# DS 7347
# High-Performance Computing (HPC) and Data Science
# Session 29

Robert Kalescky
Adjunct Professor of Data Science
HPC Research Scientist

August 2, 2022

Research and Data Sciences Services
Office of Information Technology
Center for Research Computing
Southern Methodist University

Assignments and Project

Session Questions

Distributed Computation

NVSHMEM

# Assignments and Project

- Completed project committed to project repo by August 4, 2022.
- On August 4, 2022 present for 10 minutes your work improving the performance of your workflow.

# Session Questions

Thursday
What's special about RDMA?

Tuesday
What are some benefits and detriments to using GPUs?

Thursday
Generally, what types of computations benefit from using GPUs?

- Heterogeneous computing is the use of hardware that is specialized to a particular task
- Enables hardware based acceleration of programs for specific areas of their code
- Devices are commonly referred to accelerators
- Examples:
    - Cryptographic accelerators
    - Field-programable gate arrays (FPGAs)
    - Digital signal processors (DSPs)
    - Graphics processing units (GPUs)

- General-purpose GPU computing or GPGPU computing is the use of a GPU (graphics processing unit) to do general purpose scientific and engineering computing
- Model for GPU computing is to use a CPU and GPU together for heterogeneous co-processing computing
  - Sequential part of the application runs on the CPU
  - Computationally-intensive part is accelerated by the GPU
- From the user's perspective, the application just runs faster because it is using the high-performance of the GPU to boost performance

host The node.

device The accelerator (GPUs).

host memory The node's memory.

device memory The accelerator's (GPU's) memory.

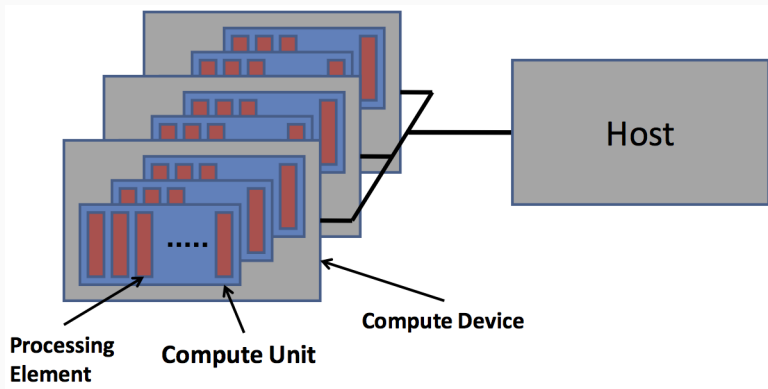kernels GPU function launched by the host and executed on the device

device function GPU function executed on the device which can only be called from the device (i.e. from a kernel or another device function)
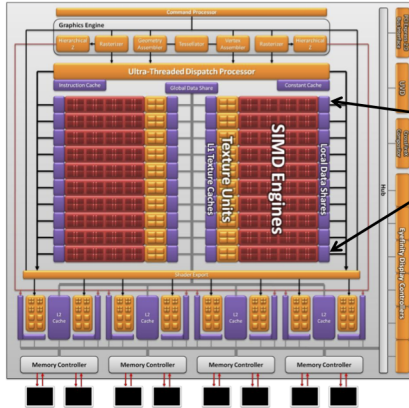
- The GPU is viewed as a compute device that:
    - Is a coprocessor to the CPU or host
    - Has its own DRAM (device memory)
    - Runs many threads in parallel
- Data-parallel portions of an application are executed on the device as kernels which run in parallel on many threads
- Differences between GPU and CPU threads
    - GPU threads are extremely lightweight (very little creation overhead)
    - GPU needs 1000s of threads for full efficiency
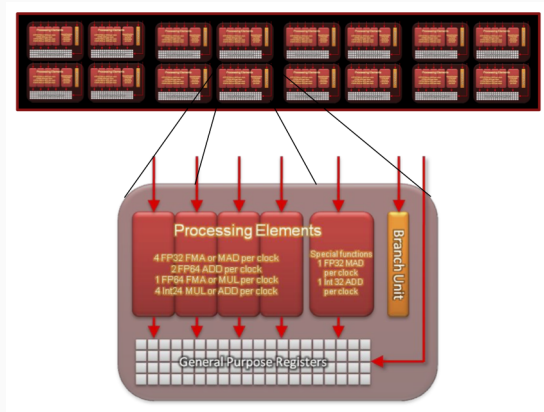    - Multi-core CPU needs only a few

- A kernel is executed as a grid of thread blocks
  - All threads share data memory space
- A thread block is a batch of threads that can cooperate with each other
  - Synchronizing their execution
  - Hazard-free shared memory accesses
  - Efficiently sharing data through a low latency shared memory
- Two threads from two different blocks cannot cooperate

**20 Compute Units**

### Libraries

- There are many libraries available that implement commonly used functions on GPGPUs.

