

An Introduction to GPU Acceleration

Research Computing Summer Camp 2023

24 May

Layla Freeborn, Ph.D.

layla.freeborn@colorado.edu

Session Overview

1. CPUs vs GPUs
2. Heterogeneous Systems
3. Criteria for GPU Acceleration and Factors Affecting GPU Speedup
4. Alpine GPU Partitions and Requesting GPUs with Slurm
5. GPU Programming Tools
6. GPU Monitoring Tools
7. Self-guided hands-on tutorial

The background of the slide is a photograph of a majestic mountain range, likely the Dolomites, featuring sharp, light-colored rock formations against a backdrop of a bright blue sky with scattered white and grey clouds.

slido

Join at
slido.com
#AlpineGPU



CPUs vs GPUs

Mythbusters Demo- GPU vs CPU



Join at
slido.com
#AlpineGPU

How do you feel about GPU acceleration?

Score: 0.0

0%

0%

0%

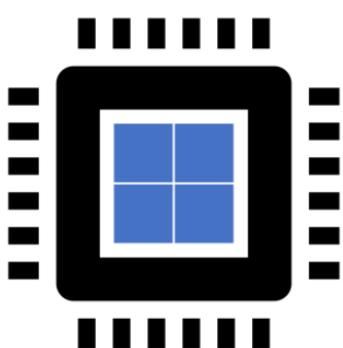
0%

0%



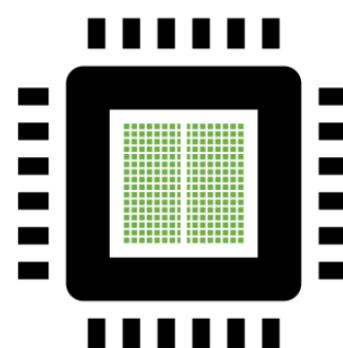
CPUs vs GPUs

Central Processing Unit



- General workhorse
- Contains **dozens** of cores
- Good for **serial processing** and handling multiple

Graphics Processing Unit



- Specialized
- Contains **thousands** of cores
- Good for **parallel processing** and handling specific tasks quickly

