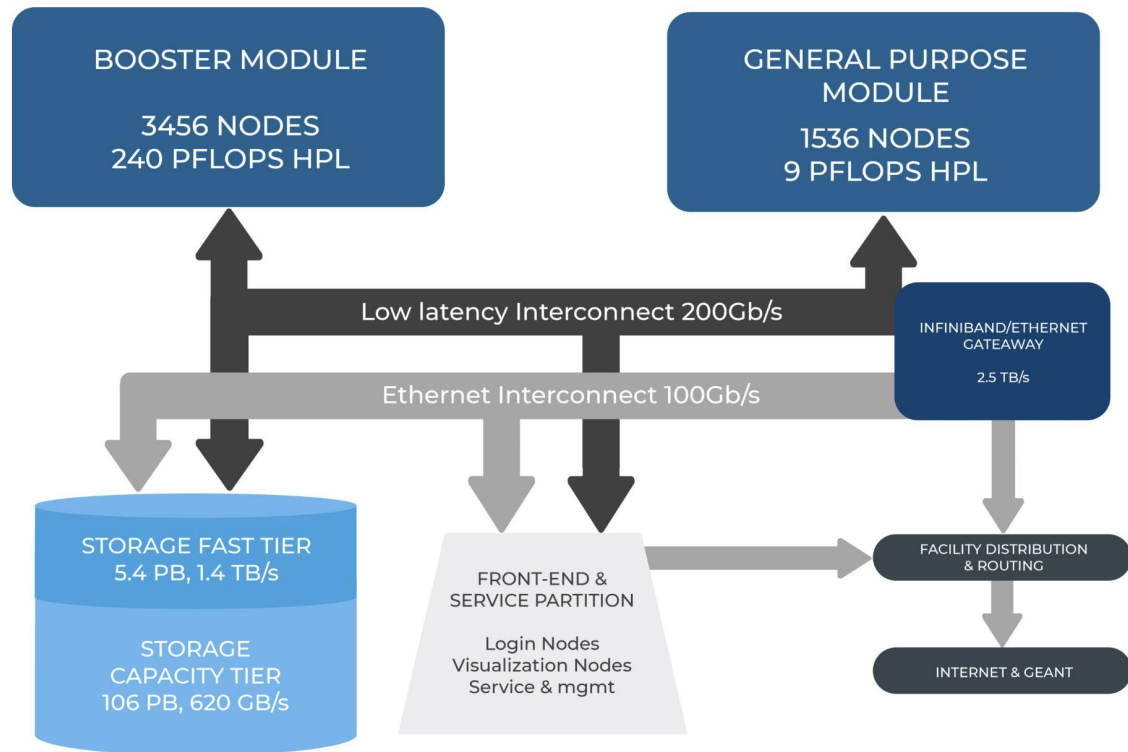
- **Official Documentation:** <https://docs.hpc.cineca.it/index.html>
- **Support E-Mail:** [superc@cineca.it](mailto:superc@cineca.it)

