

# NYU HPC Cloud Bursting Tutorial

Course: CSCI-UA.0480-051 (Fall 2025)

## Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Your Resource Allocation</b>                                 | <b>3</b>  |
| 1.1      | Available Partitions . . . . .                                  | 3         |
| <b>2</b> | <b>Prerequisites</b>  | <b>3</b>  |
| 2.1      | Network Connection . . . . .                                    | 3         |
| <b>3</b> | <b>Quick Start Guide</b>  | <b>3</b>  |
| 3.1      | Method 1: Using Open OnDemand (Recommended) . . . . .           | 3         |
| 3.1.1    | Step 1: Access Open OnDemand . . . . .                          | 3         |
| 3.1.2    | Step 2: Upload Your Files . . . . .                             | 3         |
| 3.1.3    | Step 3: Open a Terminal . . . . .                               | 4         |
| 3.2      | Method 2: Using Terminal/SSH . . . . .                          | 4         |
| 3.2.1    | SSH Connection . . . . .  | 4         |
| 3.3      | Method 3: Using VS Code Remote SSH . . . . .                    | 4         |
| 3.3.1    | Step 1: Install VS Code Extension . . . . .                     | 4         |
| 3.3.2    | Step 2: Configure SSH . . . . .                                 | 4         |
| 3.3.3    | Step 3: Connect via VS Code . . . . .                           | 5         |
| <b>4</b> | <b>Running Interactive Jobs</b>                                 | <b>5</b>  |
| 4.1      | CPU Job for MPI Programs . . . . .                              | 5         |
| 4.2      | GPU Jobs for CUDA Programs . . . . .                            | 6         |
| <b>5</b> | <b>Setting Up Singularity for MPI Programs</b>                  | <b>7</b>  |
| 5.1      | One-Time Setup . . . . .  | 7         |
| 5.1.1    | Step 1: Export PATH and Navigate to Scratch Directory . . . . . | 7         |
| 5.1.2    | Step 2: Download Singularity Image . . . . .                    | 7         |
| 5.1.3    | Step 3: Download and Extract Overlay File . . . . .             | 8         |
| 5.1.4    | Step 4: Launch Singularity Container . . . . .                  | 8         |
| 5.1.5    | Step 5: Install Miniconda (One-Time Only) . . . . .             | 8         |
| 5.1.6    | Step 6: Install OpenMPI via Conda (One-Time Only) . . . . .     | 8         |
| 5.2      | Regular Usage After Setup . . . . .                             | 8         |
| <b>6</b> | <b>Running MPI Programs</b>                                     | <b>9</b>  |
| 6.1      | Complete MPI Workflow . . . . .                                 | 9         |
| 6.2      | Important Notes for MPI . . . . .                               | 9         |
| <b>7</b> | <b>Running CUDA Programs</b>                                    | <b>10</b> |
| 7.1      | Basic CUDA Workflow . . . . .                                   | 10        |
| 7.2      | Available CUDA Modules . . . . .                                | 10        |

|           |   |           |
|-----------|---|-----------|
| <b>8</b>  | <b>Example Programs</b>   | <b>10</b> |
| 8.1       | Example 1: Basic MPI Hello World (WITH Singularity) . . . . .             | 10        |
| 8.2       | Example 2: CUDA Vector Addition . . . . .                                 | 11        |
| 8.3       | Example 4: Complete Singularity Setup from Scratch . . . . .              | 13        |
| <b>9</b>  | <b>Batch Job Scripts</b>  | <b>14</b> |
| 9.1       | MPI Batch Job . . . . .   | 14        |
| 9.2       | CUDA Batch Job . . . . .  | 14        |
| <b>10</b> | <b>Troubleshooting</b>  | <b>14</b> |
| 10.1      | Common MPI Errors and Solutions . . . . .                                 | 14        |
| 10.1.1    | Error: libpmix.so.2: cannot open shared object file . . . . .             | 14        |
| 10.1.2    | FATAL: while loading overlay images: failed to open overlay image . . . . | 15        |
| 10.1.3    | Error: Not enough slots available . . . . .                               | 15        |
| 10.1.4    | Error: File not found when copying Singularity image . . . . .            | 15        |
| 10.1.5    | Error: command not found: conda . . . . .                                 | 15        |
| 10.2      | For CUDA Programs . . . . .   | 15        |
| 10.2.1    | Error: nvcc: command not found . . . . .                                  | 15        |
| 10.2.2    | Error: No CUDA-capable device . . . . .                                   | 16        |
| 10.3      | SSH Host Key Issues . . . . .   | 16        |
| <b>11</b> | <b>Debugging Tips</b>   | <b>16</b> |
| 11.1      | For MPI Programs . . . . .  | 16        |
| 11.2      | For CUDA Programs . . . . .   | 16        |
| 11.3      | General Debugging . . . . .   | 17        |
| <b>12</b> | <b>Time and Resource Management</b>                                       | <b>17</b> |
| 12.1      | GPU Hour Budget . . . . .   | 17        |
| 12.2      | Best Practices . . . . .  | 17        |
| <b>13</b> | <b>Quick Reference</b>  | <b>17</b> |
| 13.1      | Essential Commands . . . . .  | 17        |
| 13.2      | MPI Workflow Summary . . . . .  | 18        |
| 13.3      | CUDA Workflow Summary . . . . .   | 18        |
| <b>14</b> | <b>Summary</b>  | <b>18</b> |
| 14.1      | Key Takeaways . . . . .   | 18        |
| 14.2      | Recommended Workflow . . . . .  | 19        |
| 14.3      | Important Links . . . . .   | 19        |