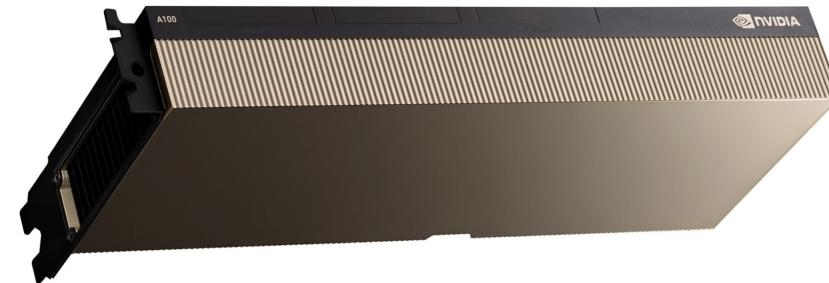
Heterogeneous Systems

MI100 GPU



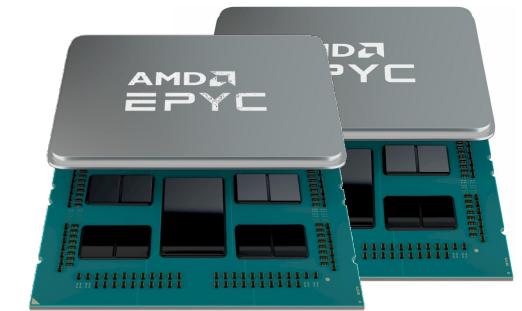
8 nodes with 3 GPUs per node

A100 GPU



12 nodes with 3 GPUs per node

Milan CPU



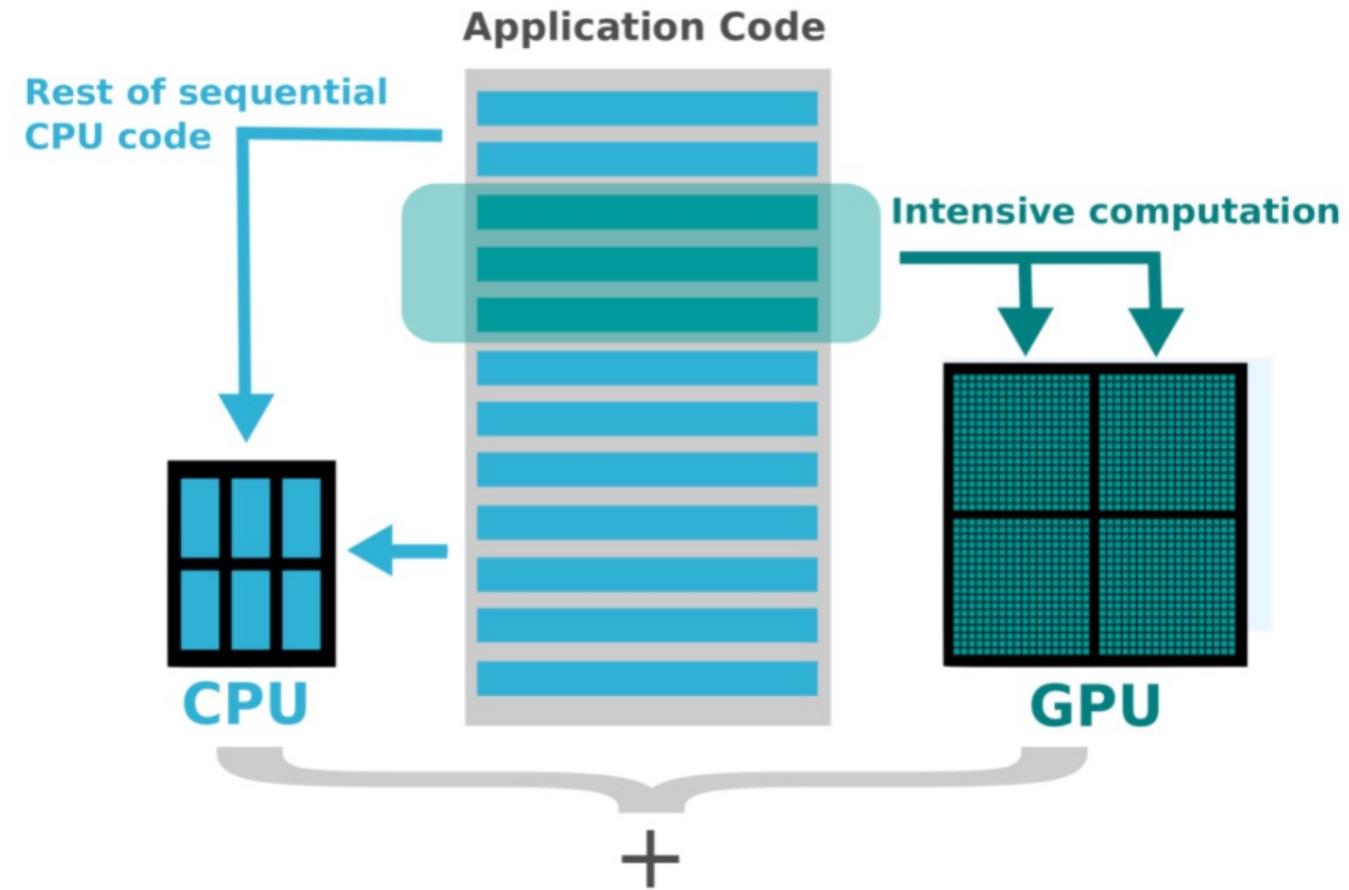
297 nodes

- 283 nodes with 256 GB RAM
- 14 nodes with 1 TB RAM

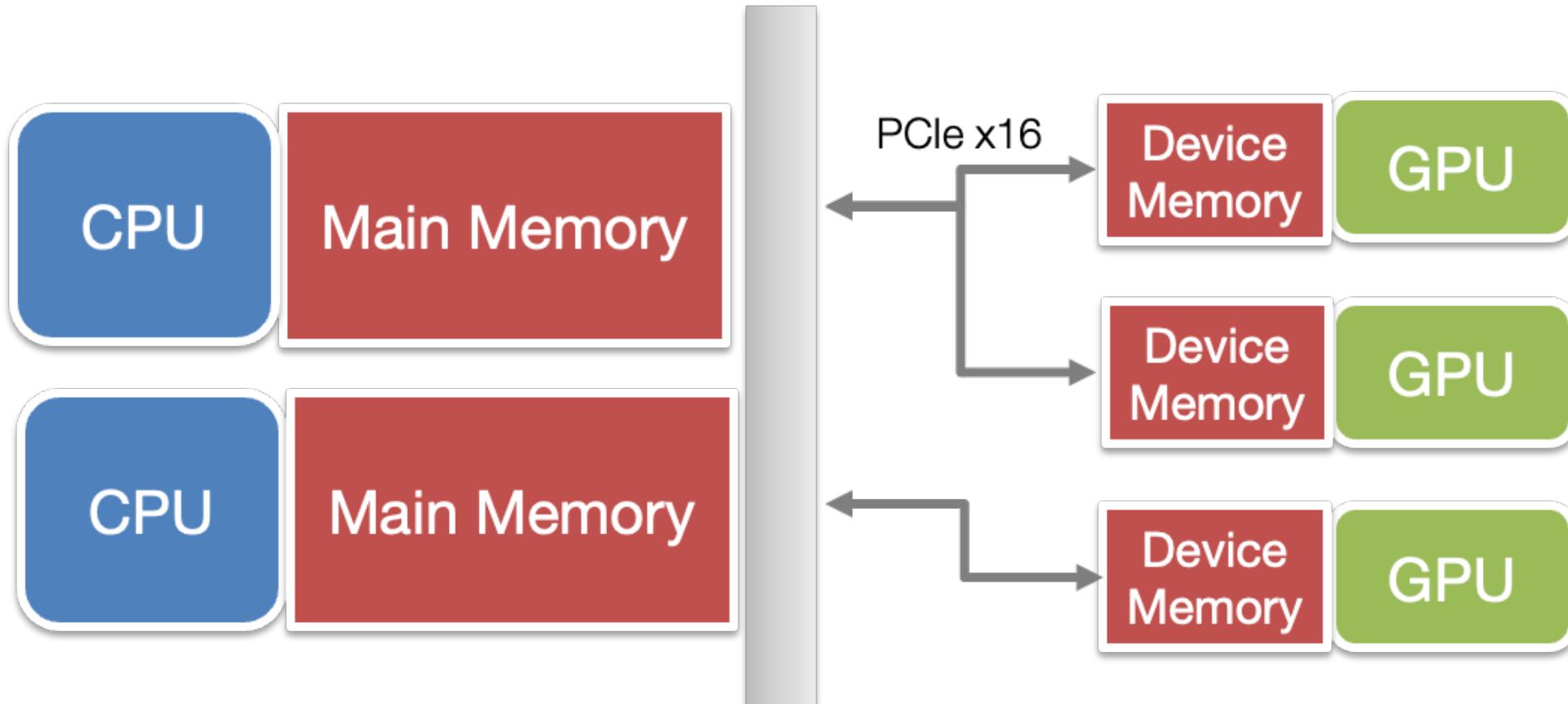
Heterogeneous Systems

GPU acceleration often follows an **offloading model**.

- Increased speed can be achieved by porting computationally intensive parts of code to GPU(s)



Heterogeneous Systems



Data needs to be copied from CPU to GPU, computation is performed on the GPU, then output is transferred back to CPU.

Criteria for GPU Acceleration

1. The time spent on computationally intensive parts of the workflow exceeds the time spent transferring data to and from GPU memory
2. Computations are massively parallel- the computations can be broken down into hundreds or thousands of independent units of work

Factors Affecting GPU Speedup

- 1** Computational Intensity
- 2** Data Dependency
- 3** Data Type
- 4** Code/Algorithmic Complexity

Factors Affecting GPU Speedup

1

Computational Intensity

GPUs perform best when there is a lot of processing compared to loading and storing data (FLOP per Byte ratio).

Factors Affecting GPU Speedup

2

Data Dependency

Data Dependency- A situation in which an instruction is dependent on a result from a sequentially previous instruction before it can complete its execution.

This should be avoided!

Factors Affecting GPU Speedup

3

Data Type

- Operations on strings are slow unless they can be treated as numbers.
- Performance per GPU can vary if workflows include 16-bit and 64-bit floats.

Factors Affecting GPU Speedup

4

Code/Algorithmic Complexity

Simple code is better ported to GPUs.

- Deeply-branched code and while-loops may perform poorly on GPUs
- Recursive functions need to be re-written

Factors Affecting GPU Speedup



The performance increases from GPUs depend on several factors related to the data and algorithm.