**LEONARDO**



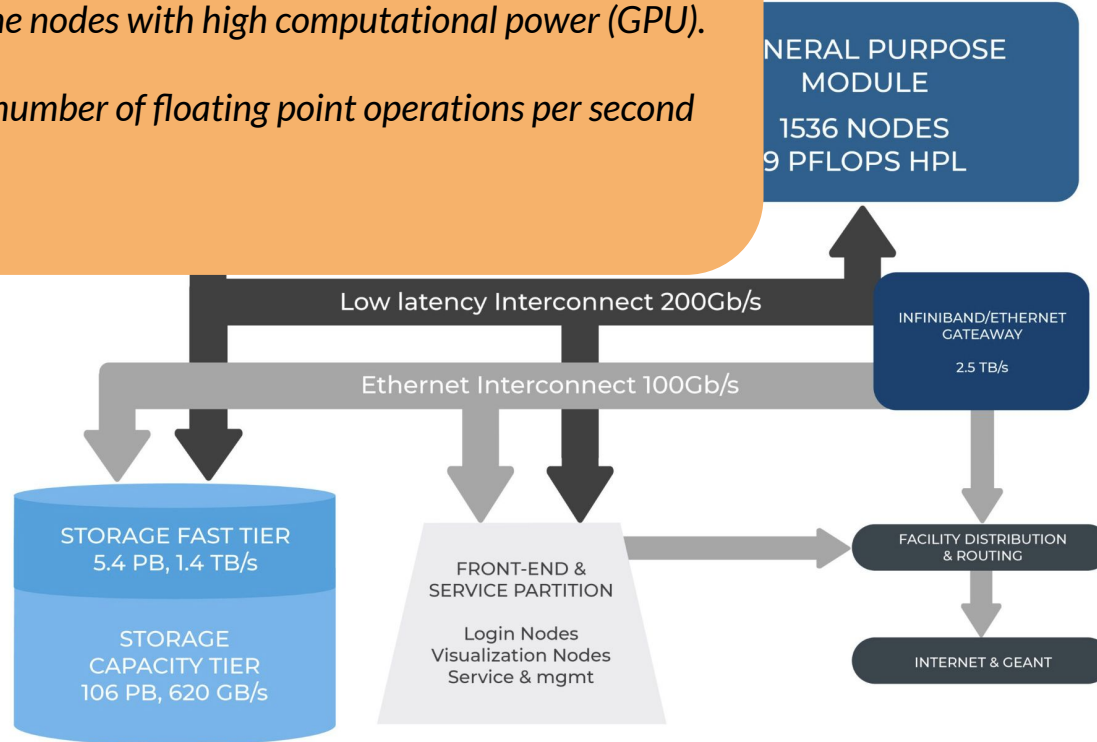
# Leonardo

## System Architecture



**Booster Modules:** the nodes with high computational power (GPU).

**PetaFlops:**  $10^{15}$  number of floating point operations per second





**Booster Modules:** the nodes with high computational power (GPU).

*PetaFlops:  $10^{15}$  number of floating point operations per second*

GENERAL PURPOSE  
MODULE  
1536 NODES  
9 PFLOPS HPL

Low

STORAGE FAST TIER  
5.4 PB, 1.4 TB/s

STORAGE  
CAPACITY TIER  
106 PB, 620 GB/s

LOGIN NODES: HAVE INTERNET CONNECTION

COMPUTE NODES: DO NOT HAVE INTERNET CONNECTION

**Booster Modules:** the nodes with high computational power (GPU).

*PetaFlops:  $10^{15}$  number of floating point operations per second*

You have to **save datasets and models locally** to use them in the compute nodes.  
e.g. Cache-HF

GENERAL PURPOSE  
MODULE  
1536 NODES  
9 PFLOPS HPL

Low

STORAGE FAST TIER  
5.4 PB, 1.4 TB/s

STORAGE  
CAPACITY TIER  
106 PB, 620 GB/s

LOGIN NODES: HAVE INTERNET CONNECTION

COMPUTE NODES: DO NOT HAVE INTERNET CONNECTION

# Leonardo: Compute Nodes

## Booster partition

**Model:** BullSequana X2135 “Da Vinci” single node GPU Blade

**Nodes:** 3456

**Processors:** single socket 32-core Intel Xeon Platinum 8358 CPU, 2.60GHz (Ice Lake)

**Cores:** 110592

**RAM:** 8×64 GB DDR4 3200 MHz (512 GB), per node!

**Accelerators:** 4x NVIDIA custom Ampere A100 GPU 64GB HBM2e, NVLink 3.0 (200GB/s), per node!

**Network:** 2 x dual port HDR100 per node (400Gbps/node)

# First Login

Documentation: <https://docs.hpc.cineca.it/general/access.html>

1. Registration on UserDB: <https://userdb.hpc.cineca.it/>.
  - a. Prof. Navigli should add you on the project: PROJECT\_NAME.
2. Correct configuration of 2FA and OTP: [LINK](#).

