

GPU Programming Basics

Northeastern University
NUCAR Laboratory

for all

Julian Gutierrez
David Kaeli

Hands-on Lab #1

Objective

- Setup computers to be able to access server.
- Learn how to use the server.
- Write our first GPU program.

NOTE: Unless stated otherwise, work in pairs.

Part 1: Setting up computer

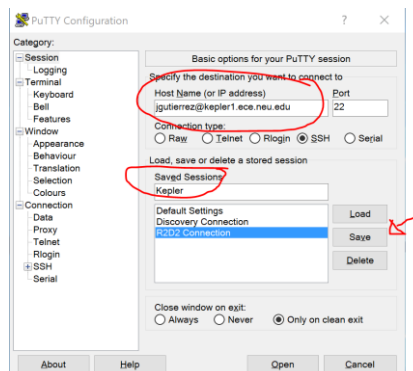
To be able to access the discovery cluster we need an ssh connection. Given most students usually have windows machines, we require additional software to be able to access the server.

NOTE: If you have Mac or Linux, you should be able to use ssh directly from the terminal.

- Install PuTTY which allows you to create a connection through a terminal:
<http://www.chiark.greenend.org.uk/~sgtatham/putty/latest.html>
- Or you can also install MobaXterm:
 - I prefer this one as it has X forwarding by default amongst other useful features
<https://mobaxterm.mobatek.net/download-home-edition.html>
- Install FileZilla Client software:
 - FileZilla will allow you to copy files back and forth from the server. It's easier if you plan on editing the codes on your machine and then copying them back to the server.
<https://filezilla-project.org/download.php?type=client>

How to use PuTTY

1. Open PuTTY and input your username along with the server address under Host Name: <username>@login.discovery.neu.edu
2. Under Saved Sessions, write "Discovery" and click on Save.
3. Now Discovery should appear on the list. Click on it and then click load.
4. Accept the dialog box that appears.
5. Input password.
6. You're all set. Connection to server should be successful.
7. **NOTE: Image is just for reference.**