# 1 Your Resource Allocation

Each student has been assigned:

- **Slurm Account:** csci\_ua\_0480\_051-2025fa
- **GPU Hours:** 300 hours (18,000 minutes)
- **CPU Time:** Sufficient for coursework

## 1.1 Available Partitions

- **interactive** - For CPU-only interactive jobs
- **n2c48m24** - CPU partition
- **g2-standard-12, g2-standard-24, g2-standard-48** - GPU partitions with L4 GPUs
- **c12m85-a100-1, c24m170-a100-2** - A100 GPU partitions
- **n1s8-t4-1** - T4 GPU partition

# 2 Prerequisites

## 2.1 Network Connection

### VPN Requirement

You need to connect to NYU VPN **only if you are NOT on NYU campus WiFi**. If you are on campus, you can skip the VPN setup.

**For off-campus access:** Download NYU VPN from <https://www.nyu.edu/life/information-technology/infrastructure/network-services/vpn.html>

# 3 Quick Start Guide

## 3.1 Method 1: Using Open OnDemand (Recommended)

This is the easiest way to get started and manage files - no complex command line setup required!

### 3.1.1 Step 1: Access Open OnDemand

- Connect to NYU VPN (if off-campus) or NYU WiFi (if on campus)
- Go to: <https://ood-burst-001.hpc.nyu.edu/>
- Log in with your NYU credentials (NetID and password)

### What You Can Do with Open OnDemand

From the OOD server, you can:

- Launch compute nodes without logging into Greene cluster
- Run Jupyter notebooks (only if needed)
- Open terminal sessions directly in your browser
- Transfer files between your local computer and HPC
- Manage files with a graphical interface

### 3.1.2 Step 2: Upload Your Files

1. From the Open OnDemand homepage, click **Files** → **Home Directory**
2. Click the **Upload** button

3. Select your MPI or CUDA source files from your computer
4. Your files will be uploaded to your home directory on HPC

### 3.1.3 Step 3: Open a Terminal

1. Click **Clusters** → **Greene Shell Access**
2. This opens a terminal in your browser
3. Run `ssh burst` to connect to Cloud Burst
4. Now you can compile and run your programs

## 3.2 Method 2: Using Terminal/SSH

For users comfortable with command line access.

### 3.2.1 SSH Connection

#### Step 1: SSH to Greene Jump Host

```
1 ssh netID@gw.hpc.nyu.edu
```

Replace `netID` with your actual NYU NetID. For example, if your NetID is `ab1234`:

```
1 ssh ab1234@gw.hpc.nyu.edu
```

#### Step 2: SSH to Greene Cluster

```
1 ssh netID@greene.hpc.nyu.edu
```

#### Step 3: SSH to Cloud Burst

```
1 ssh burst
```

#### Note

When connecting for the first time, you may see: “This key is not known by any other names. Are you sure you want to continue connecting?”

- Type `yes` and press Enter

## 3.3 Method 3: Using VS Code Remote SSH

VS Code provides an excellent development environment with syntax highlighting, IntelliSense, and integrated terminal.

### 3.3.1 Step 1: Install VS Code Extension

1. Download VS Code from <https://code.visualstudio.com/>
2. Install the “Remote - SSH” extension

### 3.3.2 Step 2: Configure SSH

Click on the icon on the bottom left corner: follow fig 1 then fig 2

Then add this to your SSH config file (`~/.ssh/config`):

```
1 # NYU HPC Gateway
2 Host greene.hpc.nyu.edu
3   HostName greene.hpc.nyu.edu
4   User <net_id>
```

```

5
6 Host greene.hpc.nyu.edu dtn.hpc.nyu.edu gw.hpc.nyu.edu
7   StrictHostKeyChecking no
8   ServerAliveInterval 60
9   ForwardAgent yes
10  StrictHostKeyChecking no
11  UserKnownHostsFile /dev/null
12  LogLevel ERROR
13
14 Host hpcgwtunnel
15   HostName gw.hpc.nyu.edu
16   ForwardX11 no
17   StrictHostKeyChecking no
18   LocalForward 8027 greene.hpc.nyu.edu:22
19   UserKnownHostsFile /dev/null
20   User <net_id>
21
22 Host greene
23   HostName localhost
24   Port 8027
25   ForwardX11 yes
26   StrictHostKeyChecking no
27   UserKnownHostsFile /dev/null
28   LogLevel ERROR
29   User <net_id>
30
31 Host burst.hpc.nyu.edu
32   HostName burst.hpc.nyu.edu
33   User <netid>
34   ServerAliveInterval 60
35   ForwardAgent yes

```

Replace <net\_id> with your actual NetID.