Then you can login into one of the 8 Login Nodes:

1. `step ssh login 'your@email.com' --provisioner cineca-hpc`
2. `ssh your-username@login01-ext.leonardo.cineca.it` [note you can use login02, ... login07]

# LEONARDO Budget and Accounting

Documentation:

[https://docs.hpc.cineca.it/hpc/hpc\\_intro.html#budget-and-accounting](https://docs.hpc.cineca.it/hpc/hpc_intro.html#budget-and-accounting)

Command: `saldo -b`

Gives you the usage and the CPU hour availability of your account, for the booster partition. [NOTE: The hours are SHARED]

```
2023, 2023, 2023,  
(.env) [lmoroni0@login01 llm-foundry]$ saldo -b
```

| account     | start    | end      | total<br>(local h) | localCluster<br>Consumed(local h) | totConsumed<br>(local h) | totConsumed<br>% | monthTotal<br>(local h) | monthConsumed<br>(local h) |
|-------------|----------|----------|--------------------|-----------------------------------|--------------------------|------------------|-------------------------|----------------------------|
| IscrB_medit | 20231116 | 20250616 | 3566080            | 3473137                           | 3473137                  | 97.4             | 185090                  | 306414                     |

# LEONARDO Budget and Accounting

Documentation:

<https://docs.hpc.cineca.it/leonardo/budget-and-accounting>

Command: `saldo -`

Gives you the usage partition. [NOTE: The monthTotal, for the booster

IMPORTANT:

**1 GPU HOUR == 8 CPU HOURS**

(.env) [lmoroni0@login01 llm-foundry]

| account     | start    | end      | consumed (local h) | consumed (local h) | consumed (local h) | consumed % | monthTotal (local h) | monthConsumed (local h) |
|-------------|----------|----------|--------------------|--------------------|--------------------|------------|----------------------|-------------------------|
| IscrB_medit | 20231116 | 20250616 | 3566080            | 3473137            | 3473137            | 97.4       | 185090               | 306414                  |

# Leonardo: Data Areas

| Name      | Area Attributes         | Quota                 | Backup | Note  |
|-----------|-------------------------|-----------------------|--------|---|
| \$HOME    | permanent,user specific | 50 GB                 | daily  |   |
| \$WORK    | permanent,shared        | 1 TB                  | no     | Large data to be shared with project's colla                                |
| \$FAST    | permanent, shared       | 1 TB                  | no     | Only on <a href="#">Leonardo</a> .<br>Faster I/O compared with outer areas. |
| \$SCRATCH | temporary,user specific | -/20 TB               | no     | files older than 40 days<br>are deleted                                     |
| \$TMPDIR  | temporary,user specific | (-)                   | no     | directory removed<br>at job completion                                      |
| \$PUBLIC  | permanent,open          | 50 GB                 | no     | Only on <a href="#">Leonardo</a> .  |
| \$DRES    | permanent,shared        | defined<br>by project | no     |   |

# Leonardo: Data Areas

| Name      | Area Attributes  |  |    |  |
|-----------|------------------|--|----|--|
| \$HOME    | permanent,use    | <div><b>HOME:</b> your home directory, where you will write your code, download the github repos and do lightweight stuff.</div> |    |  |
| \$WORK    | permanent,sha    |  |    |  |
| \$FAST    | permanent, sha   |  |    |  |
| \$SCRATCH | temporary,use    |  |    |  |
| \$TMPDIR  | temporary,use    |  |    |  |
| \$PUBLIC  | permanent,open   |  |    |  |
| \$DRES    | permanent,shared | defined<br>by project  | no |  |



# Leonardo: Data Areas

| Name      | Area Attributes  |                       |    |                      |
|-----------|------------------|-----------------------|----|----------------------|
| \$HOME    | permanent,use    |                       |    |                      |
| \$WORK    | permanent,sha    |                       |    | with project's colla |
| \$FAST    | permanent, sha   |                       |    | outer areas.         |
| \$SCRATCH | temporary,use    |                       |    |                      |
| \$TMPDIR  | temporary,use    |                       |    |                      |
| \$PUBLIC  | permanent,open   |                       |    |                      |
| \$DRES    | permanent,shared | defined<br>by project | no |                      |

**SCRATCH:** your high capable directory, non-permanent, but with a lot of space. save here the hf\_cache, models, datasets, and heavy stuffs.

Remember that the data here will be removed after 40 days.

# Environment



# Environment

Once you are logged-in, you will land in your \$HOME directory

Leonardo works with modules

Modules are containerized ready-to-use packages and tools.

