

The Graphics Processing Unit (GPU) revolution

Ramu Anandakrishnan

Outline

2

- The need for parallel processing
- Basic parallel processing concepts
- The GPU - a massively parallel processor

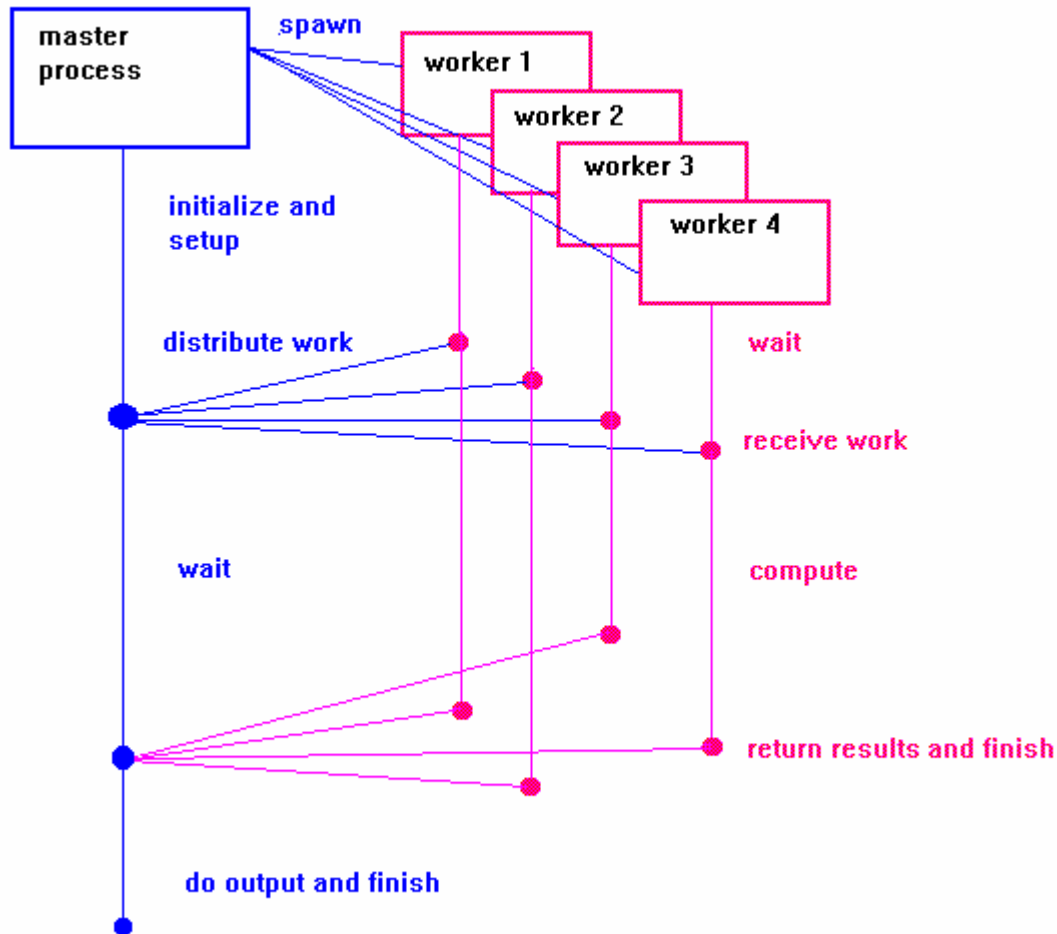
Lecture 1

-
- Introduction to GPU programming
 - Overview of GPU hardware architecture
 - Performance considerations
 - Homework assignment

