

Introduction to GPU Programming

Wim R.M. Cardoen, PhD

Center of High-Performance Computing (CHPC)
University of Utah

October 14, 2024

Outline I

1 Motivation

2 Hardware

- Streaming multiprocessor (SM)
- Warps

3 Software

- GPGPU & CUDA
- Structure of a GPU computation
- Case study: matrix multiplication
- Building CUDA applications & useful env. variables
- Profiling & debugging
- Important CUDA libraries
- Alternatives to CUDA
- Links

4 Use of GPUs at the CHPC

- GPUs available at CHPC
 - Regular env.: lp/kp/np/grn clusters

Outline II

- Protected env.: redwood cluster
- How to access the GPUs at CHPC

Motivation

Theoretical GFLOP/s: GPU vs. CPU

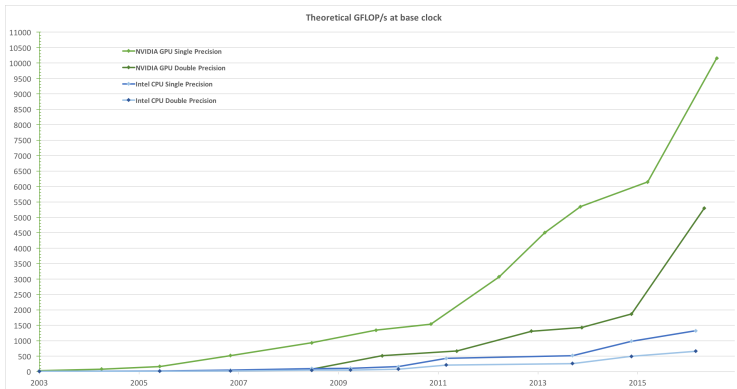
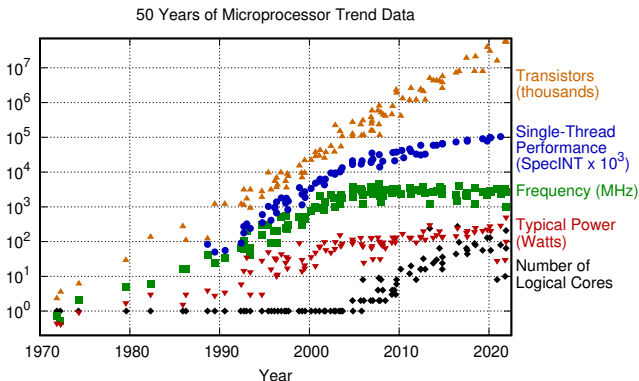


Figure: Theoretical GFLOP/s: GPUs vs. CPUs.^a

^ahttps://docs.nvidia.com/cuda/archive/9.1/pdf/CUDA_C_Programming_Guide.pdf

CPU processor trend (last 50 years)



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2021 by K. Rupp

- After the year 2000, freq./power for a single CPU core reaches a max. (**Heat dissipation!**).

Energy efficiency per job: GPU vs. CPU

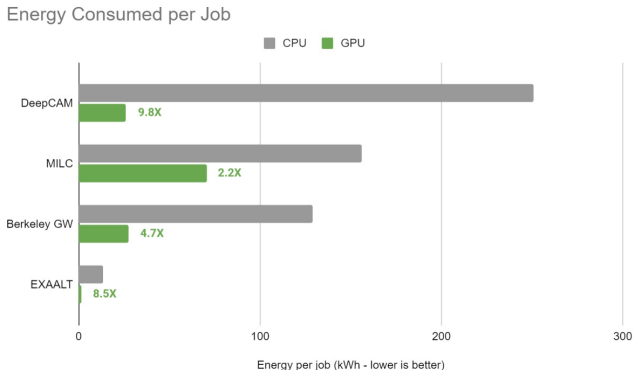


Figure: Energy efficiency per job (NERSC).^a

^a<https://blogs.nvidia.com/blog/gpu-energy-efficiency-nersc/> (05/21/2023)

Hardware

Streaming Multiprocessor (SM)

- GPU device connected to the CPU by a PCIe bus.
- each GPU device contains an array (**x**) of Streaming Multiprocessors (**SM**).
- each SM has:
 - a Single-Instruction Multiple-Thread (SIMT) Architecture.
 - contains **y** regular cores and [**z** tensor cores].
- scalable: newer generations: increase of **x**, **y** and [**z**], e.g.:
 - NVIDIA A100-PCIE-40GB (*notch293*)
 - global memory: 40 GB.
 - 108 SMs, 64 Cores/SM, 4 Tensor Cores/SM.
 - GPU Max. Clock Rate: 1.41 GHz.
 - NVIDIA H100 SXM5 NVL (*grn008*)
 - global memory: 93 GB.
 - 132 SMs, 128 Cores/SM, 4 Tensor Cores/SM.
 - GPU Max. Clock Rate: 1.78 GHz.