[https://docs.hpc.cineca.it/hpc/hpc\\_enviroment.html](https://docs.hpc.cineca.it/hpc/hpc_enviroment.html)

| Command                     | Action  |
|-----------------------------|---|
| module avail                | show the available modules on the machine   |
| module load <appl>          | load the module in the current shell session, preparing the enviroment for the application. |
| module load autoload <appl> | load the module and all dependencies in the current session                                 |
| module help <appl>          | show specific information and basic help on the application                                 |
| module list                 | show the module currently loaded in the shell session                                       |
| module purge                | unload all the loaded modules   |
| module unload <appl>        | unload a specific module  |

# Environment: install python libraries

To install and use a Python library you have to load the following modules:

- module load profile/deeplrn
- module load cuda/12.1 # available even 12.3

Once you have loaded the modules you can create a virtual environment and install the libraries:

```
$ python -m venv .env          # create the environment
$ source .env/bin/activate     # activate the environment
$ pip install transformers     # install library
```

# Execution

Once you have created the environment (usually in your \$HOME or whatever) you can execute your Python script

You can execute the scripts with the the SLURM system

SLURM -> allows you to execute programs on the compute nodes

[https://docs.hpc.cineca.it/hpc/hpc\\_scheduler.html](https://docs.hpc.cineca.it/hpc/hpc_scheduler.html)

# SLURM Scheduling

There are different SLURM directives:

- `sbatch <slurm_script.sh>`: submit the job to the scheduler
- `squeue --me`: monitoring your submitted jobs
- `scancel job_id`: cancel the job

# SLURM Scripts

You can create a .sh script with the following **SBATCH** instructions:

```
#!/bin/bash

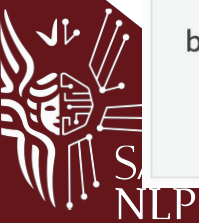
#SBATCH --nodes=1           # 1 node
#SBATCH --ntasks-per-node=32 # 32 tasks per node
#SBATCH --time=1:00:00      # time limit: 1 hour
#SBATCH --error=myJob.err   # standard error file
#SBATCH --output=myJob.out  # standard output file
#SBATCH --account=<Project Account> # project account
#SBATCH --partition=<partition_name> # partition name
#SBATCH --qos=<qos_name>        # quality of service

./my_application
```



# SLURM partitions

| Partition                   | QOS                | #Cores/#GPU per job               | Walltime   | Max Nodes/cores/GPUs/...                            |
|-----------------------------|--------------------|-----------------------------------|------------|---|
| lrd_all_serial<br>(default) | normal             | 4 cores<br>(8 logical cores)      | 04:00:00   | 1 node / 4 cores<br>(30800 MB RAM)                  |
| boost_usr_prod              | normal             | 64 nodes                          | 24:00:00   |   |
|                             | boost_qos_dbg      | 2 nodes                           | 00:30:00   | 2 nodes / 64 cores / 8 GP                           |
|                             | boost_qos_bprod    | min = 65 nodes<br>max = 256 nodes | 24:00:00   | 256 nodes   |
|                             | boost_qos_lprod    | 3 nodes                           | 4-00:00:00 | 3 nodes / 12 GPUs                                   |
| boost_fua_dbg               | normal             | 2 nodes                           | 00:10:00   | 2 nodes / 64 cores / 8 GP                           |
| boost_fua_prod              | normal             | 16 nodes                          | 24:00:00   | 4 running jobs per user ac<br>32 nodes / 3584 cores |
|                             | boost_qos_fuabprod | min = 17 nodes<br>max = 32 nodes  | 24:00:00   | 32 nodes / 3584 cores                               |
|                             | qos_fualowprio     | 16 nodes                          | 08:00:00   |   |





# LLM generation



# Example: LLM generation

A simple step-by-step guide:

1. **Create your environment**
2. **Load the environment and install the Transformers library**
3. **Create the Python script with your code**
4. **Create the slurm script with all the correct SBATCH directives**

# Python Script

my\_fancy\_generation.py

```
my_fancy_generation.py > ...
1  import transformers
2  import torch # type: ignore
3
4  model_id = "meta-llama/Meta-Llama-3.1-8B-Instruct" # <- YOUR FANCY MODEL HERE!
5
6  pipeline = transformers.pipeline(
7      "text-generation",
8      model=model_id,
9      model_kwargs={"torch_dtype": torch.bfloat16},
10     device_map="auto",
11 )
12
13 messages = [
14     {"role": "system", "content": "You are a pirate chatbot who always responds in pirate speak!"},
15     {"role": "user", "content": "Who are you?"},
16 ]
17
18 outputs = pipeline(
19     messages,
20     max_new_tokens=256,
21 )
22 print(outputs[0]["generated_text"][-1])
```