Lecture 2

Application level parallelism

3

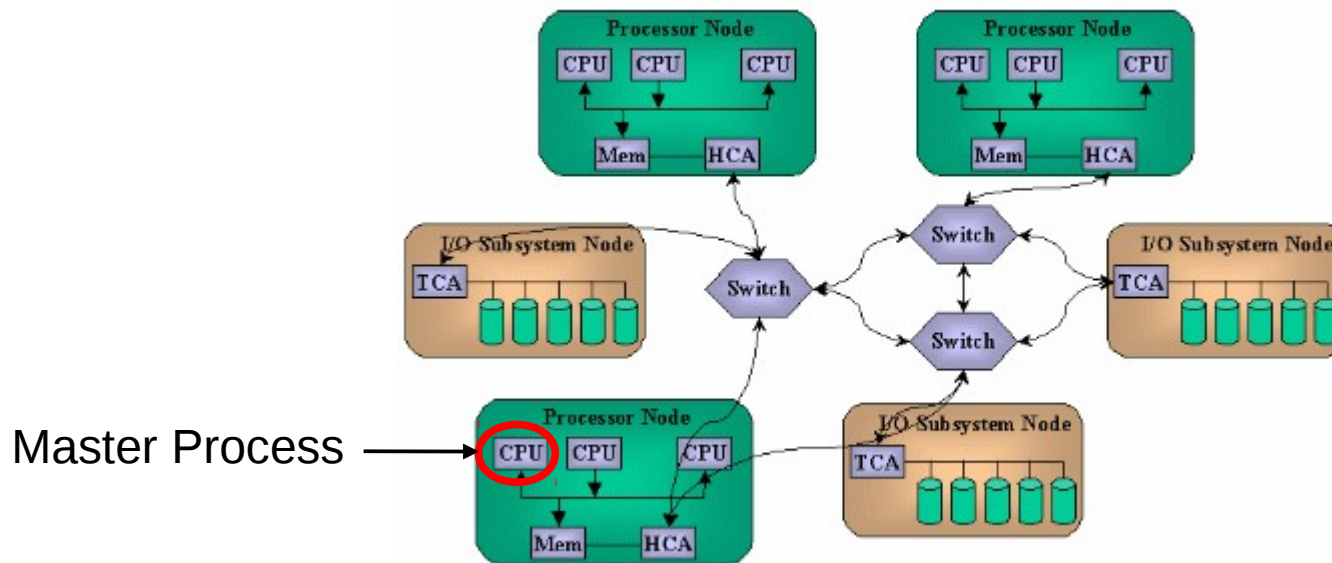


A key limitation of supercomputer architecture – inter node communication

4



(extremetech.com)

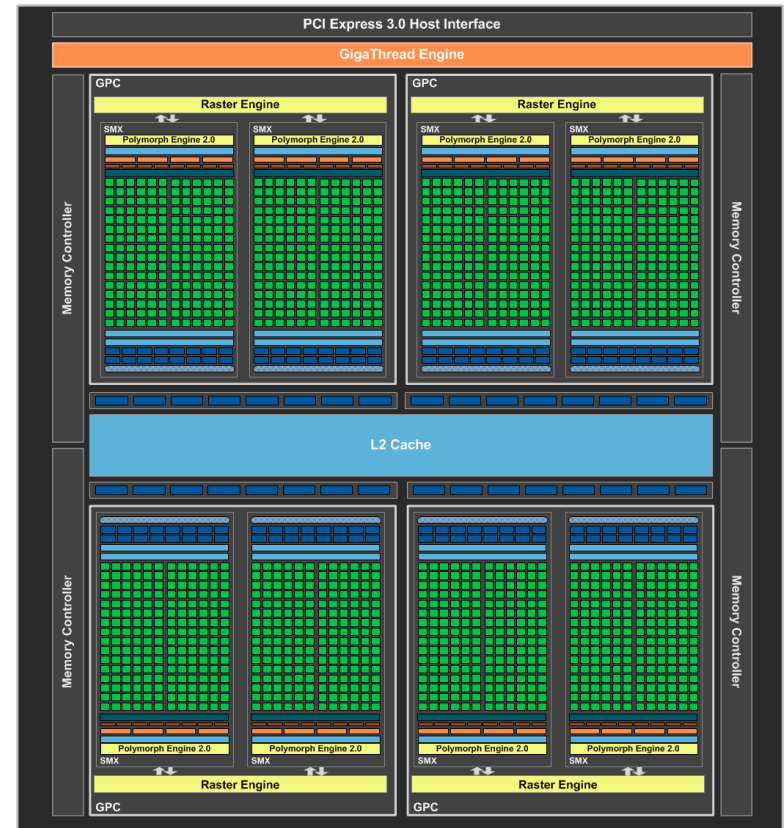
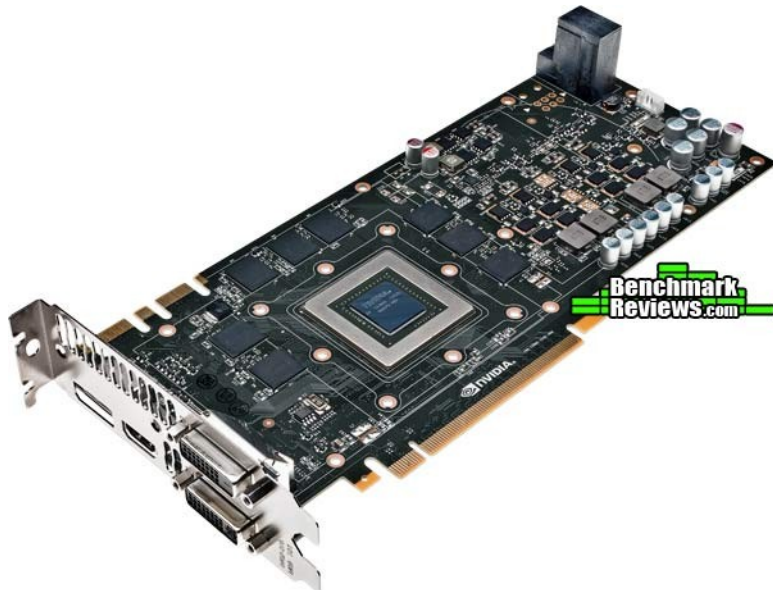