### 3.3.3 Step 3: Connect via VS Code

1. Press F1 or Ctrl+Shift+P
2. Type “Remote-SSH: Connect to Host”
3. Select `greene.hpc.nyu.edu`
4. Enter your NYU password
5. Open folder: `/scratch/your_netid`
6. Start new terminal and run `ssh burst.` look at fig 3 for reference.
7. Edit, compile, and run code directly in VS Code

## 4 Running Interactive Jobs

### 4.1 CPU Job for MPI Programs

#### IMPORTANT: Specify CPUs

When running MPI programs, you **MUST** specify the number of CPUs you need using `--cpus-per-task=N` or `--ntasks=N`, where N is the number of MPI processes you plan to run. Otherwise, you’ll only get 1 CPU and cannot run parallel programs.

**For 8 MPI processes:**

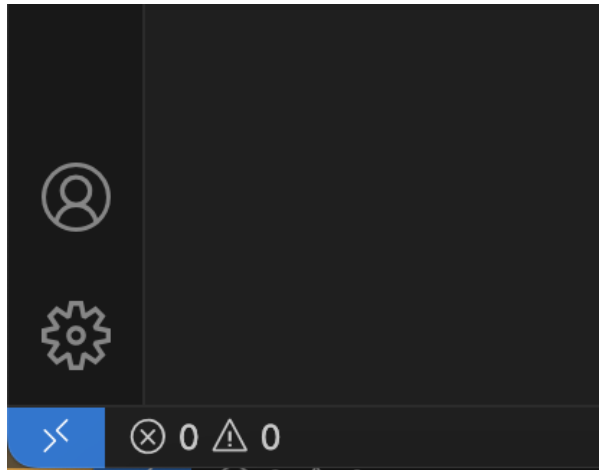


Figure 1: Click this icon to open a dialogue box

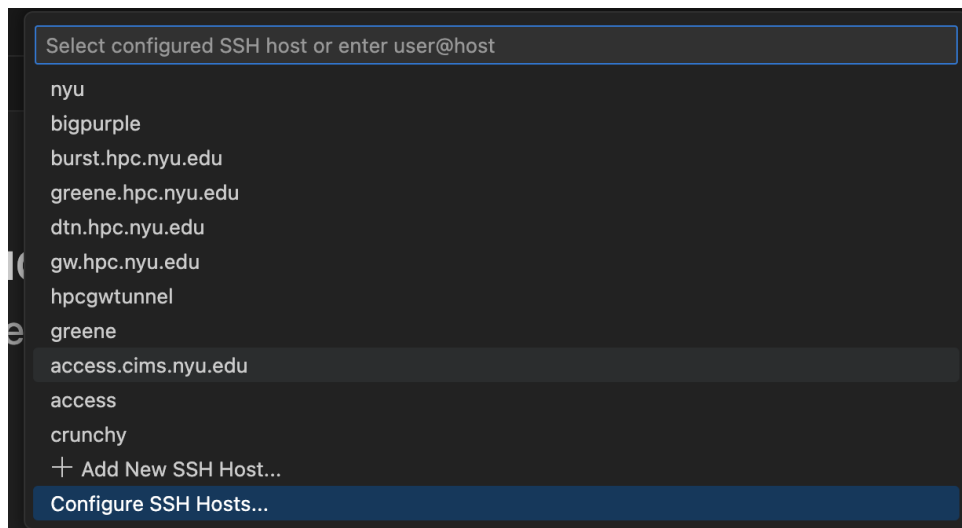


Figure 2: click on configure SSH Hosts (only for first time setup)

```
1 srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --cpus-per-task=8 --time
  =04:00:00 --pty /bin/bash
```

**Alternative with ntasks:**

```
1 srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --ntasks=2 --time=04:00:00
  --pty /bin/bash
```

## 4.2 GPU Jobs for CUDA Programs

### 1 L4 GPU for 4 Hours (Recommended):

```
1 srun --account=csci_ua_0480_051-2025fa --partition=g2-standard-12 --gres=gpu:1 --time
  =04:00:00 --pty /bin/bash
```

### 1 A100 GPU for 4 Hours:

```
1 srun --account=csci_ua_0480_051-2025fa --partition=c12m85-a100-1 --gres=gpu --time
  =04:00:00 --pty /bin/bash
```