Alpine GPU Partitions

Try these commands.

```
1 ssh <username>@login.rc.colorado.edu
2 sinfo --Format Partition
3 sinfo --partition aa100,ami100 --Format Partition,Nodes,Time,
4 scontrol show partition aa100
5 scontrol show partition ami100
6 scontrol show node c3gpu-c2-u17
7 scontrol show node c3gpu-a9-u29-1
```

Requesting Alpine GPUs with Slurm

Slurm flags needed to request 1 AMD GPU node with 2 GPUs and 20 CPU cores

```
--partition=ami100  
--gres=gpu:2  
--ntasks=20
```

in a job script submitted with *sbatch* command

```
#SBATCH --partition=ami100  
#SBATCH --gres=gpu:2  
#SBATCH --ntasks=20  
#SBATCH --job-name=gpu_test  
#SBATCH --output=gpu_test_%j.out  
#SBATCH --error=gpu_test_%j.err  
#SBATCH --mail-type=ALL  
#SBATCH --mail-user=<email>
```

in an interactive job

```
sinteractive --partition=ami100 --gres=gpu:2 --ntasks=20
```

Requesting Alpine GPUs with Slurm

```
#request one AMD GPU with default time and default CPU cores  
sinteractive --partition=ami100 --gres=gpu:1
```

```
#request two AMD GPUs with 64 CPU cores for 24 hours  
sinteractive --partition=ami100 --gres=gpu:2 --ntasks=64 --time=24:00:00
```

```
#request one AMD GPU for 7 days (requires long QoS)  
sinteractive --partition=ami100 --gres=gpu:1 --time=7-00:00:00 --qos=long
```

```
#request three 80GB NVIDIA A100 GPUs with 20 CPU cores for default time  
sinteractive --partition=aa100 --ntasks=20 --gres=gpu:3 --constraint=gpu80
```



Pay attention to defaults!



Join at
slido.com
#AlpineGPU

How do you feel about GPU acceleration?

Score: 0.0

0%

0%

0%

0%

0%



GPU Programming Tools

GPU
Libraries

“Drop-in”
Acceleration

GPU
Directives

Easily
Accelerate
Applications

GPU
Languages

Maximum
Flexibility

GPU Programming Tools

Definitions

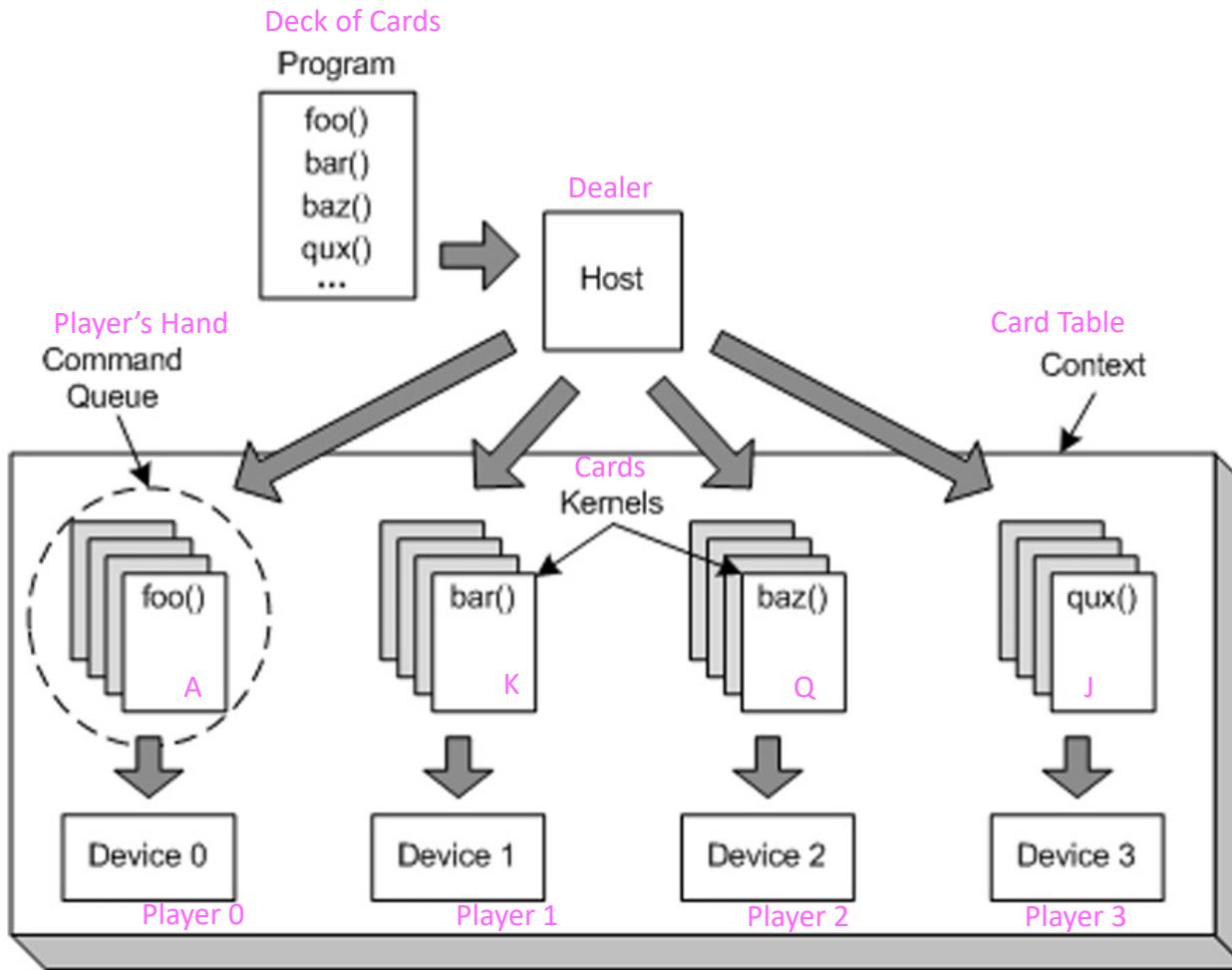
- **Host**- CPU
- **Device**- GPU
- **Kernel**- functions launched to the GPU