GPUs – thousands of processors on a single node

5

GeForce GTX 690 Specifications

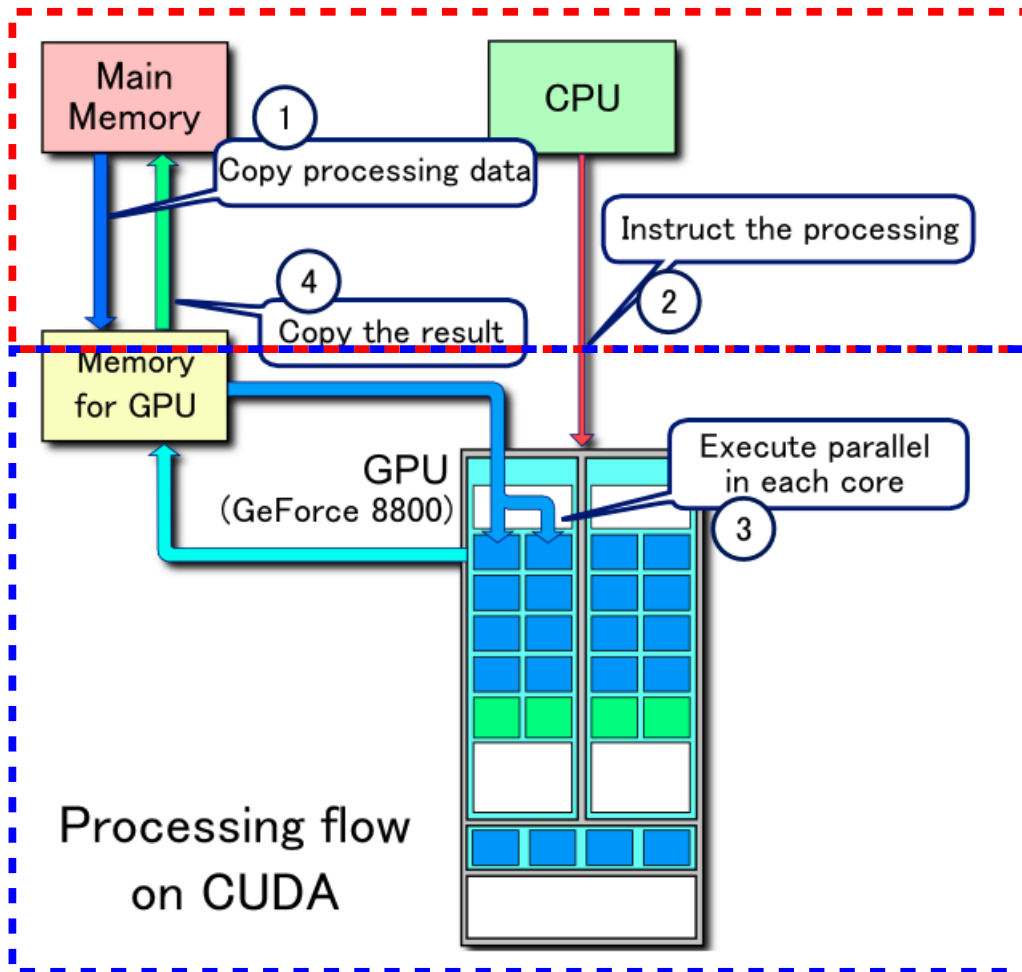
CUDA Cores	3072
Base Clock	915 MHz
Boost Clock	1019 MHz
Memory Config	4GB / 512-bit GDDR5
Memory Speed	6.0 Gbps
Power Connectors	8-pin + 8-pin
TDP	300W
Outputs	3x DL-DVI Mini-Displayport 1.2
Bus Interface	PCI Express 3.0



(www.nvidia.com)