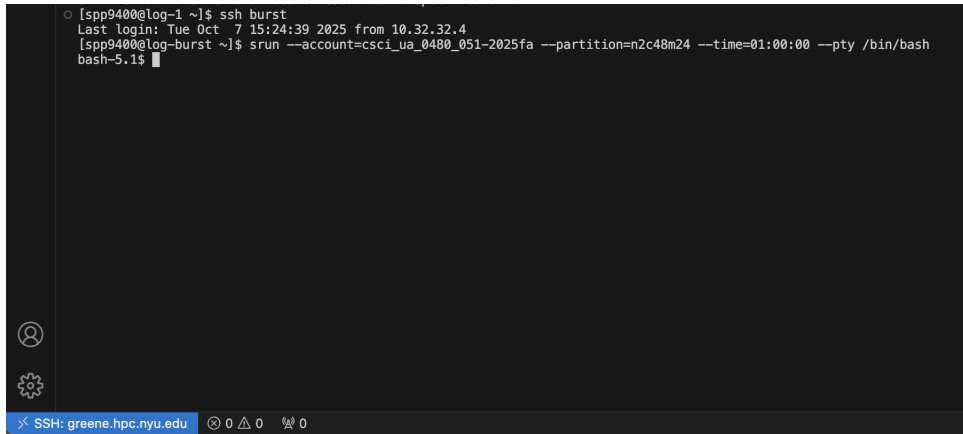


Figure 3: This is how it will look like

## 5 Setting Up Singularity for MPI Programs

### REQUIRED for MPI Programs

**Singularity is REQUIRED for MPI programs on burst nodes!** The burst compute nodes have an incomplete MPI installation that causes `libpmix.so.2` errors. You must use Singularity containers to run MPI programs successfully.

### 5.1 One-Time Setup

You only need to do this setup once. After that, you can reuse the same Singularity container and overlay.

#### 5.1.1 Step 1: Export PATH and Navigate to Scratch Directory

```

1 # Check if OpenMPI packages are available via the system package manager
2 which mpicc
3 ls /usr/lib64/openmpi/bin/
4 # If it exists in /usr/lib64/openmpi/bin/, you can add it to your PATH:
5 export PATH=/usr/lib64/openmpi/bin:$PATH
6 export LD_LIBRARY_PATH=/usr/lib64/openmpi/lib:$LD_LIBRARY_PATH
7
8 # Then try:
9 which mpicc
10 which mpirun
11 # After getting your compute node, go to scratch
12 cd /scratch/<your_netid> or cd ../../scratch/<your_netid>
13 # Example: cd /scratch/ab1234, based on your location: use pwd to check

```

#### 5.1.2 Step 2: Download Singularity Image

```

1 # Download the Ubuntu Singularity image
2 scp -rp greene-dtn:/scratch/work/public/singularity/ubuntu-20.04.3.sif .

```

This will take a few minutes as the file is several GB.

### 5.1.3 Step 3: Download and Extract Overlay File

```
1 # Download overlay template
2 scp -rp greene-dtn:/scratch/work/public/overlay-fs-ext3/overlay-15GB-500K.ext3.gz .
3
4 # Extract it
5 gunzip overlay-15GB-500K.ext3.gz
```

This will take a few minutes as the file is several GB.

### 5.1.4 Step 4: Launch Singularity Container

```
1 # Start the container with overlay
2 singularity exec --overlay overlay-15GB-500K.ext3:rw ubuntu-20.04.3.sif /bin/bash
```

Your prompt will change to show you're inside Singularity (e.g., Singularity>).

### 5.1.5 Step 5: Install Miniconda (One-Time Only)

```
1 # Inside Singularity container
2 # Download Miniconda
3 wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
4
5 # Install to /ext3 (the overlay filesystem)
6 bash Miniconda3-latest-Linux-x86_64.sh -b -p /ext3/miniconda3
7
8 # Activate conda
9 source /ext3/miniconda3/bin/activate
```

### 5.1.6 Step 6: Install OpenMPI via Conda (One-Time Only)

```
1 # Still inside Singularity, with conda activated
2 # accepting the terms and conditions the first time
3 conda tos accept --override-channels --channel https://repo.anaconda.com/pkgs/main
4 conda tos accept --override-channels --channel https://repo.anaconda.com/pkgs/r
5 conda install gcc_linux-64 gxx_linux-64 openmpi -y
```

This installs a working MPI implementation that doesn't have the libpmix issues.

## 5.2 Regular Usage After Setup

After the one-time setup, your workflow for each session is:

```
1 # 1. Get compute node with enough CPUs
2 srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --cpus-per-task=8 --time
   =04:00:00 --pty /bin/bash
3
4 # 2. Go to your scratch directory
5 cd /scratch/<your_netid>
6 # or
7 cd ../../scratch/<your_netid> if you are in home directory
8
9 # 3. Launch Singularity (note: using :ro for read-only)
10 singularity exec --overlay overlay-15GB-500K.ext3:ro ubuntu-20.04.3.sif /bin/bash
11
12 # 4. Activate conda
```



```

13 source /ext3/miniconda3/bin/activate
14
15 # 5. Now you can compile and run MPI programs!
16 mpicc -o hello_mpi hello_mpi.c
17 mpirun -np 2 ./hello_mpi

```

### Pro Tip

Use `:ro` (read-only) instead of `:rw` when you don't need to install new software. This is safer and prevents accidental modifications.