GPU Programming Tools

GPU Languages (Language Extensions)

- OpenCL (NVIDIA & AMD GPUs), HIP (NVIDIA & AMD), & CUDA (NVIDIA only)
- Require in-depth knowledge of hardware
- May involve substantial changes to code
- Code must be maintained to keep up with hardware changes & performance guidelines

GPU Programming Tools



Dealer distributes cards to players.
Host distributes kernels to devices.

Player receives cards from dealer.
Device receives kernels from host.

Dealer selects cards from a deck.
Host selects kernels from a program

Each player receives cards as part of
a hand.
Each device receives kernels through
the command queue.

Card table makes it possible for
players to transfer cards to each other.
OpenCL Context allows devices to
receive kernels and transfer data.

1	Obtain OpenCL platform	<code>clGetPlatformIDs(1, &platform, NULL)</code>
2	Obtain device id for at least one device (accelerator)	<code>clGetDeviceIds(platform, CL_DEVICE_TYPE_GPU, 1, &device, NULL)</code>
3	Create context for device	<code>context = clCreateContext(NULL, 1, &device, NULL, NULL, &err)</code>
4	Create accelerator program from source code	<code>program = clCreateProgramWithSource (context, 1, (const char**) &program_buffer, &program_size, &err)</code>
5	Build the program	<code>clBuildProgram(program, 0,..)</code>
6	Create kernel(s) from program functions	<code>kernel = clCreateKernel(program, "kernel_name", &err)</code>
7	Create command queue for target device	<code>queue = clCreateCommandQueue(context, device, 0, &err)</code>
8	Allocate device memory / move input data to device memory	<code>memObject = clCreateBuffer (context, NULL, SIZE_N, NULL, &err) clEnqueueWriteBuffer(command_queue, memObject, ..., TOTAL_SIZE, hostPointer, ...)</code>
9	Associate arguments to kernel with kernel object	<code>cl_int clSetKernelArg (kernel, arg_index, arg_size, *arg_value)</code>
10	Deploy kernel for device execution	<code>global_size = TOTAL_NUM_THREADS; local_size = WORKGROUP_SIZE; clEnqueueNDRangeKernel(command_queue, kernel, 1, NULL, &global_size, &local_size, 0, NULL, NULL);</code>
11	Move output data to host memory	<code>clEnqueueReadBuffer(command_queue, memObject, blocking_read, offset, TOTAL_SIZE, hostPointer, 0, NULL, NULL)</code>
12	Release context/program/kernels/memory	<code>clReleaseMemObject(memObject) / clReleaseKernel(kernel) / clReleaseProgram(program) / clReleaseContext(context)</code>

