# Introduction to Supercomputing

## How to Use an HPC Machine

Chris Stylianou

Research Engineer

CaSToRC

eurocc.cyi.ac.cy

- Research Engineer at CaSToRC, PhD

- Collaborations Task Leader for EuroCC2

- Area of expertise in High Performance Computing (HPC)

- Contact & Info:
  - Email: c.stylianou@cyi.ac.cy
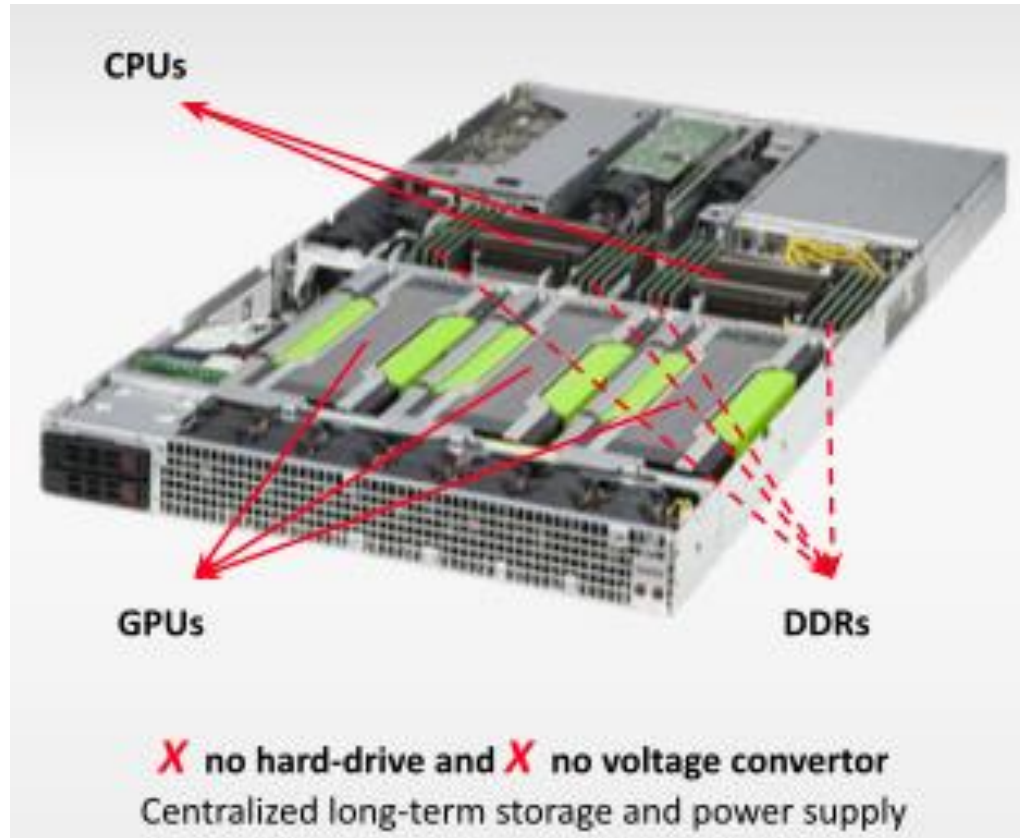  - Website: cstyl.github.io

# Compute Node    vs    Desktop



A node is
~5x faster
than a desktop

**Compute Node**



**Rack**



**Supercomputer**

Shared across **all** users!
Not for executing
heavy/long workloads!

Terminal
> cmd

UNIX shell on
your computer

The HPC has 3 parts:
Head node
Worker nodes
File system

data

Head node

submit
jobs

For executing
workloads.
Resources **scheduled**
explicitly per user.

file system

data

Worker nodes

## Linux/Mac:

Open a new Terminal Session.

In `~/.ssh` directory create a `config` file:

```
$ cd ~/.ssh
$ vim config
```

Configure file as shown in image. **Note** you will need to change the `User` field to your own and `IdentityFile` to the path of your ssh key generated for the hackathon.

## Windows:

Open a new PowerShell.