CUDA – for running code on the GPU

6



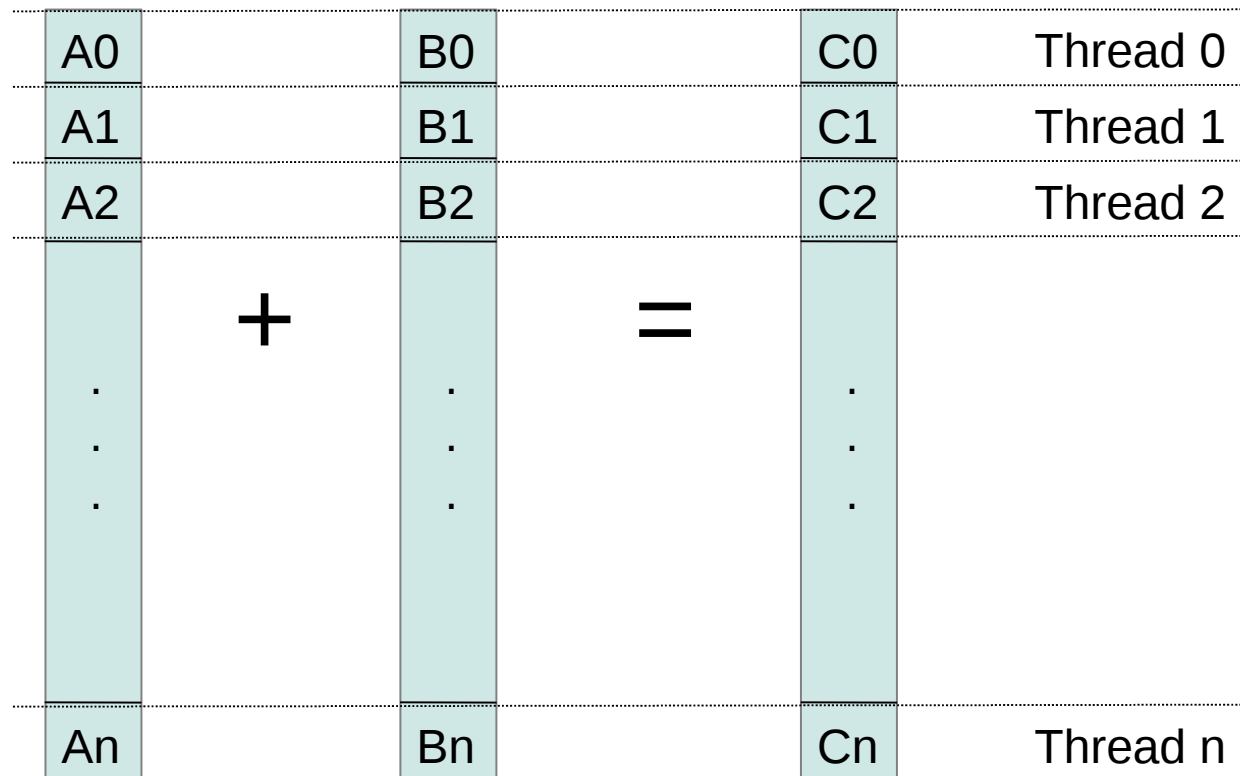
Example of CUDA process flow

1. Copy data from main mem to GPU mem
2. CPU initiates threads on the GPU
3. GPU execute parallel code
4. Copy the result from GPU mem to main mem

Example: Vector addition

7

$$A_i + B_i = C_i$$



1. Host: Copy data from main memory to GPU memory

8

```
// Host code
int main()
{
    int N = ...;
    size_t size = N * sizeof(float);
    // Allocate input vectors h_A and h_B in host memory
    float* h_A = (float*)malloc(size);
    float* h_B = (float*)malloc(size);
    // Initialize input vectors
    ...
    // Allocate vectors in device memory
    float* d_A;
    cudaMalloc(&d_A, size);
    float* d_B;
    cudaMalloc(&d_B, size);
    float* d_C;
    cudaMalloc(&d_C, size);
    // Copy vectors from host memory to device memory
    cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);
    ...
}
```


