

hcc.unl.edu

Hanying Chen, Yi Liu, Caughlin Bohn

Showmic Islam, Natasha Pavlovikj, Caughlin Bohn

Introduction to AI and ML using HCC

April 2025

Schedule

timestamp	Session			
12:00 ~ 12:15pm	Setup and Support			
12:15 ~ 1:00pm	Introduction to HCC resources for ML and AI (45 Min)			
1:00 ~ 1:15pm	Short Break (15 Min)			
1:15 ~ 2:00pm	ML and AI workflows (45 Min)			
2:00 ~ 2:10pm	Short Break (10 Min)			
2:10 ~ 2:50pm	Introduction to PyTorch (40 Min)			
2:50 ~ 3:00pm	Break (10 Min)			
3:00 ~ 3:50pm	Introduction to PyTorch cont. (50 Min)			
3:50 ~ 4:00pm	Break (10 Min)			
4:00 ~ 4:15pm	Introduction to National Research Platform (NRP)			
4:15 ~ 4:30pm	Open Questions			

Welcome!

- Available tools and software on HCC
- Requesting resources
- Monitoring model performance

How can the Holland Computing Center help you?









Takeaways

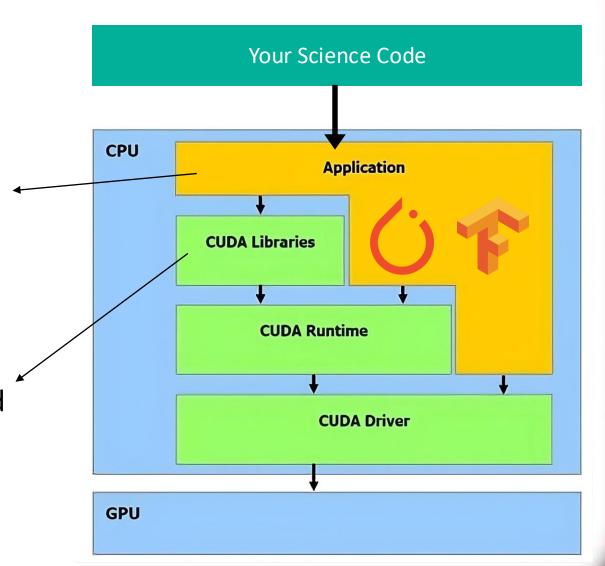
In this session, we will introduce the resources available on HCC's Swan. You will learn:

- The available software on Swan for machine learning and AI
- Options to develop and run experiments using Swan
- Options to monitor and manage your experiments on Swan

Available Software

GPU acceleration ready framework:
PyTorch, TensorFlow

Custom implementation: dedicated but complicated



Available Software

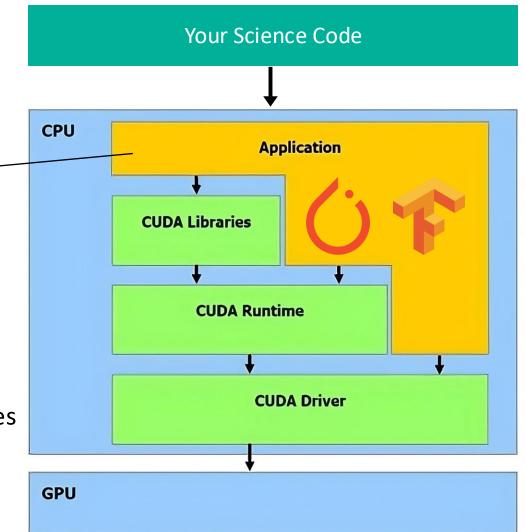
GPU acceleration ready framework:
PyTorch, TensorFlow

PyTorch:

Fast prototyping, Ease of developing

TensorFlow:

Production-ready deployment features



Overview

- **HCC** also offers
 - Pure PyTorch/TensorFlow framework via Module
 - Jupyter Lab kernels for development
 - Options to create a <u>custom</u> Anaconda environment or **Docker** image to support your research
 - LM Studio and Ollama for running local Large Language Models (LLMs)













Pre-Installed Software



Multiple PyTorch versions are installed on Swan, including both CPU and GPU versions

```
To find other possible module matches execute:

$ module -r spider '.*pytorch.*'
```