In `<path\to\home>\.ssh` directory create a `config` file.

**Remove** any extention (e.g., .txt) from the `config` file!

Open file and configure file as shown in image.

**Note** you will need to change the `User` field to your own and `IdentityFile` to the path of your ssh key generated for the hackathon.

```
≡ config
1    Host *
2        AddKeysToAgent yes
3
4    Host cyclone
5        Hostname cyclone.hpcf.cyi.ac.cy
6        User hack24cs1
7        IdentityFile ~/.ssh/id_rsa
8        ForwardAgent yes
```

# Linux/Mac/Windows:

Open a Terminal Session.

```
$ ssh cyclone
```

**If that didn't work**, possibly `config` wasn't picked up. Try:

```
$ ssh-add <path/to/ssh/key>
$ ssh <Username>@cyclone.hpcf.cyi.ac.cy
```

# *Windows*:

If you have created your key using puTTY you will need to save it using the **OpenSSH format**.

More on https://hpcf.cyi.ac.cy/documentation/login.html

**Definition:**

- Simple Linux Utility for Resource Management (SLURM) is an open-source **job scheduler** used in Linux clusters **to manage and allocate resources**.

**Key Features:**

- Job scheduling and management.
- Resource allocation.
- Monitoring and accounting.

**Basic Workflow:**

**1. Submit a Job:** Use `sbatch` to submit a job script.

**2. Monitor a Job:** Use `squeue` to view job status.

**3. Job Completion:** Job finishes and results are available.

| Category | Command | Action |
|---|---|---|
| Job Submission | sbatch <job_script> | Submit a batch job |
| | srun <command> | Run a command in parallel |
| Job Monitoring | squeue | View queued jobs |
| | squeue --user=$USER | View MY queued jobs |
| | scontrol show job <job_id> | Detailed job info |
| Job Management | scancel <job_id> | Cancel a job |
| | scontrol hold <job_id> | Hold a job |
| | scontrol release <job_id> | Release a held job |
| Resource Allocation | sinfo | Information about nodes and partitions |
| | scontrol show node <node_id> | Show Node details |

```bash
#!/bin/bash -l

#SBATCH --job-name=hackathon
#SBATCH --partition=gpu # Partition
#SBATCH --nodes=1 # Number of nodes
#SBATCH --gres=gpu:1 # Number of GPUs
#SBATCH --ntasks-per-node=1  # Number of tasks
#SBATCH --cpus-per-task=10
#SBATCH --output=~/slurm-output/job.%j.out # Stdout (%j=jobId)
#SBATCH --error=~/slurm-output/job.%j.err # Stderr (%j=jobId)
#SBATCH --time=24:00:00 # Walltime
#SBATCH -A hackathon # Accounting project
#SBATCH --reservation=hackathon

# Load any necessary modules
module load CUDA
nvidia-modprobe -c 0 -u

# Launch the executable exe.out
srun apptainer run --writable-tmpfs --nv \
                          -B /nvme/scratch/hackathon/ollama_models:/nvme/scratch/hackathon/ollama_models \
                          /nvme/scratch/hackathon/hackathon.sif
```

**One** GPU per team!

Reservation to get jobs through faster!

**Only one user:** Launch the job from Cyclone:

```
$ cd ~/
$ sbatch hackathon/launch_hackathon.sub
sbatch: Jobs under the 'hackathon' reservation will have the GPU
option set to 1.
Submitted batch job 968110
```

Check job status. After a while the status (ST) should be set to running (R):

```
$ squeue --user=$USER
JOBID   PARTITION   NAME    USER    ST  TIME  NODES NODELIST(...)
968110       gpu   hackatho hack24cs R  0:00    1       gpu03
```

Once the job starts running, a `connection_info.txt` should appear in your `$HOME` directory containing info to:

- Setup SSH tunneling

- URL for accesing notebook on your browser
  - Password

- Model endpoint.

```
[hack24cs1@front02 ~]$ sbatch hackathon/launch_hackathon.sub
sbatch: Jobs under the 'hackathon' reservation will have the GPU option set to 1.
Submitted batch job 968110
[hack24cs1@front02 ~]$ squeue --user=$USER
            JOBID PARTITION     NAME     USER ST      TIME  NODES NODELIST(REASON)
            968110       gpu hackatho hack24cs  R      0:02      1 gpu03
[hack24cs1@front02 ~]$ ls
hackathon   job.968110.err  job.968110.out  jupyter.log  ollama.log
[hack24cs1@front02 ~]$ ls
connection_info.txt  job.968110.err  jupyter.log
hackathon                  job.968110.out  ollama.log
[hack24cs1@front02 ~]$ cat connection_info.txt
    ==================================================================
    Run this command to connect on your jupyter notebooks remotely
    ssh -N -J hack24cs1@cyclone.hpcf.cyi.ac.cy hack24cs1@gpu03 -L 23239:localhost:23239

    Jupyter Notebook is running at: http://localhost:23239
    Password to access the notebook: YuZpYX53cpLWW4ME
    Ollama is serving models on: http://gpu03:23229
    ==================================================================
```

SSH Tunneling

URL & Password

Model Endpoint

- Launch a new terminal session

$ ssh -N -J <Username>@cyclone.hpcf.cyi.ac.cy <username>@gpu03 -L 23239:localhost:23239

- If prompted about establishing authenticity of the host, press "Y" and enter.

```
[cstyl@cstyl-2:~$ ssh -N -J hack24cs1@cyclone.hpcf.cyi.ac.cy hack24cs1@gpu03 -L 23239:localhost:23239
[The authenticity of host 'gpu03 (<no hostip for proxy command>)' can't be established.
ED25519 key fingerprint is SHA256:sLgZ1EDZEwruWtCE5dZ2T0y2Q2o6HiKiE5l5jk0s2Yg.
This host key is known by the following other names/addresses:
    ~/.ssh/known_hosts:7: front02.hpcf.cyi.ac.cy
    ~/.ssh/known_hosts:34: cyclone.hpcf.cyi.ac.cy
    ~/.ssh/known_hosts:35: sim01
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'gpu03' (ED25519) to the list of known hosts.
```

- After the SSH Tunnel has been established, open your browser and enter the link and password to Jupyter Notebook. E.g.,:
  - `http://localhost:23239`
  - `YuZpYX53cpLWW4ME`

- **Open** `hackathon > example.ipynb`

- Run first three cells to prompt your local LLM!
  - First time it runs might take a while until it loads the model.

- Share with each team member the URL and Password to open the model on their own browser!



jupyter  example

File  Edit  View  Run  Kernel  Settings  Help                                                          Not Trusted

Markdown ∨                          JupyterLab ⬈  ⚙  Python 3 (ipykernel) ○ ☰  notebook is read-only

### Configure Ollama for inference

```
[3]: import os

     OLLAMA_ENDPOINT="http://gpu03:23229"
     os.environ['OLLAMA_HOST']=OLLAMA_ENDPOINT
```

```
[4]: from langchain_community.chat_models import ChatOllama

     local_llm = "gemma2:27b"
     num_ctx=25000
     llm = ChatOllama(model=local_llm, base_url=OLLAMA_ENDPOINT, temperature=0, num_ctx=num_ctx)
```

### Prompt your local LLM

```
[5]: from langchain_core.messages import AIMessage

     messages = [
         (
             "system",
             "You are a helpful assistant that translates English to French. Translate the user sentence. Don't respond in json",
         ),
         ("human", "I love programming. Tell me if you also love programming"),
     ]
     ai_msg = llm.invoke(messages)
     ai_msg.pretty_print()
```

```
================================ Ai Message ================================

J'adore programmer. Oui, j'aime beaucoup programmer ! C'est très gratifiant de pouvoir créer des choses avec du code.

What do you like to program?
```

More information:

https://castorc.cyi.ac.cy/
https://eurocc.cyi.ac.cy/

Contact us at:

eurocc-contact@cyi.ac.cy