2. Host: Initiate parallel threads (device kernels) on the GPU

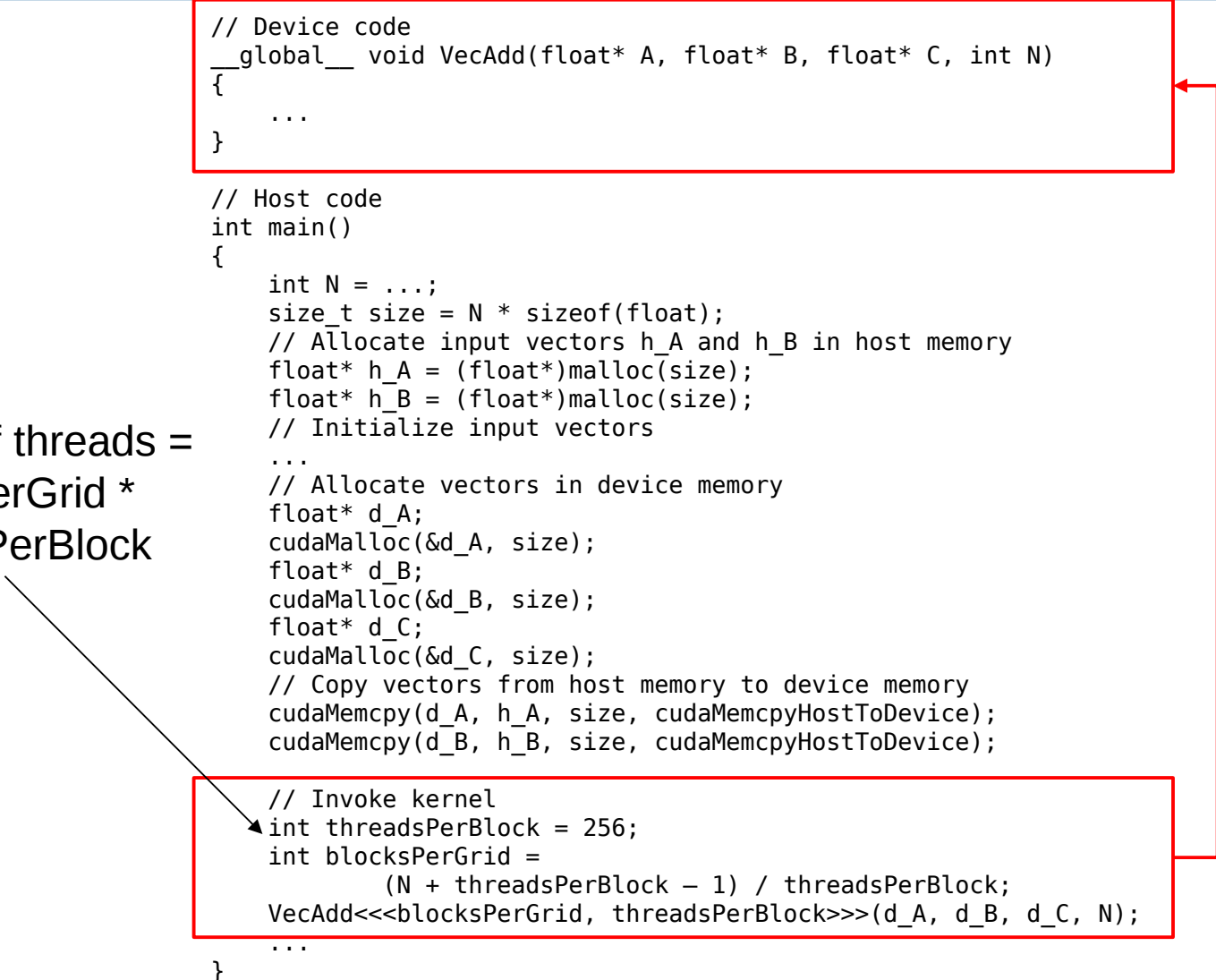
9

```
// Device code
__global__ void VecAdd(float* A, float* B, float* C, int N)
{
    ...
}
```

```
// Host code
int main()
{
    int N = ...;
    size_t size = N * sizeof(float);
    // Allocate input vectors h_A and h_B in host memory
    float* h_A = (float*)malloc(size);
    float* h_B = (float*)malloc(size);
    // Initialize input vectors
    ...
    // Allocate vectors in device memory
    float* d_A;
    cudaMalloc(&d_A, size);
    float* d_B;
    cudaMalloc(&d_B, size);
    float* d_C;
    cudaMalloc(&d_C, size);
    // Copy vectors from host memory to device memory
    cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);
```

```
// Invoke kernel
int threadsPerBlock = 256;
int blocksPerGrid =
    (N + threadsPerBlock - 1) / threadsPerBlock;
VecAdd<<<blocksPerGrid, threadsPerBlock>>>(d_A, d_B, d_C, N);
...
}
```

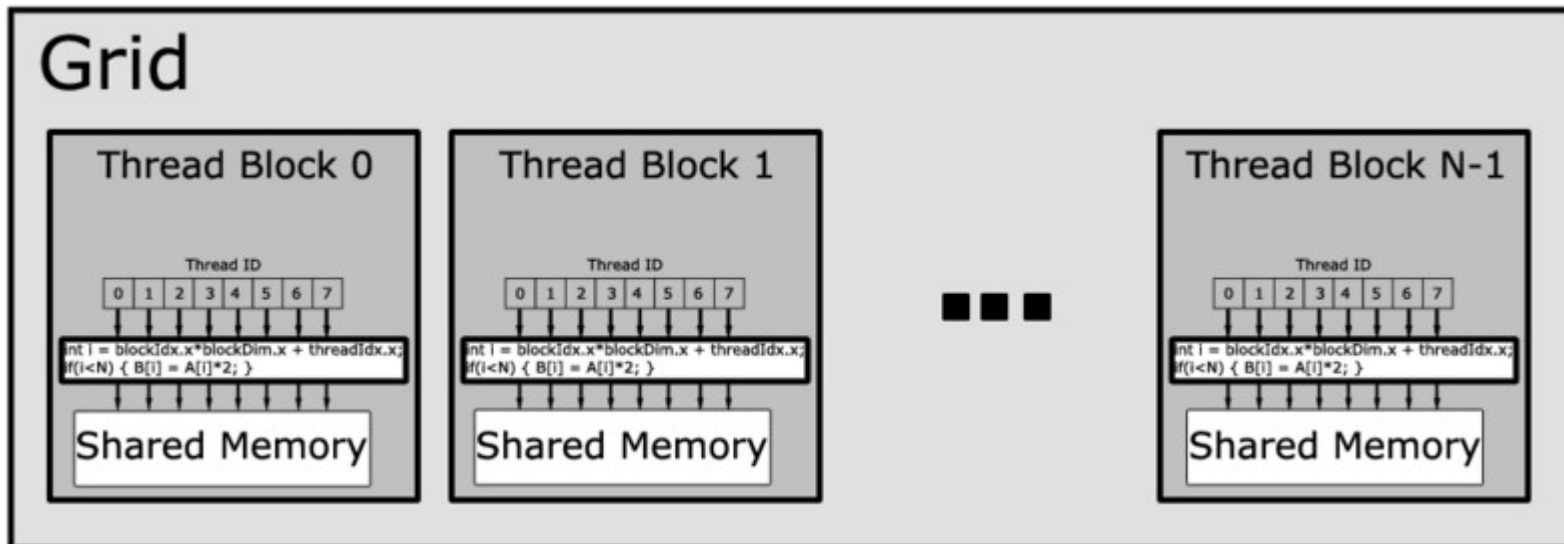
Number of threads =
blocksPerGrid *
threadsPerBlock