GPU Programming Tools

GPU Compiler Directives

- Easier to learn and involve fewer changes to code than GPU programming languages
- Better portability among devices and platforms
- Require little understanding of GPU hardware, but more than GPU libraries
- OpenACC is the most commonly used

GPU Programming Tools

OpenACC

GPU Compiler Directives

- A collection of compiler directives that specify loops and regions of code in standard C, C++, and Fortran to be offloaded from a host CPU to an attached parallel accelerator
developer.nvidia.com
- Lets compiler guess data allocation and movement or control it with directive clauses
- Can be used with GPU libraries, OpenMP

GPU Programming Tools

OpenACC

Basic program structure

```
#include "openacc.h"
#pragma acc <directive> [clauses [[,] clause]...] new-line
<code>
```

Kernel directives tell the compiler to generate parallel accelerator kernels for the loop nests following the directive.

Data directives tell the compiler to create code that performs specific data movements and provides hints about data usage.

Compile C code for NVIDIA GPU

```
pgcc -acc -ta=nvidia -c your_program_acc.c
```

Compile C++ code for NVIDIA GPU

```
nvcc --acc -Minfo=accel your_program_acc.c
```

GPU Programming Tools

OpenACC

Kernel directives tell the compiler to generate parallel accelerator kernels for the loop nests following the directive.

```
//Hello_World_OpenACC.c
void Print_Hello_World()
{
#pragma acc kernels
    for(int i = 0; i < 5; i++)
    {
        printf("Hello World!\n");
    }
}
```

Data directives tell the compiler to create code that performs specific data movements and provides hints about data usage.

```
#pragma acc data copy(a)
{
    #pragma acc kernels
    {
        for(int i = 0;
i < n; i++)
        {
            a[i] = 0.0
        }
    }
}
```

GPU Programming Tools

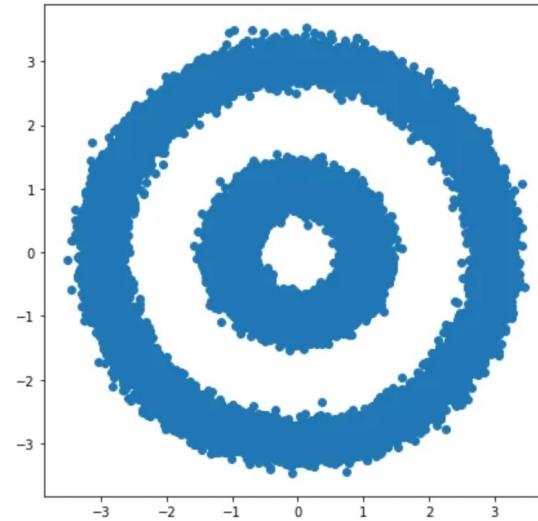
GPU Libraries

- Require little understanding of GPU hardware
- Usually involve minor changes to code
- Accelerated versions of many scientific libraries are available, e.g.
 - LAPACK → MAGMA, libFLAME
 - BLAS → cuBLAS, cUBLAS
 - NumPy & SciPy → cuPy

DBSCAN on CPU

```
#create dataset with 100,000 points using make_circles function
from sklearn.datasets import make_circles
X, y = make_circles(n_samples=int(1e5), factor=.35, noise=.05)

#run DBSCAN clustering algorithm
from sklearn.cluster import DBSCAN
db = DBSCAN(eps=0.6, min_samples=2)
y_db = db.fit_predict(X)
```