## 6 Running MPI Programs

### 6.1 Complete MPI Workflow

```

1 # 1. SSH to burst
2 ssh <netid>@gw.hpc.nyu.edu
3 ssh <netid>@greene.hpc.nyu.edu
4 ssh burst
5
6 # 2. Request compute node with CPUs matching your MPI process count
7 srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --cpus-per-task=8 --time
   =04:00:00 --pty /bin/bash
8
9 # 3. Go to scratch and launch Singularity
10 cd /scratch/<netid>
11 singularity exec --overlay overlay-15GB-500K.ext3:ro ubuntu-20.04.3.sif /bin/bash
12
13 # 4. Activate conda
14 source /ext3/miniconda3/bin/activate
15
16 # 5. Create your MPI program
17 vi hello_mpi.c
18
19 # 6. Compile with mpicc
20 mpicc -o hello_mpi hello_mpi.c
21
22 # 7. Run with mpirun
23 mpirun -np 2 ./hello_mpi
24
25 # 8. When done, exit
26 exit # Exit Singularity
27 exit # Exit compute node

```

### 6.2 Important Notes for MPI

- **Match CPUs to processes:** If you run `mpirun -np 8`, you need `--cpus-per-task=8` or more
- **Always use Singularity:** Running `mpicc/mpirun` directly on burst nodes will fail with `libpmix` errors
- **Activate conda:** Don't forget `source /ext3/miniconda3/bin/activate` inside Singularity
- **Work in /scratch:** Store your code in `/scratch/<netid>` for better performance

## 7 Running CUDA Programs

### 7.1 Basic CUDA Workflow

```
1 # 1. Request GPU node
2 srun --account=csci_ua_0480_051-2025fa --partition=g2-standard-12 --gres=gpu:1 --time
   =04:00:00 --pty /bin/bash
3
4 # 2. Load CUDA module
5 module load cuda/11.8.0
6
7 # 3. Verify GPU
8 nvidia-smi
9
10 # 4. Create your CUDA program
11 vi hello_cuda.cu
12
13 # 5. Compile
14 nvcc -o hello_cuda hello_cuda.cu
15
16 # 6. Run
17 ./hello_cuda
```

### 7.2 Available CUDA Modules

Check available CUDA versions:

```
1 module avail cuda
```

Common versions:

- cuda/11.8.0 - Recommended for most programs
- cuda/12.0 - For newer GPU features

## 8 Example Programs

### 8.1 Example 1: Basic MPI Hello World (WITH Singularity)

This is the complete, corrected workflow that actually works on burst nodes.

**Step 1: Create the MPI program**

```
1 // hello_mpi.c
2 #include <mpi.h>
3 #include <stdio.h>
4
5 int main(int argc, char** argv) {
6     MPI_Init(&argc, &argv);
7
8     int world_size, world_rank;
9     MPI_Comm_size(MPI_COMM_WORLD, &world_size);
10    MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
11
12    char processor_name[MPI_MAX_PROCESSOR_NAME];
13    int name_len;
14    MPI_Get_processor_name(processor_name, &name_len);
15
16    printf("Hello from processor %s, rank %d out of %d\n",
```

```

17     processor_name, world_rank, world_size);
18
19     MPI_Finalize();
20     return 0;
21 }

```

## Step 2: Complete execution workflow

```

1  # Connect to burst
2  ssh ab1234@gw.hpc.nyu.edu
3  ssh ab1234@greene.hpc.nyu.edu
4  ssh burst
5
6  # Get compute node with 8 CPUs
7  srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --cpus-per-task=8 --time
   =04:00:00 --pty /bin/bash
8
9  # Navigate to scratch
10 cd /scratch/ab1234
11
12 # Launch Singularity (assuming you did the one-time setup)
13 singularity exec --overlay overlay-15GB-500K.ext3:ro ubuntu-20.04.3.sif /bin/bash
14
15 # Activate conda
16 source /ext3/miniconda3/bin/activate
17
18 # Create the program
19 vi hello_mpi.c
20 # (paste the code above)
21
22 # Compile
23 mpicc -o hello_mpi hello_mpi.c
24
25 # Run with 2 processes
26 mpirun -np 2 ./hello_mpi
27
28 # Expected output:
29 # Hello from processor b-9-77, rank 0 out of 8
30 # Hello from processor b-9-77, rank 1 out of 8
31 # Hello from processor b-9-77, rank 2 out of 8
32 # ... (8 lines total)
33
34 # Exit when done
35 exit # Exit Singularity
36 exit # Exit compute node

```

## 8.2 Example 2: CUDA Vector Addition

```

1 // vector_add.cu
2 #include <stdio.h>
3 #include <cuda_runtime.h>
4
5 __global__ void vectorAdd(float *a, float *b, float *c, int n) {
6     int i = blockDim.x * blockIdx.x + threadIdx.x;
7     if (i < n) {
8         c[i] = a[i] + b[i];
9     }