GPU Programming Basics

Northeastern University
NUCAR Laboratory

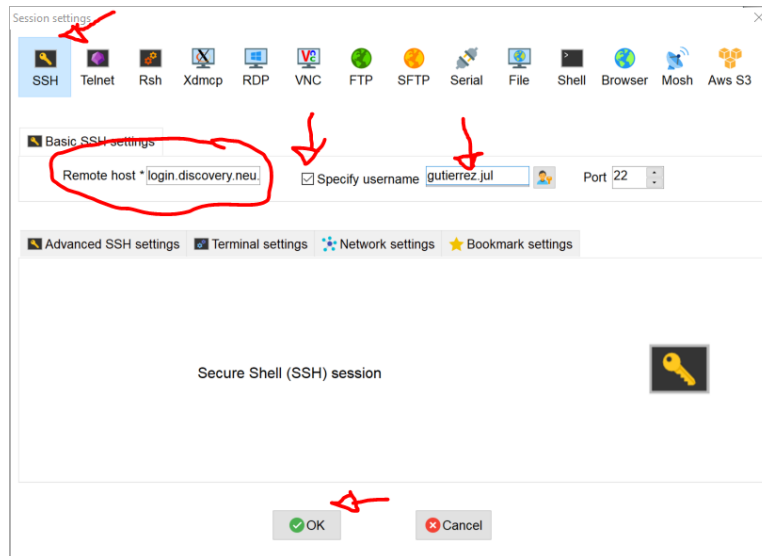
for all

Julian Gutierrez
David Kaeli

Hands-on Lab #1

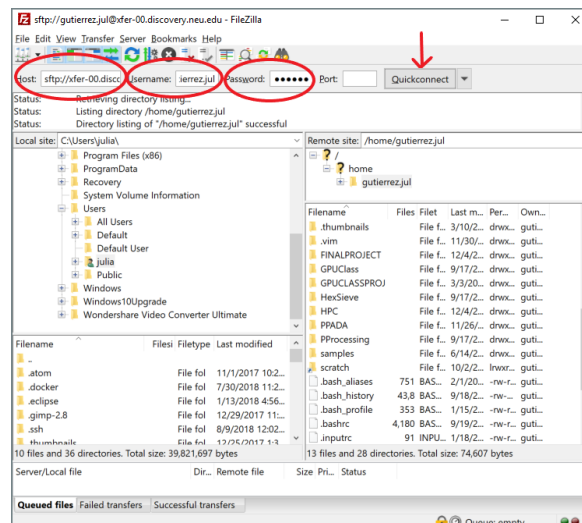
How to use MobaXterm

1. Open MobaXterm and click on Session (Top left) -> SSH
2. Remote Host: login.discovery.neu.edu
3. Click on Specify username
4. Write your username in the box
5. Click ok and accept any pop-up that might appear.



How to use FileZilla

1. Open FileZilla and input the following:
 - a. Host: xfer-00.discovery.neu.edu
 - b. Username: <Username>
 - c. Password: <Password>
 - d. Port: 22
 - e. Click on Quick-connect and it should be able to connect to the data transfer server.



GPU Programming Basics

Northeastern University
NUCAR Laboratory

for all

Julian Gutierrez
David Kaeli

Hands-on Lab #1

2. Click Quickconnect
3. And you should be done. Now you should be able to copy files from your computer (Local site to the left) to the server (Remote site to the right).

How to use SCP (On Linux)

1. Open your terminal and find the file you want to transfer to the server (the command would be very similar if you want to copy a file from the server to your computer):
 - a. The general command for scp is:
 - i. `scp <file to transfer> <where to transfer it too>`
 - b. The host that we need to use to transfer data to and from the discovery cluster is:
 - i. **Host: xfer-00.discovery.neu.edu**
 - ii. To point to your home directory on the server, the address would be:
 1. `<username>@xfer-00.discovery.neu.edu:~/`
 2. The “~/” part means your home directory
 - c. An example of this command:

```
scp /home/julian/file gutierrez.jul@xfer-00.discovery.neu.edu:~/folder/
```

GPU Programming Basics

Northeastern University
NUCAR Laboratory

for all

Julian Gutierrez
David Kaeli

Hands-on Lab #1

How to use the Discovery Cluster

1. First time accessing the discovery cluster (or using Putty)

```
ssh -X <USERNAME>@login.discovery.neu.edu
```

2. After typing the password, this will take you to a login node. **YOU SHOULD NOT** execute any work here. These nodes are meant to control the flow of users and shouldn't be used for anything besides that (this means, no compiling/running code on them). You should only run code by requesting a node with a SLURM script or submitting an interactive node (we will go into this shortly).
3. We have been granted an exclusive partition (group of nodes) for us to work with. They are 16 GPU nodes with K20 GPUs. This partition (gpu-course) is only available from 5-9 pm during the dates when we have lab sessions. Access to this partition is only available to us.
4. Outside of these times, you will have to submit your jobs to the gpu partition (available to all discovery users where it might be slow sometimes depending on how many people are using it).
5. We are first going to setup the environment in the discovery cluster. We are going to store all our codes in a general folder called scratch. Each person **must** create their own folder in scratch by running this command:

```
mkdir /scratch/`whoami`
```

6. The /scratch directory is a 1+ PB parallel file server which all users may utilize. Always launch submitted jobs from /scratch, and transfer your data from /scratch to another resource (e.g. personal/group file server) as soon as possible. For instance if you wrote code stored in scratch, you want to copy it back to your home directory as the data in /scratch might be removed after a couple of months.
7. In order to run any code on Discovery, you must use SLURM (Simple Linux Utility for Resource Management). The idea is that you ask SLURM for resources (e.g. what type of node, how many cores, how much memory), and then your code will execute once those resources are available. We will go into this later on in the lab.
8. Some common SLURM Commands:
 - o sbatch <file name> : this will send the job to the scheduler.
 - o srun : If you would prefer to work interactively on a node, you may launch and interactive session with srun.
 - o squeue : see what jobs are waiting and running currently.
 - o scancel <job id> : remove a running or pending job from the queue
 - o scontrol <flags> : find more information about the machine configuration and job settings
 - o seff <job id> : report the computational efficiency of your calculations
9. We also want to setup our environment to have access to the tools we will use (compilers, etc). You can do this by loading the necessary modules (cuda). You can add the following lines at the end of your ~/.bashrc file (in your home directory). This is a file which will be loaded every time you connect to the discovery cluster.
 - a. Open file using an editor

```
vim ~/.bashrc
```

- b. Go to the end of the file and copy these lines.