# SLURM Script

my\_fancy\_generation.slurm.sh

```
$ my_fancy_generation.slurm.sh
1  #!/bin/bash
2  #SBATCH --job-name=my_fancy_gen          # Job name
3  #SBATCH --output=/path/to/your/log/folder/log_name-%x-%j.out  # Name of stdout output file. %x job_name, %j job_number
4  #SBATCH --error=/path/to/your/log/folder/log_name-%x-%j.err   # Name of stdout output file. %x job_name, %j job_number
5
6  #SBATCH -A <Account_Name>                # account name
7  #SBATCH -p boost_usr_prod                 # partition usually production
8  #SBATCH -q boost_qos_bprod               # quality of service
9  #SBATCH --time 24:00:00                   # timing: HH:MM:SS for booster maximum 24H
10 #SBATCH -N 1                             # number of nodes, one should be enough
11 #SBATCH --ntasks=1                       # number of tasks
12 #SBATCH --ntasks-per-node=1              # number of tasks per node
13 #SBATCH --cpus-per-task=1                # number of cpu per tasks
14 #SBATCH --gres=gpu:1                     # number of GPU per node, 1 should be enough
15
16
17 # load the modules
18 module load profile/deeplrn cuda/12.1
19
20 # activate your environment
21 source /path/to/your/environment/folder/bin/activate
22
23 # load and set the Huggingface CACHE, to retrieve cached informations in a no-internet environment
24
25 export HF_DATASETS_CACHE=/path/to/hf/cache
26 export HUGGINGFACE_HUB_CACHE=/path/to/hf/cache
27 export WANDB_MODE=offline # set the wandb offline (no needed for generation...)
28 # read Huggingface token from .env file
29 export HF_TOKEN=$(python -c "import huggingface_hub; print(huggingface_hub.HfFolder.get_token() or '')")
30
31 # execute the python script
32 python /path/to/your/python/script/my_fancy_generation.py
```

# SLURM Script

my\_fancy\_generation.slurm.sh

```
$ my_fancy_generation.slurm.sh
1  #!/bin/bash
2  #SBATCH --job-name=my_fancy_gen          # Job name
3  #SBATCH --output=/path/to/your/log/folder/log_name-%x-%j.out    # Name of stdout output file. %x job_name, %j job_number
4  #SBATCH --error=/path/to/your/log/folder/log_name-%x-%j.err     # Name of stdout output file. %x job_name, %j job_number
5
6  #SBATCH -A <Account_Name>                # account name
7  #SBATCH -p boost_usr_prod                 # partition usually production
8  #SBATCH -q boost_qos_bprod               # queue
9  #SBATCH --time 24:00:00                  # time for booster maximum 24H
10 #SBATCH -N 1                             # number of nodes, one should be enough
11 #SBATCH --ntasks=1                       # number of tasks per node
12 #SBATCH --ntasks-per-node=1              # tasks
13 #SBATCH --cpus-per-task=1                # cpus per node, 1 should be enough
14 #SBATCH --gres=gpu:1
15
16
17 # load the modules
18 module load profile/deeplrn cuda
19
20 # activate your environment
21 source /path/to/your/environment/folder/bin/activate
22
23 # load and set the Huggingface CACHE, to retrieve cached informations in a no-internet environment
24
25 export HF_DATASETS_CACHE=/path/to/hf/cache
26 export HUGGINGFACE_HUB_CACHE=/path/to/hf/cache
27 export WANDB_MODE=offline # set the wandb offline (no needed for generation...)
28 # read Huggingface token from .env file
29 export HF_TOKEN=$(python -c "import huggingface_hub; print(huggingface_hub.HfFolder.get_token() or '')")
30
31 # execute the python script
32 python /path/to/your/python/script/my_fancy_generation.py
```