User-friendly web-based module search:

https://hcc.unl.edu/docs/applications/modules/available_software_for_swan/

Name	Version	Module Name	Prerequisites	Туре	Domain	Description
pytorch						
pytorch-gpu	2.5.1	pytorch-gpu/py310/2.5	None	application	deep learning	Tensors and Dynamic neural networks in Python with strong GPU acceleration
pvtorch-apu	2.5.1	pvtorch-apu/pv311/2.5	None	application	deep learning	Tensors and Dynamic neural networks in

Pre-Installed Software



Similarly, multiple CPU/GPU TensorFlow versions are installed on Swan

```
To find other possible module matches execute:

$ module -r spider '.*tensorflow.*'
```

- Modules provide pure PyTorch/TensorFlow framework
- HCC offers
 - Create a custom Anaconda environment
 - https://hcc.unl.edu/docs/applications/user_software/using_ana conda_package_manager/#creating-custom-gpu-anacondaenvironment

Pre-Installed Software

Command line access:

Load PyTorch/TensorFlow module

For example, to load the latest PyTorch GPU version:

\$ module load pytorch-gpu/py312/2.5

To load the latest TensorFlow GPU version:

\$ module load tensorflow-gpu/py311/2.17

```
[yliu95@login1.swan ~]$ module -r spider '.*pytorch.*'

pytorch:

Description:
   PyTorch is an optimized tensor library for deep learning using GPUs and CPUs.

Versions:
   pytorch/py27/0.4
   pytorch/py27/1.2
   pytorch/py35/0.4
   pytorch/py36/0.4
   pytorch/py36/1.2
```

Running in Command Line

Submitting an interactive job is useful to:

- Test your command before submitting SLURM job
- Debug your experiment

For example:

```
$ srun --nodes=1 --ntasks-per-node=4 --mem-per-cpu=1024 --time=3:00:00 --pty $SHELL
```

To request GPU resources

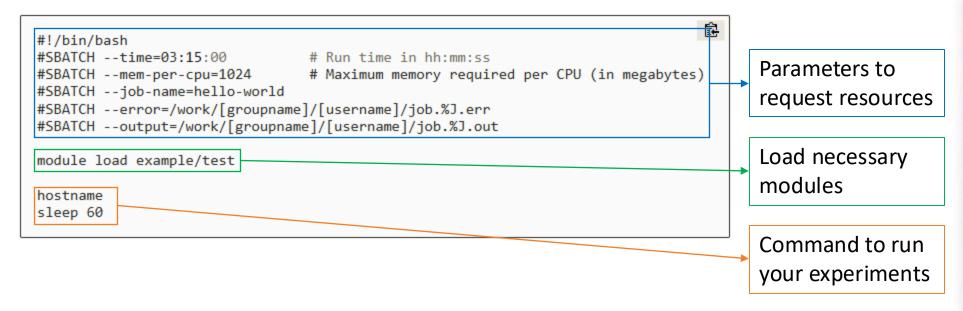
```
$ srun --nodes=1 --ntasks-per-node=4 --mem-per-cpu=1024 --time=3:00:00 --partition=gpu --gres=gpu --pty $SHELL
```

More details see:

https://hcc.unl.edu/docs/submitting jobs/creating an interactive job/

Running with SLURM

For example:



More resource requesting command see:

https://hcc.unl.edu/docs/submitting_jobs/

Running with SLURM

To request GPU resources, two additional parameters are required:

```
--partition=gpu --gres=gpu
```

For example:

```
#!/bin/bash
#SBATCH --time=03:15:00
#SBATCH --mem-per-cpu=1024
#SBATCH --job-name=cuda
#SBATCH --partition=gpu
#SBATCH --gres=gpu
#SBATCH --error=/work/[groupname]/[username]/job.%J.err
#SBATCH --output=/work/[groupname]/[username]/job.%J.out
module load cuda
./cuda-app.exe
```

More resource requesting command see:

https://hcc.unl.edu/docs/submitting_jobs/submitting_gpu_jobs/

Running with Jupyter Lab

Interactive coding environments on Open OnDemand:

- Run code part-by-part and visualize results instantly
- Develop your experiments

Open OnDemand access (https://swan-ood.unl.edu/): Interactive Apps → Jupyter Lab

Select preferred Jupyter Lab (Ver. 3.4 or 4.0)