GPU Programming Basics

Northeastern University
NUCAR Laboratory

for all

Julian Gutierrez
David Kaeli

Hands-on Lab #1

```
# Manually load the required modules
module load cuda/9.0
```

- c. Reload the file.

```
source ~/.bashrc
```

10. Each one of you will initiate an interactive (1-core non-exclusive) session on the nodes just to compile and modify the code.
11. To request an interactive node, we can use the following command. Once a resource is available, you will be connected automatically into the compute node.

```
srun --pty --export=ALL --reservation=gpu-class --partition=gpu --task
s-per-node 1 --nodes 1 --mem=2Gb --time=02:00:00 /bin/bash # Reserve t
he node
```

12. You can run the `squeue` command to guarantee that you are using a compute resource. This will print the nodes allocated for the user:

```
squeue -u `whoami` # Print the nodes that #USER has reserved
JOBID  PARTITION   NAME     USER ST       TIME  NODES NODELIST(REASON)
678595      gpu       bash     gutierre R        0:48      1 c2133
```

13. From now on, you can compile your work.
14. Remember, the node that we are connected to has a `gpu`, so a way to detect if the `gpu` is there is by running the “`nvidia-smi`” command (this prints general information about any `cuda` device in the machine):

```
gutierrez.jul@c2133:[gutierrez.jul]$ nvidia-smi
No devices were found
gutierrez.jul@c2133:[gutierrez.jul]$
```

15. The reason we get this is because even though we are connected to the node with a GPU in it, we didn't reserve the GPU in the `srun` command. We will leave access to the GPU only for the code runs that will actually be using the GPU by using batch submission, requesting 1 core and 1 GPU, each (we will look at this later on).

IMPORTANT: When you are done with everything, type the following command to make the resources available to other users.

```
exit # this will take you to the gateway computer again
```

Please, Remember to **ALWAYS** release the node. Always!

GPU Programming Basics

Northeastern University
NUCAR Laboratory

for all

Julian Gutierrez
David Kaeli

Hands-on Lab #1

Part 2: Writing our first GPU Program

Copy the baseline for our code from the following directory either to your home directory in the server or into your machine (we will be editing this file).

Using the following commands as a guideline (feel free to do a different approach):

```
cd /scratch/`whoami`/  
mkdir GPUClass18  
cd GPUClass18  
cp -r /scratch/gutierrez.jul/GPUClass18/HOL1/ .  
cd HOL1/  
vim VectorAdd.cu
```

The code we will be working on is meant to do a vector add on the CPU, do the same vector add with the GPU and compares the results at the end.

1. Open the file. Find the comments “HERE” and fill the code with the necessary commands to make it work.
2. The following are the templates for the commands that you’ll need:

```
cudaError_t cudaMalloc ( void** devPtr, size_t size )  
  
cudaError_t cudaMemcpy ( void* dst, const void* src, size_t count,  
    cudaMemcpyKind kind )  
  
    kind: cudaMemcpyHostToDevice or cudaMemcpyDeviceToHost  
  
kernel_name <<<grid, block>>>(argument list);  
  
cudaFree ( void* ptr);
```

3. Look at the kernel function and the arguments it receives. Write the code to calculate the vector add:
 - a. $C_i = A_i + B_i$
4. To compile the code, run the following command:

```
nvcc -arch=sm_35 -O3 VectorAdd.cu -o vAdd
```

5. Once your code is compiling without errors, we can run the code.
6. To run the code, we must submit a job to use the GPU. The script we will use is the one in the folder named exec.bash. Look at what it does and make sure the directory is correct, then submit it using the following command:

```
sbatch exec.bash
```

7. Once it has run, open the output files named “exec.out”. Make sure it gives the correct result. If there are any errors, fix them.
 - a. Observe that there are some commands that are preceded by a +. This is just to indicate the command that was executed.
8. Take the following observations into account and think of the answers to these questions while testing the code.