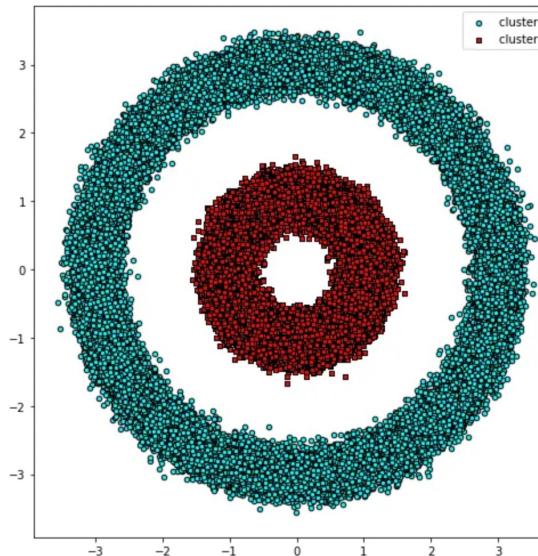


DBSCAN with RAPIDS on GPU

```
#convert data to pandas.DataFrame then create cudf.DataFrame
import pandas as pd
import cudf

X_df = pd.DataFrame({'fea%d'%i: X[:, i] for i in range(X.shape[1])})
X_gpu = cudf.DataFrame.from_pandas(X_df)

#use GPU-accelerated version of DBSCAN from cuML
from cuml import DBSCAN as cuml
DBSCANdb_gpu = cumlDBSCAN(eps=0.6, min_samples=2)
y_db_gpu = db_gpu.fit_predict(X_gpu)
```



Result of running DBSCAN on the CPU using Scikit-Learn

GPU Programming Tools

GPU Frameworks

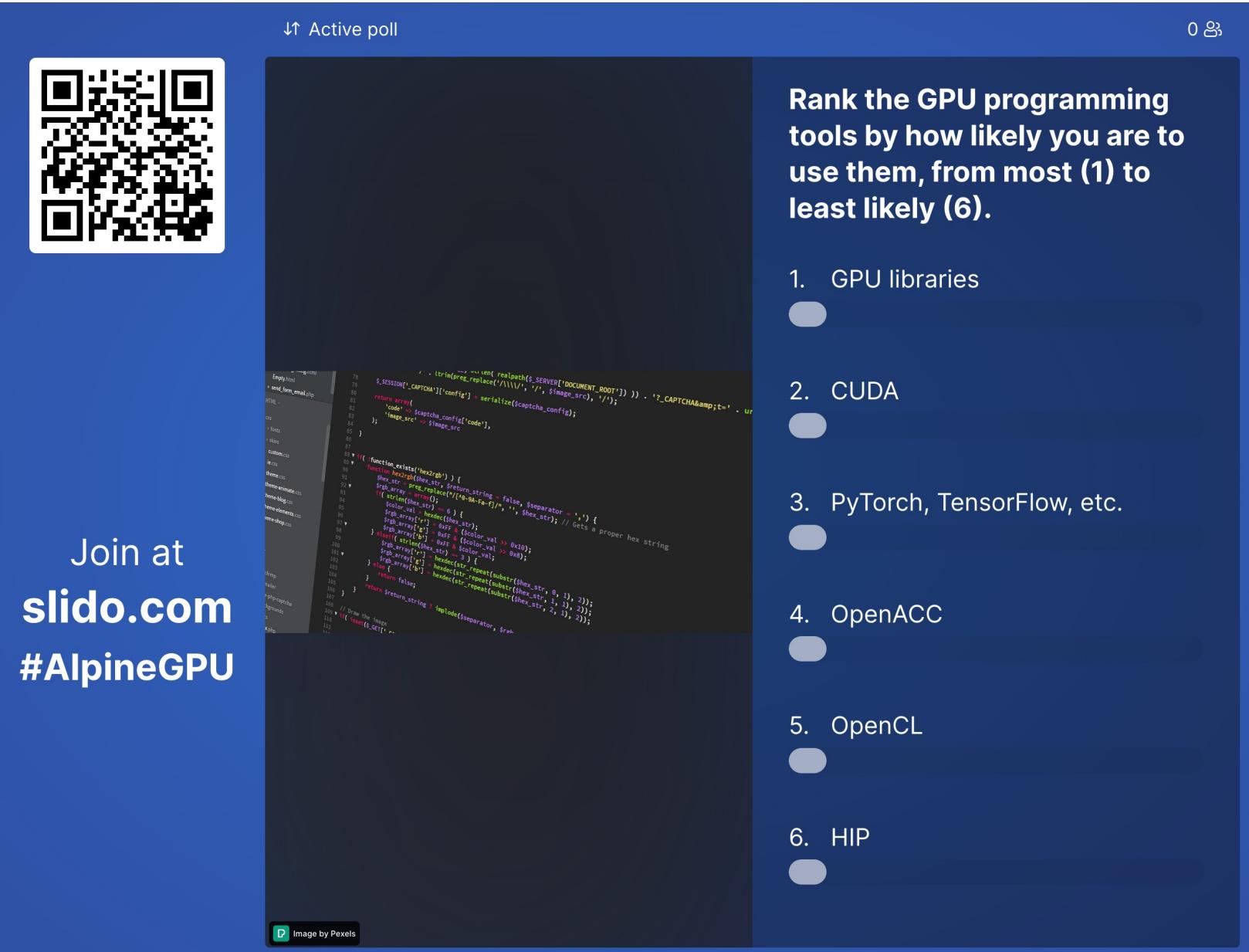
- Offer building blocks for designing, training, and validating workflows (like deep learning)
 - e.g. PyTorch, TensorFlow, Keras
- Rely on GPU-accelerated libraries