Log folder should exists  
otherwise the process  
doesn't work

# SLURM Script

my\_fancy\_generation.slurm.sh

```
$ my_fancy_generation.slurm.sh
1  #!/bin/bash
2  #SBATCH --job-name=my_fancy_gen          # Job name
3  #SBATCH --output=/path/to/your/log/folder/log_name-%x-%j.out  # Name of stdout output file. %x job_name, %j job_number
4  #SBATCH --error=/path/to/your/log/folder/log_name-%x-%j.err   # Name of stdout output file. %x job_name, %j job_number
5
6  #SBATCH -A <Account_Name>                # account name
7  #SBATCH -p boost_usr_prod                 # partition usually production
8  #SBATCH -q boost_qos_bprod                # quality of service
9  #SBATCH --time 24:00:00                   # for booster maximum 24H
10 #SBATCH -N 1                             # one should be enough
11 #SBATCH --ntasks=1
12 #SBATCH --ntasks-per-node=1               # per node
13 #SBATCH --cpus-per-task=1                 # tasks
14 #SBATCH --gres=gpu:1                      # node, 1 should be enough
15
16
17 # load the modules
18 module load profile/deeplrn cuda
19
20 # activate your environment
21 source /path/to/your/environment/folder/bin/activate
22
23 # load and set the Huggingface CACHE, to retrieve cached informations in a no-internet environment
24
25 export HF_DATASETS_CACHE=/path/to/hf/cache
26 export HUGGINGFACE_HUB_CACHE=/path/to/hf/cache
27 export WANDB_MODE=offline # set the wandb offline (no needed for generation...)
28 # read Huggingface token from .env file
29 export HF_TOKEN=$(python -c "import huggingface_hub; print(huggingface_hub.HfFolder.get_token() or '')")
30
31 # execute the python script
32 python /path/to/your/python/script/my_fancy_generation.py
```

Look at your output file to  
see the generated text!

# Llama Factory



SAPIENZA  
NLP



# Llama Factory

Llama Factory is a ready-to-use library for LLM training

You could use it to instruction-tune a LLM with new conversations

<https://github.com/hiyouga/LLaMA-Factory>

To install:

```
$ git clone --depth 1 https://github.com/hiyouga/LLaMA-Factory.git
$ cd LLaMA-Factory
$ source /path/to/env/bin/activate
$ pip install -e "[torch,metrics]" --no-build-isolation
```



# Llama Factory

After installed, run:

```
$ llamafactory-cli train /path/to/config.yaml
```

Inside a slurm script.

# Llama Factory: Config

```
examples > train_full > ! sft_example_blank.yaml
```

```
1  ### model
2  model_name_or_path: path_of_model # no hf_id but the absolute path, save it in a specific directory
3  new_special_tokens: <|start_header_id|>,<|end_header_id|>,<|eot_id|> # token added with the chat template
4
5  ### method
6  stage: sft # Supervised Fine-Tuning
7  do_train: true
8  finetuning_type: full
9  use_badam: true # use BADAM not classical full-finetuning, faster!
10 badam_mode: layer
11 badam_switch_mode: ascending
12 badam_switch_interval: 50
13 badam_verbose: 2
14 flash_attn: fa2
15 deepspeed: examples/deepspeed/ds_z3_config.json # deepspeed
16
17 ### dataset
18 dataset: magpielm-v0.1 # comma separated datasets, defined in: data/dataset_info.json
19 template: llama3 # chat_template
20 cutoff_len: 2048 #
21 max_samples: 560000
22 overwrite_cache: true
23 preprocessing_num_workers: 16
24
25 ### output
26 output_dir: /path/to/output # output directory
27 logging_steps: 10
28 save_steps: 500
29 plot_loss: true
30 overwrite_output_dir: true
```

```
32 ### train
33 per_device_train_batch_size: 4
34 gradient_accumulation_steps: 8
35 learning_rate: 1.0e-5
36 num_train_epochs: 1.0
37 lr_scheduler_type: cosine
38 warmup_ratio: 0.05
39 bf16: true
40 ddp_timeout: 180000000
41
42 ### eval
43 val_size: 0.01
44 per_device_eval_batch_size: 4
45 eval_strategy: steps
46 eval_steps: 100
47
48 ### logging
49 report_to: wandb
50 run_name: run_name # the name of the run
```