### OpenMP & OpenACC

- These are pragma and directive based tools use to decorate code to provide instructions to compiles as to what and how sections of code can be run on an accelerator.

- OpenACC and OpenMP each began their work independently, but have since begun the work of integrating their standards as they are very similar.

### CUDA

- Popular C-based programming language to program exclusively NVIDIA GPUs
- It has many features that make programming GPUs easier than OpenCL

### OpenCL

- OpenCL (Open Compute Language) is a C-based programming language
- Hardware agnostic
- Supports a wide variety of heterogeneous platforms
- Heterogeneous Programming Tools

There are a handful of high-level languages that allow for GPU-based calculations. Access to the GPU can be through function libraries or directly through lower level languages.

### Kernel Interfaces

- PyOpenCL, PyCUDA, MATLAB, and Mathematica
- Kernels themselves are still written in OpenCL or CUDA specific languages

### Numba

- Numba is a Python just-in-time compiler that can uses decorators to instruct the interpreter when to compile blocks of code.
- The compilation can be done for CPUs and GPUs.

- CUDA is NVIDIA's low level language for programming NVIDIA GPUs
- Available in in three language versions, CUDA C, CUDA C++, and CUDA Fortran
- Theses languages are similar to standard C, C++, and Fortran, but with implicit array or vector programming
- Low-level control over memory management and algorithms

- OpenCL is a vendor neutral standard low level language for programming many types of parallel devices
    - CPUs, GPUs, DSPs, and FPGAs
- Available in two language versions OpenCL C and OpenCL C++
    - OpenCL C also has C++ bindings
- Theses languages are similar to standard C, C++, and Fortran, but with implicit array or vector programming
- Low-level control over memory management and algorithms

- CUDA is designed specifically for NVIDIA hardware
  - Has special knowledge and therefore CUDA code is more performance portable across NVIDIA hardware
- OpenCL is designed to be source portable across a wide variety of hardware
  - Not necessarily performance portable

```
1  __global__ void
2  vectorAdd(const float *A, const float *B, float *C, int numElements)
3  {
4      int i = blockDim.x * blockIdx.x + threadIdx.x;
5
6      if (i < numElements)
7      {
8          C[i] = A[i] + B[i];
9      }
10 }
```

```
1   __kernel void
2   vectorAdd(__global double *a, __global double *b, __global double *c, const
    ↪  unsigned int n)
3   {
4       //Get our global thread ID
5       int id = get_global_id(0);
6
7       //Make sure we do not go out of bounds
8       if (id < n)
9           c[id] = a[id] + b[id];
10  }
```

```
1    #include <stdlib.h>
2
3    void saxpy(int n, float a, float *x, float *restrict y) {
4        #pragma acc parallel loop
5        for (int i = 0; i < n; ++i)
6            y[i] = a * x[i] + y[i];
7    }
8
9    int main(int argc, char **argv) {
10       int N = 1<<20; // 1 million floats
11       if (argc > 1) N = atoi(argv[1]);
12       float *x = (float*)malloc(N * sizeof(float));
13       float *y = (float*)malloc(N * sizeof(float));
14       for (int i = 0; i < N; ++i) {
15           x[i] = 2.0f;
16           y[i] = 1.0f;
17       }
18       saxpy(N, 3.0f, x, y);
19       return 0;
20   }
```

```
1   @cuda.jit
2   def f(a, b, c):
3       # like threadIdx.x + (blockIdx.x * blockDim.x)
4       tid = cuda.grid(1)
5       size = len(c)
6
7       if tid < size:
8           c[tid] = a[tid] + b[tid]
```
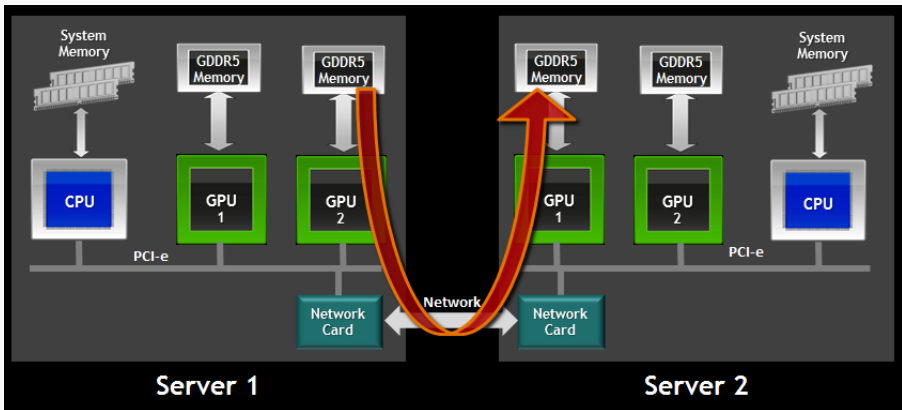
```
1    @cuda.jit
2    def matmul(A, B, C):
3        """Perform square matrix multiplication of C = A * B."""
4        i, j = cuda.grid(2)
5        if i < C.shape[0] and j < C.shape[1]:
6            tmp = 0.
7            for k in range(A.shape[1]):
8                tmp += A[i, k] * B[k, j]
9            C[i, j] = tmp
```
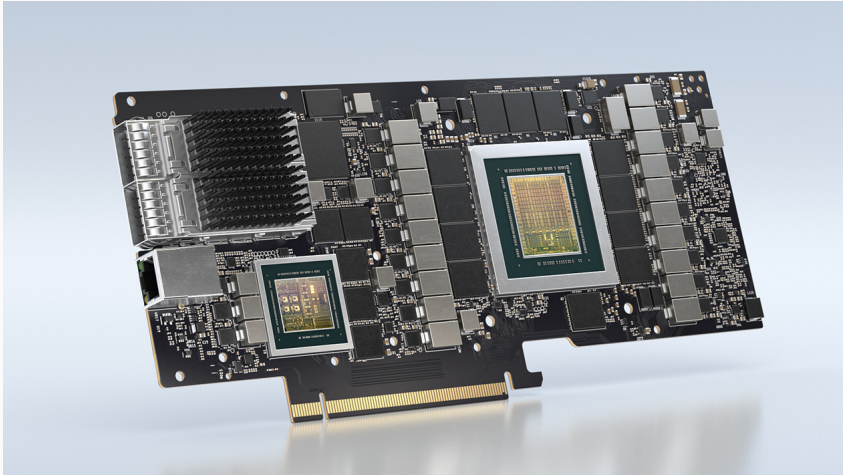
# Distributed Computation