3. Device: Execute parallel code

```
// Device code
__global__ void VecAdd(float* A, float* B,
                      float* C, int N)
{
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if (i < N)
        C[i] = A[i] + B[i];
}
```

```
// Host code
int main()
{
    int N = ...;
    size_t size = N * sizeof(float);
    // Allocate input vectors h_A and h_B in host memory
    float* h_A = (float*)malloc(size);
    float* h_B = (float*)malloc(size);
    // Initialize input vectors
    ...
    // Allocate vectors in device memory
    float* d_A;
    cudaMalloc(&d_A, size);
}
```



4. Host: Copy results from GPU memory back to main memory

11

```
// Device code
__global__ void VecAdd(float* A, float* B,
                      float* C, int N)
{
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if (i < N)
        C[i] = A[i] + B[i];
}
```

```
// Host code
int main()
{
    int N = ...;
    size_t size = N * sizeof(float);
    // Allocate input vectors h_A and h_B in host memory
    float* h_A = (float*)malloc(size);
    float* h_B = (float*)malloc(size);
    // Initialize input vectors
    ...
    // Allocate vectors in device memory
    float* d_A;
    cudaMalloc(&d_A, size);
    float* d_B;
    cudaMalloc(&d_B, size);
    float* d_C;
    cudaMalloc(&d_C, size);
    // Copy vectors from host memory to device memory
    cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);

    // Invoke kernel
    int threadsPerBlock = 256;
    int blocksPerGrid =
        (N + threadsPerBlock - 1) / threadsPerBlock;
    VecAdd<<<blocksPerGrid, threadsPerBlock>>>(d_A, d_B, d_C, N);

    // Copy result from device memory to host memory
    // h_C contains the result in host memory
    cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);
    // Free device memory
    cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
    // Free host memory
    ....
}
```

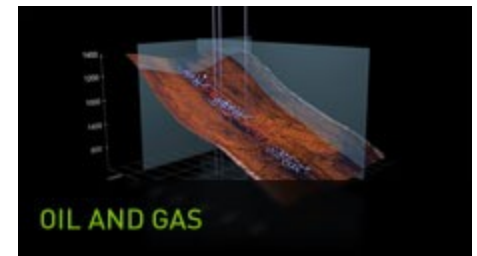
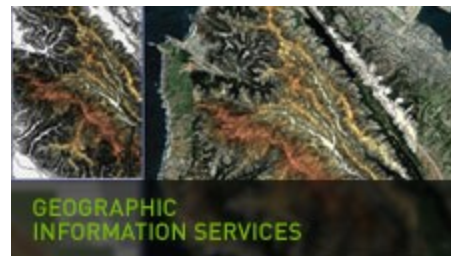
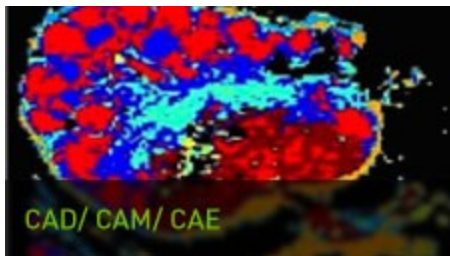
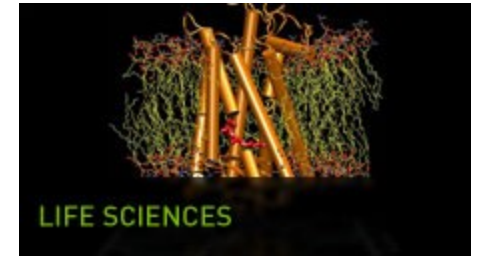
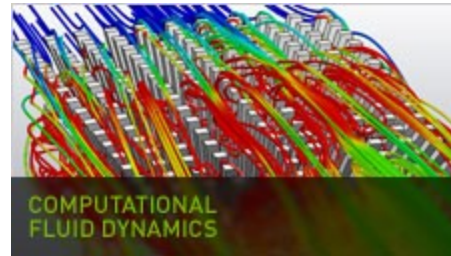
Compiling and executing GPU code

12

- `cudaSetDevice(1)` – first statement in program
- `nvcc -o vectorAdd vectorAdd.cu`
- `./vectorAdd`