NVIDIA GH100 SM

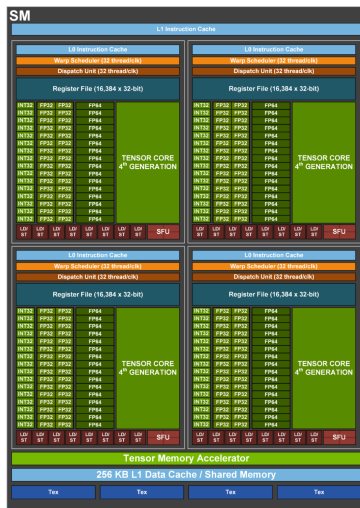


Figure: GH100 SM

NVIDIA GH100 Full Device



Figure: NVIDIA GH100 Full Device (144 SMs).

GPU Threads - Warps

- Each SM:
 - generates, schedules, executes threads in batches of 32 threads.
 - **WARP**: a batch of 32 threads
- each thread in a WARP executes the same instructions but runs its own "path".
- if threads within a WARP diverge, the threads become inactive/disabled.

Software

- GPU (Graphic Processing Unit):
originally developed for graphical applications.
- GP-GPU: General-Purpose GPU, i.e.
the use of GPUs beyond graphical applications.
CAVEAT: problem to be reformulated in terms of the graphics API.
- **2007**: NVIDIA introduces the **CUDA**¹ framework
(**C**ompute **U**nified **D**evice **A**rchitecture)
 - CUDA API: extension of the C language.
 - handles the GPU thread level parallelism.
 - deals with moving data between CPU and GPU.
 - also support for C++, Fortran and Python.

¹The **CUDA Toolkit** consists of 2 parts:

- CUDA Driver
- CUDA Toolkit (nvcc, nvprof, . . . , libraries, header files).

Structure of a GPU computation

- 1 **Allocate** memory space on the GPU device.
- 2 **Transfer** the data from the CPU to the GPU device.
- 3 Perform the **calculation** on the GPU device.
 - **kernel**: function executed on the GPU.
 - To enhance performance: keep data as long as possible on the GPU device.
- 4 **Transfer** the result back from the GPU device to the CPU.
- 5 **Deallocate** memory space on the GPU device.

Note: source code & makefile available in `./src`

Alloc. & free of global memory on the GPU

- `cudaError_t`
CUDA Error types.
- `cudaError_t cudaMalloc(void **devPtr, size_t size)`
Allocates memory on the device.
- `cudaError_t cudaFree(void *devPtr)`
Frees memory on the device.


```

double *M_d, *N_d, *P_d; // Pointers (device)
int const SZ=16;
int const SZ2=SZ*SZ;

// Allocate M on device (M_d)
if(cudaMalloc(&M_d, sizeof(double)*SZ2) != cudaSuccess)
{
    printf(" ERROR: alloc vector M on DEVICE \n");
    return 1;
}

// Deallocate matrix M_d
if(cudaFree(M_d) != cudaSuccess)
{
    printf(" ERROR: unable to deallocate M_d (DEVICE)\n");
    return 1;
}

```

Listing 1: Alloc/Free extract

Copy data between host (CPU) and device (GPU)

- Copy data between host (CPU) and device (GPU)

```
cudaError_t cudaMemcpy(void *dst, const void *src,  
                      size_t count, cudaMemcpyKind kind)
```