- Unified Communication X (UCX)
- NVIDIA Mellanox MXM
- InfiniBand Verbs

- Message Passing Interface (MPI) (UCX)
- NVIDIA Collective Communication Library (NCCL) (UCX, NVLink)
- OpenSHMEM (UCX)
- NVSHMEM (UCX, NVLink)
- Apache Spark (UCX)

- Message Passing Interface (MPI) (UCX)
- NVIDIA Collective Communication Library (NCCL) (UCX, NVLink)
- OpenSHMEM (UCX)
- NVSHMEM (UCX, NVLink)
- Apache Spark (UCX)

- mpi4py (MPI)
- Ray (MPI)
- Dask (MPI, Ray, Spark)
- Horovod (MPI, Ray, Spark)

- Goal is to make it easy to take a single-GPU training script and successfully scale it to train across many GPUs in parallel.
- Optimized functions built using standard tools, such as MPI.
- Reimplements or wraps core functions of other frameworks.
    - TensorFlow
    - Keras
    - PyTorch
    - MXNet

- Goal is to make it easy to take a single-GPU training script and successfully scale it to train across many GPUs in parallel.
- Optimized functions built using standard tools, such as MPI.
- Reimplements or wraps core functions of other frameworks.
    - TensorFlow
    - PyTorch
    - Horovod (Automatically builds cluster)

# NVSHMEM

- Programming model for efficient and scalable multi-GPU codes.
- Based on OpenSHMEM.
- Provides a partitioned global address space (PGAS) that spans the memory of all included GPUs.
- Communication is initiated directly from the GPUs, which can be more performant than CPU-bound distributed programming models.

```
1   #include <iostream>
2   #include <curand_kernel.h>
3
4   #include <nvshmem.h>
5   #include <nvshmemx.h>
6
7   inline void CUDA_CHECK (cudaError_t err) {
8       if (err != cudaSuccess) {
9           fprintf(stderr, "CUDA error: %s\n", cudaGetErrorString(err));
10          exit(-1);
11      }
12  }
13
14  #define N 1024*1024
```

Listing 1: Includes and definitions.

```
16   __global__ void calculate_pi(int* hits, int seed) {
17       int idx = threadIdx.x + blockIdx.x * blockDim.x;
18
19       // Initialize random number state (unique for every thread in the grid)
20       int offset = 0;
21       curandState_t curand_state;
22       curand_init(seed, idx, offset, &curand_state);
23
24       // Generate random coordinates within (0.0, 1.0]
25       float x = curand_uniform(&curand_state);
26       float y = curand_uniform(&curand_state);
27
28       // Increment hits counter if this point is inside the circle
29       if (x * x + y * y <= 1.0f) {
30           atomicAdd(hits, 1);
31       }
32   }
```