GPUs are great for highly parallelizable applications

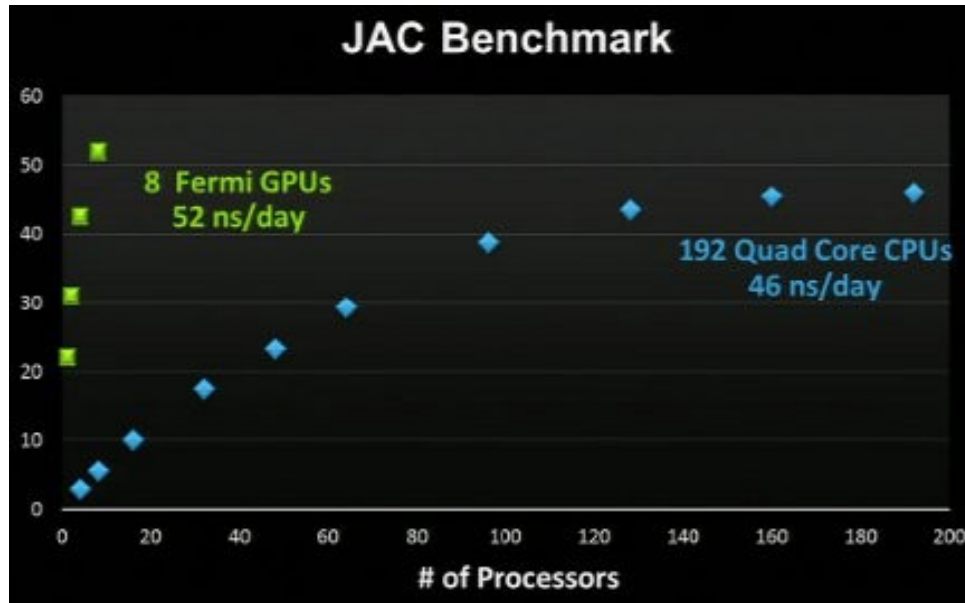
13



(www.nvidia.com)