- Direction (kind):
 - cudaMemcpyHostToHost
 - cudaMemcpyHostToDevice
 - cudaMemcpyDeviceToHost
 - cudaMemcpyDeviceToDevice

```

// Copy M from Host onto Device
if (cudaMemcpy(M_d, M_h, sizeof(double)*SZ2,
               cudaMemcpyHostToDevice) != cudaSuccess)
{
    printf(" ERROR: copy vector M: HOST -> DEVICE\n");
    return 1;
}

// Copy P (result) from the Device to the Host
if (cudaMemcpy(P_h, P_d, sizeof(double)*SZ2,
               cudaMemcpyDeviceToHost) != cudaSuccess)
{
    printf(" ERROR: copy vector N: DEVICE -> HOST\n");
    return 1;
}

```

Listing 2: cudaMemcpy extract

- **CUDA kernel**: alias for a function which may run on a GPU device.
- **Kernel declaration** syntax:

```
funcspec void kernelName(args){ body }
```

where:

- *funcspec*: function type qualifier, i.e.
`__global__`, `__host__`, `__device__`
 - *kernelName*: name of the kernel/CUDA function.
 - *args*: argument list of the kernel/CUDA function.
 - *body*: body of the kernel/CUDA function (your code).
- **Kernel call** syntax:

```
kernelName<<<gridSize,blockSize>>>(args)
```

where:

- *gridSize*: size of the grid of thread blocks.
- *blockSize*: size of a thread block.

Function type qualifiers

Qualifier	Called from	Executed on
<code>__global__</code>	host	device
<code>__host__</code>	host	host
<code>__device__</code>	device	device

Table: Function type qualifiers

Note:

- You can have to different versions of a function i.e.: one with `__host__` & one with `__device__`

Grids, Blocks and Threads

We have a hierarchical (software) implementation.

- `uint3,dim3`:
 - CUDA defined structures of unsigned integer `x,y,z`
 - `dim3`: based on `uint3` but unspecified components are initialized to 1.
- **Grid**: each Grid consists of Blocks
 - `dim3 gridDim`: dimensions of the Grid.
 - `uint3 blockIdx`: block index within the Grid.
- **Block**: each Block consists of Threads
 - `dim3 blockDim`: dimensions of the Block
 - `uint3 threadIdx`: thread index within the block.

Matrix Mul.: kernel (v.1)

Case study: $P = M \times N$ where $P, M, N \in \mathbb{R}^{n \times n}$

```
#include <mul.h>

__global__ void MatrixMulKernel1(double *M_d, double *N_d,
                                double *P_d, int const SZ)
{
    double Pval=0;
    int tx=threadIdx.x;
    int ty=threadIdx.y;
    for(int k=0; k<SZ; k++)
        Pval+=M_d[ty*SZ+k]*N_d[tx+k*SZ];
    P_d[tx+ty*SZ]=Pval;
    return;
}
```

Listing 3: Kernel (v.1)

Invoking kernel (v. 1)

- Invoking 1 Block of Threads

```
int main(void)
{
    int const SZ=16;
    // ...
    // Invoke Kernel to generate P=MxN
    dim3 dimBlock(SZ,SZ,1);
    dim3 dimGrid(1,1,1);
    MatrixMulKernel1<<<<dimGrid , dimBlock>>>>(M_d,N_d,P_d,SZ);
    //..
}
```

Listing 4: Invoking kernel (v. 1)

Mat. Mul (v.2): Grid of 2D Blocks

- `int tx = blockIdx.x*blockDim.x + threadIdx.x;`
- `int ty = blockIdx.y*blockDim.y + threadIdx.y;`

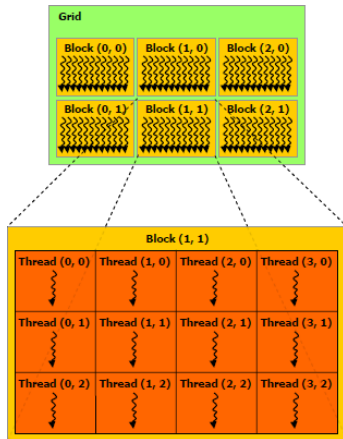


Figure: 2D-Grid of 2D-Blocks of Threads

Matrix Mul.: visualization (v. 2)

