# Georgia Institute of Technology

ECE 8803: Hardware-Software Co-Design for Machine Learning Systems (HML)

Spring 2025

Lab 4B

Due: Monday, April 7, 2025 @ 11:59 pm EST

Created by: Changhai Man (cman8@gatech.edu), William Won (william.won@gatech.edu)
Last Revised: Mar 27, 2025

## Instructions
**Please read the following instructions carefully.**
- The lab is distributed in 3 parts, each assigned with 3 points.
- You will need to modify the code in various files, submit generate file in a pack. At the final part of the notebook there will be codes to help you pack your code.
- Rename the zip according to the format:
  **LastName_FirstName_ECE_8803_HML_sp25_lab4B.tar.gz**
- It is encouraged for you to discuss homework problems amongst each other, but any copying is strictly prohibited and will be subject to Georgia Tech Honor Code.
- Late homework is not accepted unless arranged otherwise and in advance.
- Comment on your codes.
- For all problems, please post queries on piazza. If you add a comment to an answered
- query, make sure to change the comment to "Unresolved".

## Lab Description
In lab 3B, you've experienced executing and profiling GPT training jobs **over a real cluster**, using a few different parallelization strategies on GPT. However, real clusters might not always be available, especially when you're studying futuristic ML platforms. In such scenarios, **simulation-based evaluations of ML execution are necessary.**

In this lab, you'll familiarize yourself with distributed ML simulation frameworks, ASTRA-sim. Specifically, you'll first be asked to *implement a new network topology* for the simulator (this kind of approach is inevitable since building a physical network topology is hard). Then, you will use STAGE+AstraSim to *simulate* distributed LLM training, and try to *find out what is the optimal parallelization strategy* for different target workloads and systems.
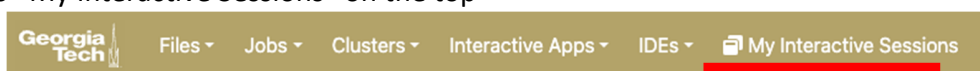
## Lab Layout

- Part 0: Environment Setup [0 pts]

- Part 1: Implementing New Network Topology [3 pts]

  o Task 1.1: Implement physical connectivity of 2D Mesh [1 pt]

  o Task 1.2: Implement xy routing over the 2D mesh [2 pts]

- Part 2: Using STAGE+AstraSIM for basic simulation [3 pts]

  o Task 2.1: Generate Chakra Workload with STAGE [1 pt]

  o Task 2.2: Generate AstraSim System/Network Configs [1 pt]

  o Task 2.3: Run AstraSim with generated workloads/configs, and Extract results [1 pt]

- Part 3: Find optimal parallel stratrgies for different system/models

  o Task 3.1: Generate Parallel Strategies Design Space [1 pt]

  o Task 3.2: Doing Design Space Exploration in Batch [1 pt]

  o Task 3.3: Find optimal parallel strategies for each setup, and why? [2 pt]
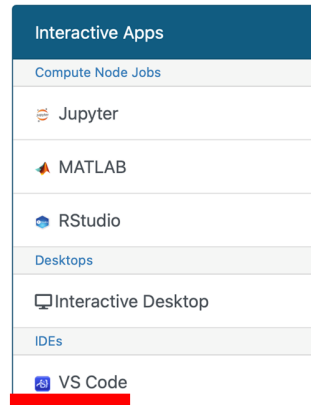
## Lab Environment Setup

In this lab, you will follow the instructions and update the code in a Jupyter Notebook file named "lab4b.ipynb", similar to Lab 3. There are several tools being used, especially they contain some customization for this lab. Therefore, we require you use the PACE-ICE machine when working on this lab. Please follow the instructions below to launch this Jupyter Notebook file on the PACE-ICE OnDemand session:

**Pace OnDemand Setup**
1. Configure GlobalProtect VPN. Refer to the following link.
   https://vpn.gatech.edu/
2. Use the following link to access to OnDemand ICE cluster.
   Link:
   https://ondemand-ice.pace.gatech.edu/pun/sys/dashboard
3. Go to "My Interactive Sessions" on the top
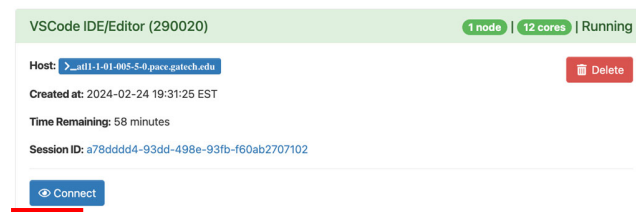
4. Select "VSCode IDE/Editor version" under IDE.



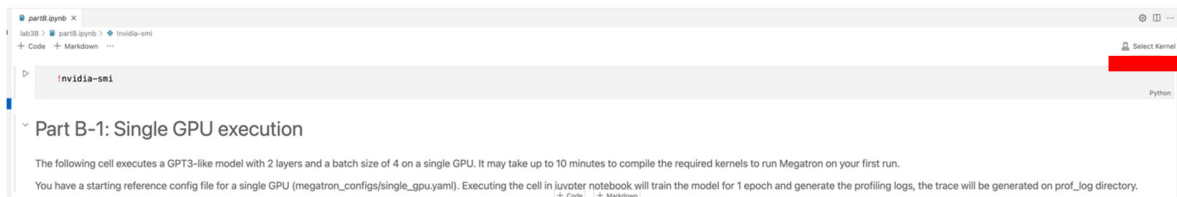5. Configure the setting as follows:
    - **Modules:** Custom Pre-load
    - **Custom Pre-load Commands:**
      *module load anaconda3 gcc*
      *bash /storage/ice-shared/ece8803hml/lab4B_setup.sh*
      **Quality of Service:** Default (none)
    - **Node Type:** Intel CPU
    - **Cores:** 6
    - **Memory (GB):** 48
    - **Number of hours:** 4 # Or whichever hours you need. Please be respectful to others; delete instances immediately after completing the lab.
6. Once configured launch the session. Your session will start soon (This can take some time if there is heavy traffic).
7. Click "connect" to open the session.



**VSCode Setup**
1. Click "open folder" on the left, and open:
   */home/hice1/$YOURGATECHID/ece8803_hml_lab4B*
    - If the folder does not appear, verify your configuration in step 5.
2. Open "lab4B.ipynb" in directory
3. Click "Select Kernel" on top-right

4. Select Kernel → Python Environment → */home/hice1/$YOURGATECHID/ece8803_hml_lab4B/.conda/bin/python*
   or
   *.conda/bin/python*

5. Now you're ready to work on the actual lab assignments. Please follow the instructions in the Jupyter Notebook.

## Lab Submission

- Submission due: Monday, Apr 7, 2025, by 11:59 pm EDT.
- Submission format: 1 tar.gz with following contents
  - **Mesh2D.cpp** with finished codes
  - **lab4b.ipynb** with finished codes
  - contents generated in
    */home/hice1/$YOURGATECHID/ece8803_hml_lab4B/submission*