```

```

10 }
11
12 int main() {
13     int n = 1000;
14     size_t size = n * sizeof(float);
15
16     // Allocate host memory
17     float *h_a = (float*)malloc(size);
18     float *h_b = (float*)malloc(size);
19     float *h_c = (float*)malloc(size);
20
21     // Initialize vectors
22     for (int i = 0; i < n; i++) {
23         h_a[i] = i;
24         h_b[i] = i * 2;
25     }
26
27     // Allocate device memory
28     float *d_a, *d_b, *d_c;
29     cudaMalloc(&d_a, size);
30     cudaMalloc(&d_b, size);
31     cudaMalloc(&d_c, size);
32
33     // Copy to device
34     cudaMemcpy(d_a, h_a, size, cudaMemcpyHostToDevice);
35     cudaMemcpy(d_b, h_b, size, cudaMemcpyHostToDevice);
36
37     // Launch kernel
38     int threadsPerBlock = 256;
39     int blocksPerGrid = (n + threadsPerBlock - 1) / threadsPerBlock;
40     vectorAdd<<<blocksPerGrid, threadsPerBlock>>>(d_a, d_b, d_c, n);
41
42     // Copy result back
43     cudaMemcpy(h_c, d_c, size, cudaMemcpyDeviceToHost);
44
45     // Verify result
46     printf("First 5 results:\n");
47     for (int i = 0; i < 5; i++) {
48         printf("%f + %f = %f\n", h_a[i], h_b[i], h_c[i]);
49     }
50
51     // Cleanup
52     cudaFree(d_a);
53     cudaFree(d_b);
54     cudaFree(d_c);
55     free(h_a);
56     free(h_b);
57     free(h_c);
58
59     return 0;
60 }

```

**To run:**

```

1 # Get GPU node (NO Singularity needed for CUDA)
2 srun --account=csci_ua_0480_051-2025fa --partition=g2-standard-12 --gres=gpu:1 --time
   =02:00:00 --pty /bin/bash
3
4 # Load CUDA

```

```

5 module load cuda/11.8.0
6
7 # Compile and run
8 nvcc -o vector_add vector_add.cu
9 ./vector_add

```

### 8.3 Example 4: Complete Singularity Setup from Scratch

This shows the complete one-time setup process:

```

1 # 1. Connect to burst and get compute node
2 ssh ab1234@gw.hpc.nyu.edu
3 ssh ab1234@greene.hpc.nyu.edu
4 ssh burst
5 srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --cpus-per-task=8 --time
   =02:00:00 --pty /bin/bash
6
7 # 2. Go to scratch directory
8 cd /scratch/ab1234
9
10 # 3. Download Singularity image
11 scp -rp greene-dtn:/scratch/work/public/singularity/ubuntu-20.04.3.sif .
12
13 # 4. Download and extract overlay
14 scp -rp greene-dtn:/scratch/work/public/overlay-fs-ext3/overlay-15GB-500K.ext3.gz .
15 gunzip overlay-15GB-500K.ext3.gz
16
17 # 5. Launch Singularity with read-write overlay
18 singularity exec --overlay overlay-15GB-500K.ext3:rw ubuntu-20.04.3.sif /bin/bash
19
20 # 6. Inside Singularity: Install Miniconda
21 wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
22 bash Miniconda3-latest-Linux-x86_64.sh -b -p /ext3/miniconda3
23
24 # 7. Activate conda
25 source /ext3/miniconda3/bin/activate
26
27 # 8. Install OpenMPI
28 conda install -c conda-forge openmpi -y
29
30 # 9. Test it works
31 echo '#include <mpi.h>'
32 #include <stdio.h>
33 int main(int argc, char** argv) {
34     MPI_Init(&argc, &argv);
35     int rank;
36     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
37     printf("Hello from rank %d\n", rank);
38     MPI_Finalize();
39     return 0;
40 }' > test_mpi.c
41
42 mpicc -o test_mpi test_mpi.c
43 mpirun -np 4 ./test_mpi
44
45 # 10. If you see 4 lines of output, setup is complete!
46 # Exit and use :ro for future sessions
47 exit

```

---

## 9 Batch Job Scripts

### 9.1 MPI Batch Job

```
1 #!/bin/bash
2 #SBATCH --job-name=mpi_job
3 #SBATCH --account=csci_ua_0480_051-2025fa
4 #SBATCH --partition=n2c48m24
5 #SBATCH --nodes=1
6 #SBATCH --ntasks-per-node=32
7 #SBATCH --time=01:00:00
8 #SBATCH --output=%j_%x.out
9 #SBATCH --error=%j_%x.err
10 #SBATCH --requeue
11
12 # Navigate to scratch
13 cd /scratch/$USER
14
15 # Launch Singularity with conda and run program
16 singularity exec --overlay overlay-15GB-500K.ext3:ro \
17     ubuntu-20.04.3.sif /bin/bash << 'EOF'
18 source /ext3/miniconda3/bin/activate
19 mpirun -np 32 ./my_mpi_program
20 EOF
```

### 9.2 CUDA Batch Job

```
1 #!/bin/bash
2 #SBATCH --job-name=cuda_job
3 #SBATCH --account=csci_ua_0480_051-2025fa
4 #SBATCH --partition=g2-standard-12
5 #SBATCH --gres=gpu:1
6 #SBATCH --time=01:00:00
7 #SBATCH --output=%j_%x.out
8 #SBATCH --error=%j_%x.err
9 #SBATCH --requeue
10
11 module load cuda/11.8.0
12
13 echo "GPU Information:"
14 nvidia-smi
15 echo ""
16
17 echo "Running CUDA program..."
18 ./my_cuda_program
```