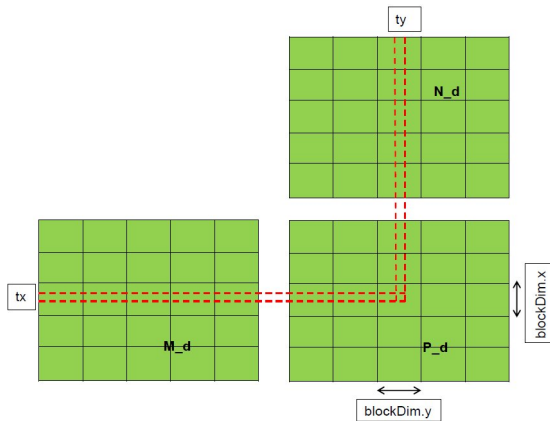


Figure: Matrix Mul. (2D Grid)

Matrix Mul.: kernel (v.2)

```
#include <mul.h>

__global__ void MatrixMulKernel2(double *M_d, double *N_d,
                                double *P_d, int const SZ)
{
    int tx = blockIdx.x * blockDim.x + threadIdx.x;
    int ty = blockIdx.y * blockDim.y + threadIdx.y;

    if ((tx < SZ) && (ty < SZ))
    {
        double Pval=0.0;
        for(int k=0; k<SZ; k++)
            Pval += M_d[tx*SZ+k]*N_d[k*SZ+ty];
        P_d[tx*SZ+ty]=Pval;
    }
    return;
}
```

Listing 5: Kernel (v.2)

Invoking kernel (v.2)

- Invoking a grid of blocks of threads

```
int main(void)
{
    int const SZ=500;

    // ..
    int const THREADX=16;
    int const THREADY=16;
    dim3 dimBlock(THREADX,THREADY,1);
    int numBlocksX=(SZ%THREADX==0 ? SZ/THREADX : SZ/THREADX +1);
    int numBlocksY=(SZ%THREADY==0 ? SZ/THREADY : SZ/THREADY +1);
    dim3 dimGrid(numBlocksX,numBlocksY,1);

    MatrixMulKernel2<<<<dimGrid,dimBlock>>>>(M_d,N_d,P_d,SZ);
    // ..
}
```

Listing 6: Invoking kernel (v.2)

Types of GPU memory

- **global memory**: (largest, slowest and often the bottleneck).
- **shared memory**: allocated per thread block & low latency
- **constant memory**: cached, read-only
- **registers**: fast, on-chip memory (exclusive to each thread).

Matrix Mul.: use of shared memory (v.3)

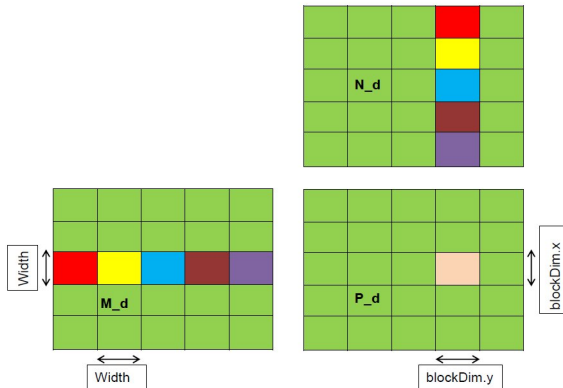


Figure: Matrix Mul.: use of shared memory

Matrix Mul.: kernel (v.3) - use of shared memory

```
#include <mul.h>
#define WIDTH 16
__global__ void MatrixMulKernel3(double *M_d, double *N_d, double *P_d, int const SZ)
{
    int tx = blockIdx.x * blockDim.x + threadIdx.x;
    int ty = blockIdx.y * blockDim.y + threadIdx.y;

    __shared__ double M_s[WIDTH][WIDTH];
    __shared__ double N_s[WIDTH][WIDTH];

    double Pval=0.0;
    int nslices=(SZ%WIDTH==0)?(SZ/WIDTH):(SZ/WIDTH+1);
    for(int islice=0; islice<nslices; islice++)
    {
        M_s[threadIdx.x][threadIdx.y]=M_d[tx*SZ + islice*WIDTH + threadIdx.y];
        N_s[threadIdx.x][threadIdx.y]=N_d[islice*WIDTH*SZ + threadIdx.x*SZ + ty];
        __syncthreads();

        for(int k=0; k<WIDTH; k++)
            Pval += M_s[threadIdx.x][k]*N_s[k][threadIdx.y];
        __syncthreads();
    }

    if(tx<SZ && ty<SZ)
        P_d[tx*SZ+ty]=Pval;
    return;
}
```