Jupyter Lab version:	
vv1.0.0_16_g3837770	
This app will launch a Jupyter Lab server using Python.	
Jupyter Lab version	
4.0	~
This defines the version of Jupyter Lab you want to load.	
Working Directory	

Command Line vs. Open OnDemand

Command line:

- Allow extended running time using SLURM (7 days maximum)
- Better for running experiments
- Submit your job to Swan using SLURM https://hcc.unl.edu/docs/submitting_jobs/

Jupyter Lab in Open OnDemand:

- GUI IDE
- Limited running time (8 hours maximum)
- Better for developing code
- Create interactive app using Open OnDemand
 https://hcc.unl.edu/docs/open_ondemand/virtual_desktop_and_interactive_apps/

Custom Environment

ML/AI is data-driven

- Additional package often required data processing and post-analysis
- For example:
 - OpenCV image processing
 - NLTK natural language processing
 - SciPy computational analysis

HCC provides support for custom environments using:

- Anaconda easier environment build
- Docker more consistency across machines

Custom Environment

Anaconda:

- Start by cloning one of the pre-installed kernel
- Tutorial:

https://hcc.unl.edu/docs/applications/user_software/using_anaconda_package_manager/#creating-custom-gpu-anaconda-environment

Docker:

- Similarly, starting from a GPU-enabled Docker image on HCC's Docker hub (https://hub.docker.com/u/unlhcc/)
- Tutorial:

https://hcc.unl.edu/docs/applications/user_software/using_apptainer/

Large Language Models

Large language models (LLMs) are large models that are pre-trained on vast amounts of data. LLMs can be used on different ways on HCC resources.

CLI Tool: Ollama and Ollama-GPU

Name	Version	Module Name
= ↑		
ollama		
ollama-gpu	0.11.4	ollama-gpu/0.11
ollama	0.11.4	ollama/0.11

Open OnDemand Tool: LM Studio

LM Studio: Swan

This app will launch the LM Studio GUI. You will be able to interact with the LM Studio GUI through a VNC session.

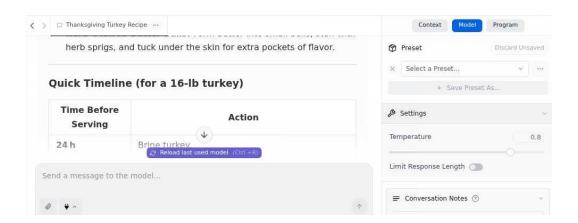
LM Studio version

0.3.23

This defines the version of LM Studio you want to load.

LM Studio

- Beginner-friendly, no-code setup with Graphical user interface (GUI)
- Ideal for:
 - Local testing and prototyping
 - GUI-based prompt engineering
 - Lightweight experiments
 - Non-developer use



Ollama

- Command-line interface (CLI) and flexible API for developers
- Ideal for:
 - Scripting and chaining models in developer workflows
 - Seamless integration into backend systems
 - Fine-tuning and customizing LLM pipelines

More info: https://hcc.unl.edu/docs/applications/app_specific/using_llm/

Monitor Resource Usage

Monitoring Resources Usage is important in optimizing machine learning and deep learning models

- Hyper-parameter decision based on memory usage
 - batch size larger batch size needs more memory
 - training precision higher precision (32/16/8-bit) needs more memory
 - model architecture deeper model (AlexNet vs. ResNet) needs more memory
- Computing resource usage (%CPU/GPU)
 - long GPU idle time large data loading overhead

Monitor Resource Usage

For more tips for monitoring the running jobs see:

Monitoring CPU Usage:

https://hcc.unl.edu/docs/submitting_jobs/monitoring_jobs/

Monitoring GPU Usage:

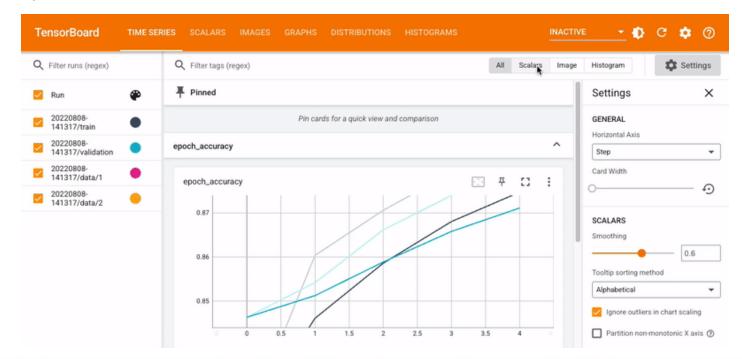
https://hcc.unl.edu/docs/submitting_jobs/monitoring_gpu_usage/

```
top - 14:22:31 up 7 days, 21:16,  0 users,  load average: 1.85, 2.04, 1.64
Tasks: 812 total, 3 running, 809 sleeping,
                                             0 stopped,
%Cpu(s): 8.1 us, 1.0 sy, 0.0 ni, 89.9 id, 0.8 wa, 0.1 hi, 0.1 si, 0.0 st
MiB Mem : 256738.5 total, 247767.9 free, 4959.9 used,
MiB Swap:
              0.0 total.
                              0.0 free,
   PID USER
1736758 hccdemo
                      0 5940976
                                 1.1q 191832 S 135.9
1736757 hccdemo
                      0 5938428 1.1g 192432 S 135.2
1736756 hccdemo
                      0 6004472
                                 1.1a 192460 S 134.9
```

Every	1.0s: r	ıvidia-smi		c2420.swan	.hcc.unl.edu:	Fri Sep 27	14:28:23 2024
Fri Se	ep 27 14	4:28:23 2024					
NVID	IA-SMI	550.107.02	Driver	Version: 5	50.107.02	CUDA Versio	on: 12.4
GPU Fan	Name Temp	Perf	Persistence-M Pwr:Usage/Cap	Bus-Id 	Disp.A Memory-Usage		Uncorr. ECC Compute M. MIG M.
0 N/A 	Tesla 34C	V100-PCIE-32GB P0	======================================		0:5E:00.0 off В / 32768мів	+======== 4% 	0 Default N/A
+ 1 N/A 	Tesla 51C	V100-PCIE-32GB P0	On 46w / 250w		0:D8:00.0 off B / 32768MiB	+ 27% 	0 Default N/A

TensorBoard is useful to analyze the model performance during or post training

- Two steps in general
 - add logging code and write to file
 - open a TensorBoard session to check



Add logging code and write to file

For example, write the model graph in PyTorch:

```
Python

from torch.utils.tensorboard import SummaryWriter
writer = SummaryWriter()
writer.add_graph(model, dummy_input)
```

You can add metrics recording such as loss, accuracy, recall, etc. in training loop, see more in:

https://pytorch.org/docs/stable/tensorboard.html

https://pytorch.org/tutorials/intermediate/tensorboard_tutorial.ht

ml

Similarly, write the model graph in TensorFlow, for example:

```
Python

/ • □

writer = tf.summary.create_file_writer(log_dir)

# Trace and export the graph
tf.summary.trace_on(graph=True, profiler=False)
forward_pass(sample_input)
with writer.as_default():
    tf.summary.trace_export(
        name="model_graph",
        step=0,
        profiler_outdir=log_dir)
```

Find more tutorials on logging models at:

https://www.tensorflow.org/tensorboard/get_started

Open an interactive TensorBoard session to check the logged information

Open OnDemand:
 Interactive Apps → TensorBoard

TensorBoard: Swan version:										
vv1.0.0_13_g6ed6278										
This app will launch TensorBoard, TensorFlow's Visualization Toolkit.										
Additional modules (optional)										
Space separated list of additional modules to load (eg.										
tensorflow/py38/2.4)										
Full list for Crane or Rhino	_									
Tensorboard logdir] _	c	20	۰if	. , I	00	. fil	0 r	\a+	h
_		<u>ာ</u>	pe	CII	уι	Ug	, 111	e p	Jat	[] —
Select the directory that contains data to visualize; defaults to										
\$HOME.										
Select Path										

Setup for Hands on

In this session, we will use Jupyter Lab via Swan's Open OnDemand portal for hands-on practice.

- 1. Login to Swan's Open OnDemand portal: https://swan-ood.unl.edu/
- 2. Copy the practice notebook and data to your \$WORK directory:
 - using terminal:
 - Select "Clusters" → "Swan Shell Access"

```
$ cd $WORK
$ git clone https://github.com/unlhcc/hcc-ai-ml-workshop-2025.git
$ ls hcc-ai-ml-workshop-2025
```

If you **are done**, please put up your yellow sticky note.

If you **need help**, please put up your red sticky note.

