

Georgia Institute of Technology

ECE 8803 HML - Spring 2025

Lab3 B, C

PartB Due Dates:

- Last Name A-H : Monday, March 3, 2025 @ 11:59 pm EST
- Last Name J-P : Monday, March 10, 2025 @ 11:59 pm EST
- Last Name Q-Z : Monday, March 17, 2025 @ 11:59 pm EST

PartC Due Date:

- Monday, March 17, 2025 @ 11:59 pm EST

Instructions

Please read the following instructions carefully.

- The lab is divided into three parts: A, B, and C.
- Part B has different deadlines based on last name groups. Ensure submission before the deadline.
- It is encouraged for you to discuss homework problems with each other, but any copying is strictly prohibited and will be subject to the Georgia Tech Honor Code.
- Late homework is not accepted unless arranged otherwise and in advance.
- 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".

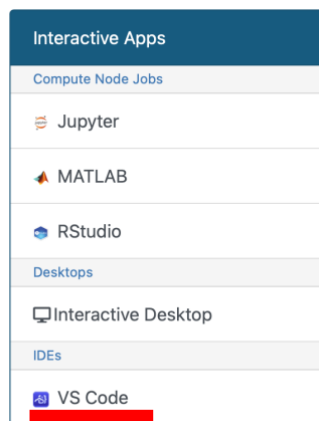
Lab Setup

Pace OnDemand Setup

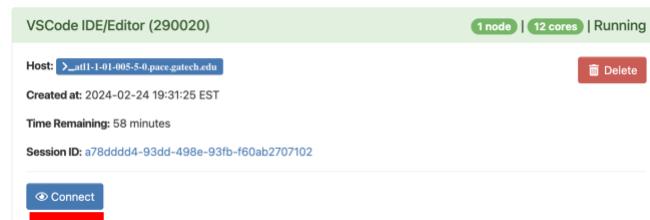
1. Configure GlobalProtect VPN. Refer to the following link.
https://gatech.service-now.com/home?id=kb_article_view&sysparm_article=KB0026837
2. Use the following link to access to OnDemand ICE cluster.
[Recommended browser: Google Chrome. Disable any contents/ad blockers. TensorBoard may not work properly on other browsers such as Safari.]
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.

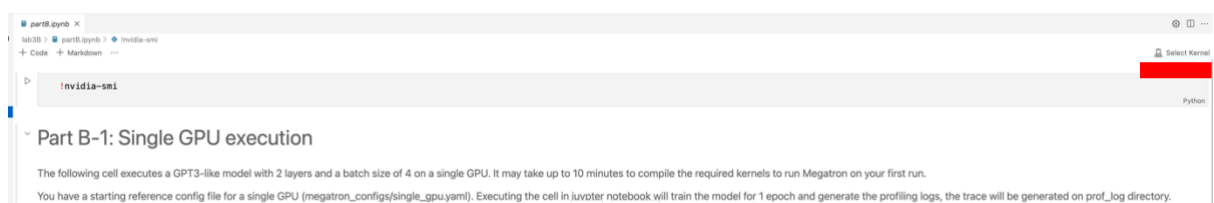


5. Configure the setting as follows:
 - **Modules:** Custom Pre-load
 - **Custom Pre-load Commands:**
`sh /storage/ice-shared/ece8803hml/setup.sh`
`module load anaconda3/2023.03`
 - **Quality of Service:** Default (none)
 - **Node Type:** NVIDIA GPU H100 HGX
 - **CPUs:** 16
 - **Memory (GB):** 256
 - **GPUs:** 2
 - **Number of hours:** 2 # 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/ lab3B
 - If the folder does not appear, verify your configuration in step 5.
2. Open "partB.ipynb" in directory
3. Click "Select Kernel" on top-right



4. Click "Install suggested extensions"
5. Press ctrl+shift+p (or cmd+shift+p for Mac) to open Command Palette, then :

- Type ">Python: Select Interpreter" and press enter.
 - Click "Select at workspace level" - "Enter Interpreter Path" - "Find"
 - Enter "/storage/ice-shared/ece8803hml/envs/lab3/bin/python3"
 - Click "Ok"
6. Click "Select Kernel" - "Python Environments" - "lab3 (Python 3.10.6)"
 7. Run the first cell in the notebook. Check if it prints the message "Profiling done successfully". If there is any error, check the OnDemand configurations
 - Sometimes Jupyter Notebook will fold the outputs, and you might miss the part you are looking for. Please view the outputs of a cell by clicking "in a text editor".

Lab Description

As part of this lab, we will experiment with training a GPT-like model. We will do this using a single and multiple GPUs with various types of parallelisms. We use [Megatron-LM](#) for generating distributed training runs.

PartB-1. Analyzing Multi-GPU training. [4 points]

1. Next, generate the proper Megatron configuration file in the "megatron_configs" directory for the following parallelism strategy.
 - (a) Single GPU. (Already Provided)
 - (b) Tensor Parallelism (TP=2) on 2 GPUs. (Already Provided)
 - (c) Pipeline Parallelism (PP=2) on 2 GPUs.
 - (d) Data Parallelism (DP=2) on 2 GPUs.
 - (e) Tensor Parallelism (TP=2) + [Activation Recomputation.](#)
 - (f) Pipeline Parallelism (PP=2) + [Activation Recomputation..](#)
 - (g) Data Parallelism (DP=2) + [Activation Recomputation..](#)

The degree of data parallelism is not explicitly specified but is automatically inferred as follows:

$$DP = \text{num_processes} / (PP * TP)$$

2. Using [TensorBoard](#), identify various characteristics of different training runs and fill the provided Excel sheets answer.xlsx for all 7 configurations.

In case tensorboard doesn't work:

1. Open the terminal and activate the 'lab3' environment with conda activate.
2. Run "python3 -m pip install -U torch-tb-profiler -user"
3. Now rerun the cell to start the tensorboard

Navigating tensorboard

You can access various analyses using the panels on the left.

- **Runs:** You can change between different traces
- **Views:**
 - o **Trace:** Find the latency of “Single Iteration”, “Forward Pass”, “Backward Pass”, and “CoreAttention” time. Use the search bar in the top-right corner of TensorBoard to find the names of the relevant operators. For “CoreAttention”, sum the latency of all “CoreAttention” operators. Use “wall duration” for reporting time.
 - o **Operator:** Find the operator with the highest “host self time” here.
 - o **Memory:** Find the peak memory usage here
- **Workers:** You can switch between different workers (GPUs) to toggle between multiple GPUs. The “GPU Summary” in the Overview displays the current GPU number (e.g., GPU 0).

[PartB-2. Training a large model \[2 point\]](#)

Train the largest possible GPT-like model with a batch size of 4 on 2 H100 GPUs by modifying the following parameters:

- **Model Size:** Edit the "model_configs/gpt3_27.json" file. You may only modify the number of layers ("n_layer": 24). Set the number of layers to a multiple of 24.
- **Distributed Training Configuration:** Choose any of the six configurations from Part B-1.

Report the largest possible model size and used distributed training configuration in answer.xlsx.

[PartC. Discussion \[2 points\]](#)

Answer the 5 questions in the discussion.txt file provided, based on the results of Part B. Each question is worth 0.4 points.

[Lab Submission](#)

Submit the following 2 files.

- answer.xlsx for part B
- discussion.txt for part C