Listing 7: Kernel (v.3)

Building/Compiling CUDA applications

General scheme:

- Source code for CUDA applications:
 - C/C++ host code with extensions to deal with the device(s).
 - Other programming languages are allowed e.g. Fortran
- Primo: **separate** device functions from host code.
- **Device code**: preprocessing, compilation with the NVIDIA compiler (`nvcc`).
- **Host code**: preprocessed, compiled with a host (C/C++) compiler e.g. (`gcc`, `g++`, `icc`, `icpc`, ...)
- Compiled device functions are **embedded** as fatbinary images in the host object file.
- **Linking stage**: adding CUDA runtime libraries to the host object file to create an executable.

Further concepts

- .cu : Suffix for CUDA source file (host code (C,C++) & device code).
- .cuf: Suffix for CUDA source file (host code (Fortran) & device code).
- .ptx: Suffix for **P**arallel **T**hread **E**xecution (PTX) files. An intermediate representation (similar to assembly for a **virtual GPU architecture**²
- .cubin: Suffix for the **C**UDA device **b**inary file pertaining to a **real GPU architecture**³
- fatbin: Multiple PTX [& cubin] files are merged into a fatbin file.

²Virtual architectures bear the “compute_” **prefix** e.g. “compute_70”.

³Real (physical) architectures bear the “sm_” **prefix** e.g. “sm_70”.

Memento: “sm” stands for the physical streaming multiprocessor.

Some recent CUDA Architectures

Architecture	Year	compute_	sm_
Maxwell	2014	50, 52, 53	50, 52, 53
Pascal	2016	60, 61, 62	60, 61, 62
Volta	2017	70, 72	70, 72
Turing	2018	75	75
Ampere	2020	80, 86, 87	80, 86, 87
Ada Lovelace	2022	89	89
Hopper	2022	90, 90a	90, 90a

Table: Some of the recent CUDA architectures (10/08/2024)

Compilation trajectory (cont.)

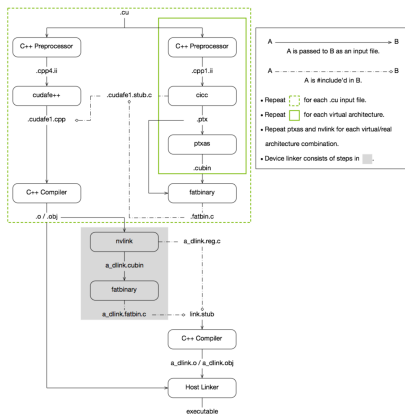


Figure: Compilation trajectory

- What are the existing CUDA architectures?
- How find the architecture of a machine?
- How to build for a particular architecture?
- How to build for multiple architectures?

CUDA SDK comes with:

- its own profiler: `nvprof`.
- its own debugger: `nvsight`

Profiling mu13 using nvprof

```
[u0253283@notch001:3]$ nvprof ./mul3
==2424062== NVPROF is profiling process 2424062, command: ./mul3
Calling Kernel ...
Kernel Call Finished ...

Frob. Norm(P-P_h): 0.0000000000

==2424062== Profiling application: ./mul3
==2424062== Profiling result:
```

Type	Time(%)	Time	Calls	Avg	Min	Max	Name
GPU activities:	59.70%	1.1707ms	1	1.1707ms	1.1707ms	1.1707ms	MatrixMulKernel3(double*, double*, double*, int)
	32.02%	627.87us	2	313.93us	307.13us	320.73us	[CUDA memcpy HtoD]
	8.28%	162.33us	1	162.33us	162.33us	162.33us	[CUDA memcpy DtoH]
API calls:	95.82%	138.45ms	3	46.151ms	109.78us	138.23ms	cudaMalloc
	2.88%	4.1668ms	3	1.3889ms	494.15us	3.1420ms	cudaMemcpy
	0.82%	1.1901ms	1	1.1901ms	1.1901ms	1.1901ms	cudaLaunchKernel
	0.27%	395.44us	3	131.81us	94.455us	202.38us	cudaFree
	0.17%	250.22us	114	2.1940us	291ns	103.89us	cuDeviceGetAttribute
	0.01%	16.428us	1	16.428us	16.428us	16.428us	cuDeviceGetPCIBusId
	0.01%	12.713us	1	12.713us	12.713us	12.713us	cuDeviceGetName
	0.00%	1.5660us	3	522ns	321ns	905ns	cuDeviceGetCount
	0.00%	1.2730us	2	636ns	338ns	935ns	cuDeviceGet
	0.00%	613ns	1	613ns	613ns	613ns	cuModuleGetLoadingMode
	0.00%	556ns	1	556ns	556ns	556ns	cuDeviceTotalMem
	0.00%	547ns	1	547ns	547ns	547ns	cudaGetLastError
	0.00%	337ns	1	337ns	337ns	337ns	cuDeviceGetUuid