GPU Programming Basics

Northeastern University
NUCAR Laboratory

for all

Julian Gutierrez
David Kaeli

Hands-on Lab #1

- a. Grid Size and Block Size should be chosen based on the variables already available.
- b. Notice cudaMalloc receives a ****variable**. What does this mean?
- c. Why are we allocating an array called `c_gpu` on the CPU?
- d. Test different vector sizes and block sizes to see if the result is still correct.
 - i. Compare the performance when you increase the block size and when you decrease it.
 - ii. Do this by modifying the `exec.bash` script and adding more lines to execute different scenarios. E.g.:
 1. `./vadd 10000000 32`
 2. `./vadd 10000000 64`
 3. `./vadd 10000000 128`
 4. `./vadd 10000000 256`
 5. `./vadd 10000000 512`
 6. Submit the script again and wait for it to complete
 - iii. What could be affecting the behavior?
 - iv. Does GPU perform better than CPU? Try different vector sizes (100, 100 000, 100 000 000)
 - v. Try running using different block sizes for a big vector. How does that affect performance? Remember the 1024 maximum threads per block limit.
 - vi. Should the execution time include the time it takes to allocate and copy data?
9. Was it easier to implement the GPU version compared to writing an equivalent program on the CPU?

NOTE: The GPU reservation will be available from 5 pm to 9 pm on lab days. If you want to work on this lab outside of these hours, you need to remove the `--reservation` from the scripts and your `srun` command. The job will now go to the `gpu` partition which is shared with other researches/students at NEU. This means your jobs might take longer to run depending on how many jobs are running. You can check for this information by using the following commands:

```
sinfo -p gpu  
squeue -p gpu
```

Part 2: CUDA-MEMCHECK

CUDA memory checker is a very useful tool to check if your code is working correctly. It looks for illegal accesses in memory (accessing an illegal pointer, or when you accessed an array beyond its size, etc). In addition, it will let you know if there were any issues with any Cuda command.

To run the tool, you use the following command:

```
cuda-memcheck <command used to execute the code>
```

Use the `memcheck` bash script to submit a job to run `cuda-memcheck` with your executable. When you run it with your code, are there any errors? If so, fix them.

```
sbatch memcheck.bash
```

GPU Programming Basics

Northeastern University
NUCAR Laboratory

for all

Julian Gutierrez
David Kaeli

Hands-on Lab #1

Part 3: NVPROF

NVPROF is a great tool to profile your application and understand what is happening and how to improve it. Use the executable you have compiled and works correctly to do the following tests.

We will use the `nvprof.bash` script to submit a job that runs `nvprof`. Make sure the run is using a vector size of 100 000 000 and the block size is set to 512.

As an example, this is the command the script should execute to run NVPROF:

```
nvprof ./VectorAdd 100000000
```

This will print out a lot of information. When you submit the script and it executes, look for “NORMAL RUN” results. Look at the percentages and answer these questions:

- What information do you find useful?
- What insights do you get from this?
- How would you improve this code? What would you improve first?
- Notice that the code runs slower with the profiler, why is that?

Now, another way of running `nvprof` is using the following command (the script runs it as well):

```
nvprof --print-gpu-trace ./VectorAdd 100000000
```

- From the previous run outputs, look for “GPU TRACE”. Any useful information displayed running it like this?
- Imagine your code was much more complex. Would this information give you an insight as to how your code is behaving?

Now, the final command we used in the script is to recollect performance metrics. Look for “METRICS” in the same file output. Look at some of the metrics by running the following command (It will calculate all of the metrics):

```
nvprof -m all ./VectorAdd 100000000
```

Open the output file and look at the results.

- Which metrics do you think are useful?
- If you look at `gld_efficiency` and `l2_l1_read_hit_rate`, what do you think is happening on the memory? Is there something we can do to improve this?
- How about `gst_efficiency`?
- There are many metrics. Try to find which ones are the most useful.
- If you want to compare different runs, try setting the block size to a small value (32) and collecting the metrics, and then run it with a larger value (1024).
 - Can you notice any differences in the metrics? In the speed?

Note: Please remember to ask as many questions as possible. I am here to help as much as we can.