Setup for Hands on

- 3. Select "Interactive Apps" → "Jupyter Lab"
- 4. Navigate to the "hcc-ai-ml-workshop-2025" folder in left-hand sidebar
- 5. Open the "00_prepare_datasets.ipynb" notebook and run all cells

Open OnDemand Settings

Parameter	Value
Jupyter Lab version	4.0
Working Directory	/work/groupname/username
Number of cores	8
Running time in hours	3
Requested RAM in GBs	32
Partition selection	guest_gpu
Reservation	Located on the back of your name tag
GRES	gpu
Job Constraints	Leave BLANK

Note that, since this is an intermediate workshop, we assume that participants:

- Have basic knowledge of Python, so we will not explain Python code.
- Are familiar with using Swan, including terminal access and interactive apps.

Open OnDemand access (https://swan-ood.unl.edu/): Interactive Apps → Jupyter Lab

Select preferred Jupyter Lab (Ver. 3.4 or 4.0)

Jupyter Lab version:

vv1.0.0_16_g3837770

This app will launch a Jupyter Lab server using Python.

Jupyter Lab version

4.0

This defines the version of Jupyter Lab you want to load.

Working Directory

Parameters to request resources:

- Jupyter Lab version
- Path to working directory

Jupyter Lab version

4.0

This defines the version of Jupyter Lab you want to load.

Working Directory

/work/<groupname>/<username>

Select your Notebook directory; defaults to \$HOME

Select Path

Parameters to request resources: CPU related resources

Number of cores
4
Number of cores requested on a node (min 1, max 16)
Running time in hours
1
Maximum runtime in hours of Jupyter Lab server (min 1, max 8)
Requested RAM in GBs
16
Maximum memory requested for Jupyter Lab server (min 2GB, max 60GBs)

Parameters to request resources: GPU related resources

- Type "guest_gpu" in Partition selection and GRES to request GPU
- "jupyter" or "batch" for CPU resources

Partition selection

jupyter

- jupyter Resources reserved for Jupyter notebooks.
- batch The general batch queue.
- other partitions may be specified if your account has access

GRES

This field is used primarily for gpu submissions. You must specify a gres of at least gpu when using the GPU partition. Other possible values may be found here.

Parameters to request resources: other parameters

Reservation

Located on the back of your name tag

Submit to a specific reservation if your account has access.

Job Constraints

Additional constraints for the job. Primarily used for specifying a GPU type or node type.

Open OnDemand access:

Wait for your Jupyter Lab session to start

Jupyter Lab (9470413)

Queued

Created at: 2025-02-24 12:49:41 CST



Time Requested: 1 hour

Session ID: f40577c3-5a1b-401b-bb24-63a0d618ad2a

Please be patient as your job currently sits in queue. The wait time depends on the number of cores as well as time requested.

Setup for Hands on

- 3. Select "Interactive Apps" → "Jupyter Lab"
- 4. Navigate to the "hcc-ai-ml-workshop-2025" folder in left-hand sidebar
- 5. Open the "00_prepare_datasets.ipynb" notebook and run all cells

Open OnDemand Settings

Parameter	Value
Jupyter Lab version	4.0
Working Directory	/work/groupname/username
Number of cores	8
Running time in hours	3
Requested RAM in GBs	32
Partition selection	guest_gpu
Reservation	Located on the back of your name tag
GRES	gpu
Job Constraints	Leave BLANK

Note that, since this is an intermediate workshop, we assume that participants:

- Have basic knowledge of Python, so we will not explain Python code.
- Are familiar with using Swan, including terminal access and interactive apps.

Schedule

timestamp	Session			
12:00 ~ 12:15pm	Setup and Support			
12:15 ~ 1:00pm	Introduction to HCC resources for ML and AI (45 Min)			
1:00 ~ 1:15pm	Short Break (15 Min)			
1:15 ~ 2:00pm	ML and AI workflows (45 Min)			
2:00 ~ 2:10pm	Short Break (10 Min)			
2:10 ~ 2:50pm	Introduction to PyTorch (40 Min)			
2:50 ~ 3:00pm	Break (10 Min)			
3:00 ~ 3:50pm	Introduction to PyTorch cont. (50 Min)			
3:50 ~ 4:00pm	Break (10 Min)			
4:00 ~ 4:15pm	Introduction to National Research Platform (NRP)			
4:15 ~ 4:30pm	Open Questions			