**Listing 2:** Monte Carlo method for estimating $\pi$ CUDA kernel.

```
35    int main(int argc, char** argv) {
36        // Initialize NVSHMEM
37        nvshmem_init();
38
39        // Obtain our NVSHMEM processing element ID and number of PEs
40        int my_pe = nvshmem_my_pe();
41        int n_pes = nvshmem_n_pes();
42
43        // Each PE (arbitrarily) chooses the GPU corresponding to its ID
44        int device = my_pe;
45        CUDA_CHECK(cudaSetDevice(device));
```

Listing 3: Initialize NVSHMEM and get local processing element (PE) ID and global number of PEs.

```
47      // Allocate host and device values
48      int* hits = (int*) malloc(sizeof(int));
49      int* d_hits = (int*) nvshmem_malloc(sizeof(int));
50
51      // Initialize number of hits and copy to device
52      *hits = 0;
53      CUDA_CHECK(cudaMemcpy(d_hits, hits, sizeof(int), cudaMemcpyHostToDevice));
```

Listing 4: Allocate memory on hosts and devices and initialize the number of hits.

```
55    // Launch kernel to do the calculation
56    int threads_per_block = 256;
57    int blocks = (N / n_pes + threads_per_block - 1) / threads_per_block;
58
59    int seed = my_pe;
60    calculate_pi<<<blocks, threads_per_block>>>(d_hits, seed);
61    CUDA_CHECK(cudaDeviceSynchronize());
62
63    // Accumulate the results across all PEs
64    int* d_hits_total = (int*) nvshmem_malloc(sizeof(int));
65    nvshmem_int_sum_reduce(NVSHMEM_TEAM_WORLD, d_hits_total, d_hits, 1);
```

Listing 5: Launch kernel on the devices to perform the computations, block until all are finished, and then accumulate the results.

```
67        if (my_pe == 0) {
68            // Copy final result back to the host
69            CUDA_CHECK(cudaMemcpy(hits, d_hits_total, sizeof(int),
   ↪    cudaMemcpyDeviceToHost));
70
71            // Calculate final value of pi
72            float pi_est = (float) *hits / (float) (N) * 4.0f;
73
74            // Print out result
75            std::cout << "Estimated value of pi averaged over all PEs = " << pi_est <<
   ↪    std::endl;
76            std::cout << "Relative error averaged over all PEs = " << std::abs((M_PI -
   ↪    pi_est) / pi_est) << std::endl;
77        }
```

Listing 6: Only on the first PE, copy accumlated hits back to host and report estimated value of $\pi$.

```
79        free(hits);
80        nvshmem_free(d_hits);
81        nvshmem_free(d_hits_total);
82
83        // Finalize nvshmem
84        nvshmem_finalize();
85
86        return 0;
87    }
```

Listing 7: Free memory on hosts and device, finalize NVSHMEM, and exit.

```bash
1  #!/bin/bash
2  #SBATCH -J nvshmem_pi         # Job name
3  #SBATCH -o nvshmem_pi_%j.out  # Output file name
4  #SBATCH -p amd                # Queue (partition)
5  #SBATCH -c 1                  # Cores
6  #SBATCH --mem=6G              # Memory
7  #SBATCH --gres=gpu:1          # GPUs
8  #SBATCH -t 5                  # Time limit
```

Listing 8: Request compute resources.

```
10   echo $SLURM_JOB_PARTITION
11
12   module purge
13   module load nvhpc-22.2          # Alternatively nvhpc-21.2, nvhpc-21.9
14
15   NVSHMEM_HOME=/hpc/applications/nvidia/hpc_sdk/2022_22.2/Linux_x86_64/22.2/comm_libs/11.2/
16   GCC_HOME=/hpc/spack/opt/spack/linux-centos7-x86_64/gcc-7.3.0/gcc-9.2.0-
     ↪  6zgrndxveon2m5mjhltrqccdcewrdktx/bin
```

Listing 9: Setup the environment. $GCC_HOME is to provide full C++ 11 support.

```
18  nvcc -ccbin=$GCC_HOME -x cu -arch=sm_60 -rdc=true -I $NVSHMEM_HOME/include\
19   -L $NVSHMEM_HOME/lib -lnvshmem -lcuda -o nvshmem_pi nvshmem_pi.cpp
20  # CUDA Compute Capabilities
21  # P100 sm_60
22  # V100 sm_70
23  # A100 sm_80
24
25  srun nvshmem_pi
```

Listing 10: Build the executable for the specific PE being used and run.