# Llama Factory: Config

```
examples > train_full > ! sft_example_blank.yaml
```

```
1  ### model
2  model_name_or_path: path_of_model # no hf_id but the absolute path, save it in a specific directory
3  new_special_tokens: <|start_header_id|>,<|end_header_id|>,<|eot_id|> # token added with the chat template
4
5  ### method
6  stage: sft # Supervised Fine-Tuning
7  do_train: true
8  finetuning_type: full
9  use_badam: true # use BADAM not classical
10 badam_mode: layer
11 badam_switch_mode: ascending
12 badam_switch_interval: 50
13 badam_verbose: 2
14 flash_attn: fa2
15 deepspeed: examples/deepspeed/ds_z3_config
16
17 ### dataset
18 dataset: magpielm-v0.1 # comma separated
19 template: llama3 # chat_template
20 cutoff_len: 2048 #
21 max_samples: 560000
22 overwrite_cache: true
23 preprocessing_num_workers: 16
24
25 ### output
26 output_dir: /path/to/output # output directory
27 logging_steps: 10
28 save_steps: 500
29 plot_loss: true
30 overwrite_output_dir: true
```

The *chat\_template* is the text structure used by the chat models.

Starting from a *base models* we should add to the vocabulary the *special tokens* used to structure the *chat\_template*.

```
32  ### train
33  per_device_train_batch_size: 4
34  gradient_accumulation_steps: 8
35  learning_rate: 1.0e-5
36  num_train_epochs: 1.0
37  lr_scheduler_type: cosine
38  warmup_ratio: 0.05
39  bf16: true
40  timeout: 180000000
41
42  eval
43  eval_size: 0.01
44  device_eval_batch_size: 4
45  eval_strategy: steps
46  eval_steps: 100
47
48  logging
49  report_to: wandb
50  run_name: run_name # the name of the run
```

# Add your custom dataset

To add a conversation dataset, you can load it from HF, and add it to the  
data/dataset\_info.json LlamaFactory file

```
"conversations-it": {  
  "hf_hub_url": "sapienzanlp/conversations_italian",  
  "formatting": "sharegpt",  
  "columns": {  
    "messages": "messages"  
  },  
  "tags": {  
    "role_tag": "role",  
    "content_tag": "content",  
    "user_tag": "user",  
    "assistant_tag": "assistant",  
    "system_tag": "system"  
  }  
},
```

# Create your custom dataset

Once you have collected your beautiful and useful conversations you can upload them to your HF profile

Make sure you follow a standardized format. We suggest you to use the **SHAREGPT** format.

# Create your custom dataset

Once you have collected your beautiful and useful conversations you can upload them to your HF profile

Make sure you follow a standardized format. We suggest you to use the SHAREGPT format.

Each conversation is encoded in a list field named “messages”.

The “messages” field contains a list of dictionaries with user and assistant (LLM) messages.

# Create your custom dataset

Once you have collected your beautiful and useful conversations, you can export them to your HF profile

Make sure you follow a standardized format. We suggest the following format.

Each conversation is encoded in a list field named “messages”

The “messages” field contains a list of dictionaries with the following structure:

```
messages
list
```

```
[
{
  "content": "These instructions apply to section-based themes (Responsive 6.0+, Retina 4.0+, Parallax 3.0+ Turbo 2.0+, Mobilia 5.0+). What theme version am I using?\n\nOn your Collections pages & Featured Collections sections, you can easily show the secondary image of a product on hover by enabling one of the theme's built-in settings!\n\nYour Collection pages & Featured Collections sections will now display the secondary product image just by hovering over that product image thumbnail.\n\nDoes this feature apply to all sections of the theme or just specific ones as listed in the text material?",
  "role": "user"
},
{
  "content": "This feature only applies to Collection pages and Featured Collections sections of the section-based themes listed in the text material.",
  "role": "assistant"
}
]
```

# wandb

**wandb** is an useful tool to monitor your ML training

Some of you already used it in the first homework

you have to install it in your environment: `pip install wandb`

Create an account on [wandb.ai](https://wandb.ai) then login locally by `wandb login`

**wandb will be used offline by the compute nodes (see the slurm script) you have to sync the run online after the execution:** `wandb sync path_to_run`