Introduction to GPU Acceleration on Alpine

CURC Summer Camp 2023

GPU-Accelerated Matrix Multiplication

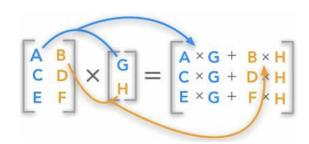
- A matrix is a rectangular array of numbers arranged in rows and columns.
- The order of a matrix is the number of rows and columns.

$$\begin{bmatrix} 6 & 2 & 4 \\ 2 & 1 & 6 \end{bmatrix}$$

order = 2x3

Matrix multiplication is used in many ML algorithms.

We will use Tensorflow's tf.matmul function.





Login to CURC systems using your credentials. CU Boulder and CSU users can login directly from the terminal. RMACC and AMC users will need to go through Open OnDemand.

https://curc.readthedocs.io/en/latest/access/rmacc.html#logging-in-to-open-ondemand

Navigate to your projects directory:

```
cd /projects/$USER/
```

If you haven't already grabbed the Summer_Camp_2023 repository, do a git clone:

```
git clone
https://github.com/ResearchComputing/Summer_Camp_2023.git
```

Otherwise, run git pull.

Make the scripts executable:

```
cd
/projects/lafr9499/Summer_Camp_2023/Day_Three/GPU_Acceleration
chmod u+x tf*
```

Start an interactive node on Alpine with one NVIDIA A100 GPU and 20 CPU cores:

```
sinteractive --time=40:00 --ntasks=20 --gres=gpu:1 --
partition=aa100 --reservation=nvidia_summercamp
```

Look at the software and compilers available on Alpine:

```
module avail
```

Activate the conda environment:

```
module load anaconda
conda activate /curc/sw/conda_env/tf-gpu-cuda11.2
```

Run the CPU-only script:

```
python tf.matmul-cpu.py
```

Run the GPU-accelerated script:

```
python tf.matmul-gpu.py
```

Discussion Questions

Try running each of these scripts a few times. Why might the first GPU script run take longer than subsequent runs?

Try adjusting the size of the matrices. What do you notice about GPU speedups when the matrices are very small? What could be limiting the size of the matrices you are able to try right now?

How could we speed up the matrix multiplication even further?

Start PyTorch on AMD GPUs

```
#get on an AMD GPU node with one GPU and 5 CPU cores
sinteractive --time=01:00:00 --ntasks=5 --gres=gpu:1 --
partition=aa100 --reservation=amd_summercamp

#load modules
module load rocm
module load pytorch

#start python
python3

#from the python shell
>>> import torch
>>> torch.cuda.is_available()
True
```

>>> print(torch.__version__)
1.13.0+rocm5.2