## 10 Troubleshooting

### 10.1 Common MPI Errors and Solutions

#### 10.1.1 Error: libpmix.so.2: cannot open shared object file

**Problem:** You tried to run MPI programs directly without Singularity.

**Solution:** Always use Singularity with conda for MPI programs:

```
1 singularity exec --overlay overlay-15GB-500K.ext3:ro ubuntu-20.04.3.sif /bin/bash
2 source /ext3/miniconda3/bin/activate
3 mpirun -np 8 ./program
```

### 10.1.2 FATAL: while loading overlay images: failed to open overlay image

**Problem:** The overlay file might have a stale lock from a previous session that didn't exit cleanly.

**Solution:** Reboot Your Session:

```
1 mv overlay-15GB-500K.ext3 overlay-15GB-500K.ext3.backup
2 scp -rp greene-dtn:/scratch/work/public/overlay-fs-ext3/overlay-15GB-500K.ext3.gz .
3 gunzip overlay-15GB-500K.ext3.gz
4 singularity exec --overlay overlay-15GB-500K.ext3:rw ubuntu-20.04.3.sif /bin/bash
```

### 10.1.3 Error: Not enough slots available

**Problem:** You didn't request enough CPUs when launching your job.

**Solution:** Match CPU count to MPI processes:

```
1 # If running mpirun -np 8, need at least 8 CPUs:
2 srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --cpus-per-task=8 --time
   =04:00:00 --pty /bin/bash
```

### 10.1.4 Error: File not found when copying Singularity image

**Problem:** Wrong file path.

**Solution:** Use the correct path:

```
1 # CORRECT path:
2 scp -rp greene-dtn:/scratch/work/public/singularity/ubuntu-20.04.3.sif .
3
4 # NOT this:
5 # scp -rp greene-dtn:/share/apps/images/ubuntu-20.04.3.sif .
```

### 10.1.5 Error: command not found: conda

**Problem:** Forgot to activate conda inside Singularity.

**Solution:** Always run after entering Singularity:

```
1 source /ext3/miniconda3/bin/activate
```

## 10.2 For CUDA Programs

### 10.2.1 Error: nvcc: command not found

**Solution:** Load CUDA module first:

```
1 module load cuda/11.8.0
```

All "module load" commands are for greene and not burst.

## 10.2.2 Error: No CUDA-capable device

**Problem:** Not on a GPU node.

**Solution:** Request GPU node:

```
1 srun --account=csci_ua_0480_051-2025fa --partition=g2-standard-12 --gres=gpu:1 --time  
   =04:00:00 --pty /bin/bash
```

## 10.3 SSH Host Key Issues

If you see “Host key verification failed” when copying from greene-dtn:

```
1 # Remove old host key  
2 ssh-keygen -R greene-dtn  
3 ssh-keygen -R greene-dtn.hpc.nyu.edu  
4  
5 # Then retry your scp command
```

# 11 Debugging Tips

## 11.1 For MPI Programs

- **Start with 1 process:** `mpirun -np 1 ./program` to verify basic functionality
- **Add rank-specific output:** Help identify which process has issues
- **Use `printf` with `fflush`:** Ensures output appears immediately
- **Check MPI return codes:** Many MPI functions return error codes

Example debugging code:

```
1 printf("Rank %d: Before MPI_Send\n", rank);  
2 fflush(stdout);  
3 int result = MPI_Send(...);  
4 if (result != MPI_SUCCESS) {  
5     printf("Rank %d: MPI_Send failed with code %d\n", rank, result);  
6 }  
7 printf("Rank %d: After MPI_Send\n", rank);  
8 fflush(stdout);
```

## 11.2 For CUDA Programs

- **Check every CUDA call:** Don’t skip error checking
- **Use `cuda-memcheck`:** Detects memory errors
- **Start with CPU version:** Verify algorithm before GPU
- **Use small test cases:** Easier to debug

Example error checking:

```
1 cudaError_t err = cudaMalloc(&d_array, size);  
2 if (err != cudaSuccess) {  
3     printf("CUDA Error: %s\n", cudaGetErrorString(err));  
4     exit(1);  
5 }  
6  
7 // After kernel launch  
8 err = cudaGetLastError();  
9 if (err != cudaSuccess) {
```



```

10 printf("Kernel launch error: %s\n", cudaGetErrorString(err));
11 }

```

## 11.3 General Debugging

- **Compile with debug symbols:** `-g` flag
- **Check you're in the right environment:** Inside Singularity for MPI, GPU node for CUDA
- **Verify file locations:** Use `ls` and `pwd`
- **Check resource allocation:** Use `squeue -u $USER`

# 12 Time and Resource Management

## 12.1 GPU Hour Budget

With 300 GPU hours for the semester:

- **Development:** 15-30 min sessions for debugging
- **Testing:** 1 hour sessions for testing
- **Production:** 4+ hour sessions for final runs
- **Monitor usage:** Check regularly with `sacct`

Check your usage:

```

1 # View recent jobs
2 sacct -u $USER --format=JobID,JobName,Partition,Elapsed,State
3
4 # View all jobs since semester start
5 sacct -u $USER --format=JobID,JobName,Elapsed,State,AllocGRES -S 2025-01-01

```

## 12.2 Best Practices

- **Test interactively first:** Debug before submitting batch jobs
- **Use appropriate partitions:** CPU for MPI, GPU for CUDA
- **Request appropriate time:** Don't request 4 hours if you only need 30 minutes
- **Use `--requeue`:** Allows jobs to restart if nodes shut down
- **Save checkpoints:** For long-running jobs

# 13 Quick Reference

## 13.1 Essential Commands

```

1 # SSH to burst
2 ssh <netid>@gw.hpc.nyu.edu
3 ssh <netid>@greene.hpc.nyu.edu
4 ssh burst
5
6 # Interactive MPI job (with Singularity)
7 srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --cpus-per-task=8 --time
  =04:00:00 --pty /bin/bash
8 cd /scratch/$USER
9 singularity exec --overlay overlay-15GB-500K.ext3:ro ubuntu-20.04.3.sif /bin/bash
10 source /ext3/miniconda3/bin/activate
11

```

```

12 # Interactive GPU job
13 srun --account=csci_ua_0480_051-2025fa --partition=g2-standard-12 --gres=gpu:1 --time
    =04:00:00 --pty /bin/bash
14 module load cuda/11.8.0
15
16 # Submit batch job
17 sbatch job_script.sh
18
19 # Check job status
20 squeue -u $USER
21
22 # Cancel job
23 scancel <job_id>
24
25 # View job history
26 sacct -u $USER

```

## 13.2 MPI Workflow Summary

```

1 # 1. Get compute node with CPUs
2 srun --account=csci_ua_0480_051-2025fa --partition=n2c48m24 --cpus-per-task=8 --time
    =04:00:00 --pty /bin/bash
3
4 # 2. Navigate to scratch
5 cd /scratch/<netid>
6
7 # 3. Launch Singularity
8 singularity exec --overlay overlay-15GB-500K.ext3:ro ubuntu-20.04.3.sif /bin/bash
9
10 # 4. Activate conda
11 source /ext3/miniconda3/bin/activate
12
13 # 5. Compile and run
14 mpicc -o program program.c
15 mpirun -np 2 ./program

```

## 13.3 CUDA Workflow Summary

```

1 # 1. Get GPU node
2 srun --account=csci_ua_0480_051-2025fa --partition=g2-standard-12 --gres=gpu:1 --time
    =04:00:00 --pty /bin/bash
3
4 # 2. Load CUDA
5 module load cuda/11.8.0
6
7 # 3. Compile and run
8 nvcc -o program program.c
9 ./program

```

# 14 Summary

## 14.1 Key Takeaways

- **Singularity is REQUIRED for MPI:** Burst nodes have broken MPI, use Singularity + conda

- **Singularity NOT needed for CUDA:** CUDA works directly with module load
- **Specify CPUs for MPI:** Use `--cpus-per-task=N` matching your MPI process count
- **Use correct file paths:** `/scratch/work/public/singularity/...`
- **Activate conda in Singularity:** `source /ext3/miniconda3/bin/activate`
- **Work in /scratch:** Better performance than `/home`
- **Open OnDemand is easiest:** <https://ood-burst-001.hpc.nyu.edu/>
- **Use `--requeue`:** In all batch scripts for spot instances
- **Monitor GPU hours:** 300 hours total for semester
- **Test interactively first:** Debug before batch submission

## 14.2 Recommended Workflow

1. Develop code locally or using VS Code Remote-SSH
2. Upload to HPC via Open OnDemand or SCP
3. For MPI: Set up Singularity once, then reuse
4. Test interactively with small inputs
5. Debug and verify correctness
6. Create batch script for production runs
7. Submit batch job with `--requeue`
8. Monitor with `squeue` and download results

## 14.3 Important Links

- **Open OnDemand:** <https://ood-burst-001.hpc.nyu.edu/>
- **NYU HPC Docs:** <https://sites.google.com/nyu.edu/nyu-hpc/>
- **Singularity Guide:** <https://sites.google.com/nyu.edu/nyu-hpc/hpc-systems/greene/software/singularity-with-miniconda>
- **HPC Support:** [hpc@nyu.edu](mailto:hpc@nyu.edu)

### Good Luck!

You now have everything you need to successfully use HPC Cloud Bursting! Remember:

- MPI requires Singularity (one-time setup, then easy to use)
- CUDA works directly without Singularity
- Always specify CPUs when running MPI programs
- Start early and test often
- Ask questions when stuck - check Brightspace or email [hpc@nyu.edu](mailto:hpc@nyu.edu)
- Have fun learning parallel computing!