GPU Programming Tools



GPU programming tools vary in terms of ease of implementation, portability, and flexibility.



5/24/23

34

GPU Monitoring Tools

- `nvidia-smi` for NVIDIA GPUs
 - Command-line utility tool for monitoring NVIDIA GPUs
 - Returns device- and process-level information
 - Available on Alpine Nvidia nodes without loading any modules

```
NVIDIA-SMI 510.47.03    Driver Version: 510.47.03    CUDA Version: 11.6
+-----+-----+-----+-----+-----+-----+-----+-----+
| GPU  Name      Persistence-M | Bus-Id      Disp.A  | Volatile Uncorr. ECC | | | | | |
| Fan  Temp     Perf  Pwr:Usage/Cap| Memory-Usage | GPU-Util  Compute M. |
|        |          |          |          |      |          |          |      |
| 0  NVIDIA A100-PCI... Off | 00000000:21:00.0 Off | 0MiB / 40960MiB | 0%       Default |
| N/A   36C     P0    40W / 250W |                  |                  |          |          |
|                               |          |          |          |      |          |          |      |
+-----+-----+-----+-----+-----+-----+-----+-----+
| GPU  Name      Persistence-M | Bus-Id      Disp.A  | Volatile Uncorr. ECC | | | | | |
| Fan  Temp     Perf  Pwr:Usage/Cap| Memory-Usage | GPU-Util  Compute M. |
|        |          |          |          |      |          |          |      |
| 1  NVIDIA A100-PCI... Off | 00000000:81:00.0 Off | 0MiB / 40960MiB | 0%       Default |
| N/A   36C     P0    40W / 250W |                  |                  |          |          |
|                               |          |          |          |      |          |          |      |
+-----+-----+-----+-----+-----+-----+-----+-----+
| GPU  Name      Persistence-M | Bus-Id      Disp.A  | Volatile Uncorr. ECC | | | | | |
| Fan  Temp     Perf  Pwr:Usage/Cap| Memory-Usage | GPU-Util  Compute M. |
|        |          |          |          |      |          |          |      |
| 2  NVIDIA A100-PCI... Off | 00000000:E2:00.0 Off | 0MiB / 40960MiB | 0%       Default |
| N/A   37C     P0    40W / 250W |                  |                  |          |          |
|                               |          |          |          |      |          |          |      |
+-----+-----+-----+-----+-----+-----+-----+-----+
Processes:
+-----+-----+-----+-----+-----+-----+-----+
| GPU  GI  CI      PID  Type  Process name      GPU Memory |
| ID   ID              ID           Usage      |
+-----+-----+-----+-----+-----+-----+-----+
| No running processes found |
+-----+-----+-----+-----+-----+-----+-----+
```

GPU Monitoring Tools

- `rocm-smi` for AMD GPUs
 - Command-line utility tool for monitoring AMD GPUs
 - Returns extensive information (use `rocm-smi --help`)
 - Available on Alpine AMD nodes without loading any modules

```
===== ROCm System Management Interface =====
===== Concise Info =====
GPU  Temp   AvgPwr   SCLK    MCLK    Fan     Perf    PwrCap   VRAM%   GPU%
0    33.0c  39.0W   300Mhz  1200Mhz  0%     auto    290.0W   0%     0%
1    35.0c  41.0W   300Mhz  1200Mhz  0%     auto    290.0W   0%     0%
2    34.0c  35.0W   300Mhz  1200Mhz  0%     auto    290.0W   0%     0%
=====
WARNING:          One or more commands failed
===== End of ROCm SMI Log =====
```

The Alpine supercomputer is funded by contributions from the University of Colorado Boulder, the University of Colorado Anschutz, Colorado State University, and the National Science Foundation.

Questions?