Figure: Profiling mu13 on notch001

Important CUDA libraries

In order to increase the performance of your code we recommend to use highly-optimized libraries. Among them, we have:

- **cuBLAS**: **B**asic **L**inear **A**lgebra **S**ubroutines on NVIDIA GPUs.
- **MAGMA**: **M**atrix **A**lgebra on **G**PU and **M**ulti-core **A**rchitectures.
- **cuRAND**: Random Number Generation library.
- **cuFFT**: CUDA **F**ast **F**ourier **T**ransform library.
- **NCCL**: **N**VIDIA **C**ollective **C**ommunications **L**ibrary.
- **cuDNN**: CUDA **D**eep **N**eural **N**etwork library.
- **cuTENSOR**: GPU-accelerated Tensor Linear Algebra.
- **DALI**: Library for decoding & augmenting images (DL applications).
- ...

Alternatives to CUDA

- Similar to CUDA
 - **ROCM** (AMD)
- **OpenACC** (use of directives (cfr. OpenMP))
 - GCC: supports OpenACC for NVIDIA & AMD GPUs.
 - NVIDIA HPC SDK (formerly PGI)
 - Sourcery Codebench (AMD GPU)
- Higher-level abstractions
 - **Kokkos** (prog. model for parallel algorithms for many-core chips)

- [CUDA Toolkit Documentation](#)
- [CUDA C++ Programming Guide Release 12.6 \(10/01/24\)](#)
- [CUDA C++ Best Practices Guide, Release 12.6 \(09/24/24\)](#)
- [NVIDIA CUDA Compiler Driver NVCC, Release 12.6 \(09/24/24\)](#)
- [PTX & ISA Release 8.5 \(09/24/24\)](#)

Use of GPUs at the CHPC

GPU devices on lp/kp/np/grn

GPU device type	compute capability
NVIDIA GeForce GTX TITAN X	5.2
Tesla P100-PCIE-16GB	6.0
Tesla P40	6.1
NVIDIA GeForce GTX 1080 Ti	6.1
NVIDIA Titan V	7.0
NVIDIA Tesla V100-PCIE-16GB	7.0
Tesla T4	7.5
NVIDIA GeForce RTX 2080 Ti	7.5
NVIDIA A100-PCIE-40GB	8.0
NVIDIA A100-SXM4-80GB	8.0
NVIDIA A800 40GB Active	8.0

Table: GPU devices on lp/kp/np/grn (10/01/2024)

GPU devices on lp/kp/np/grn (cont.)

GPU device type	compute capability
NVIDIA GeForce RTX 3090	8.6
NVIDIA A40	8.6
NVIDIA RTX A5500	8.6
NVIDIA RTX A6000	8.6
NVIDIA RTX 6000 Ada Generation	8.9
NVIDIA L40	8.9
NVIDIA L40S	8.9
NVIDIA H100 NVL/Deep Dive	9.0

Table: GPU devices on lp/kp/np/grn (10/01/2024)

GPU devices on redwood

GPU device type	compute capability
NVIDIA GeForce GTX 1080 Ti	6.1
NVIDIA A100-SXM4-40GB	8.0
NVIDIA A100 80GB PCIe	8.0
NVIDIA A30	8.0
NVIDIA A40	8.6
NVIDIA RTX 6000 Ada Generation	8.9
NVIDIA H100 NVL/Deep Dive	9.0

Table: GPU devices on redwood (10/01/2024)

- Using GPUs at the CHPC (Presentation by Martin Čuma)

Questions?

Thank you!

Any questions?