For such applications, you can get the power of a supercomputer on your desktop

14

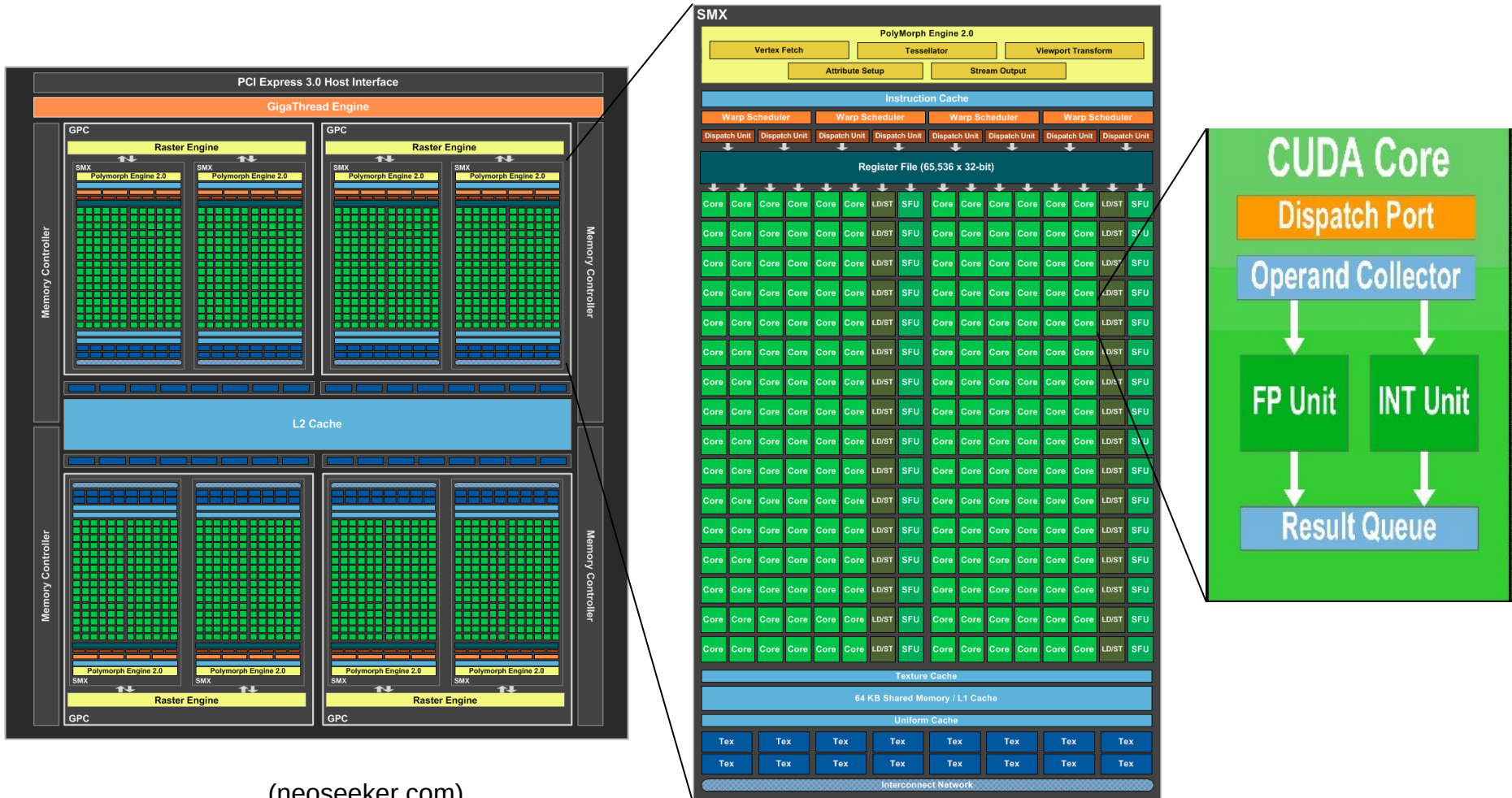


<http://hexus.net/tech/news/graphics/26553-nvidia-talks-up-plans-cuda/>



nvidia - hardware architecture

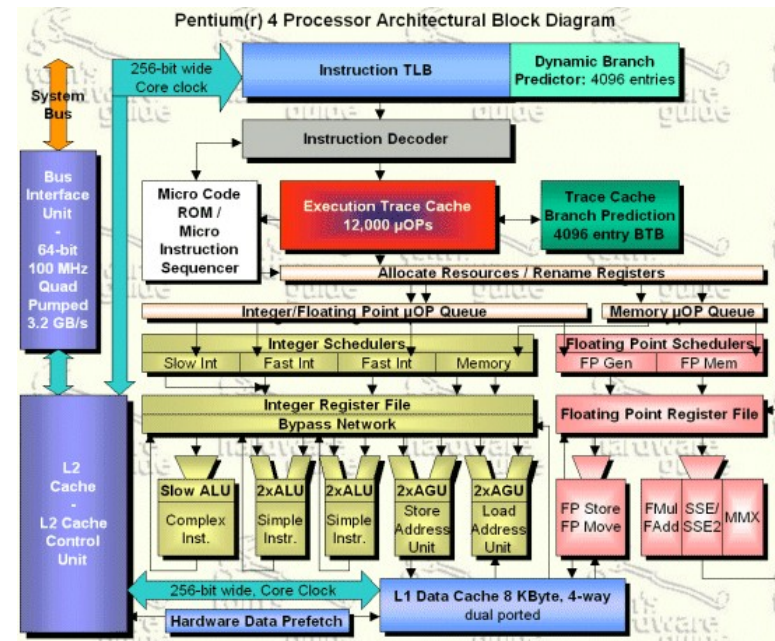
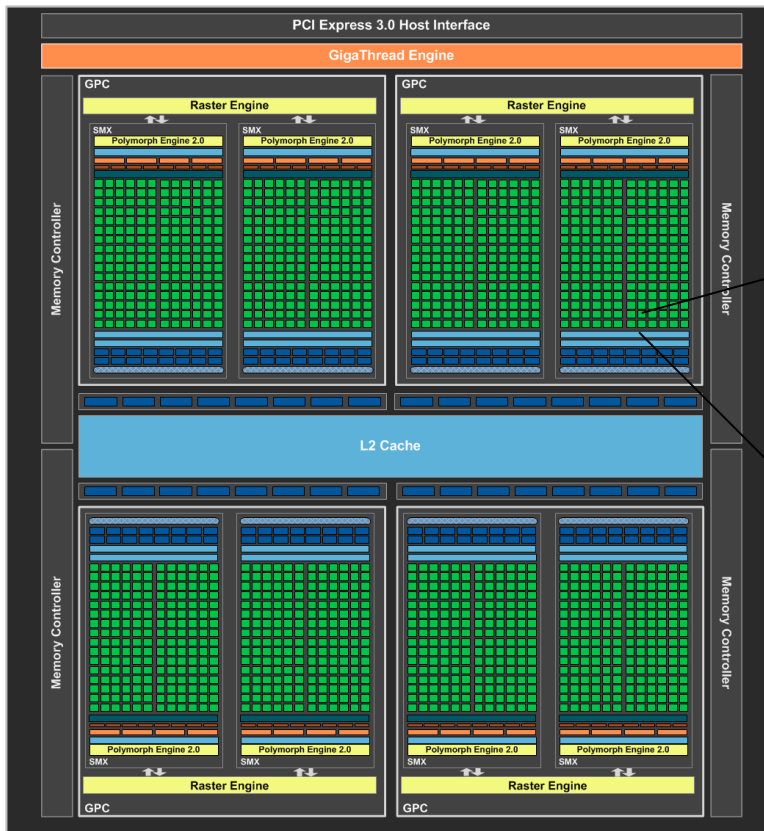
15



(neoseeker.com)

Tradeoffs between the CPU and GPU

16

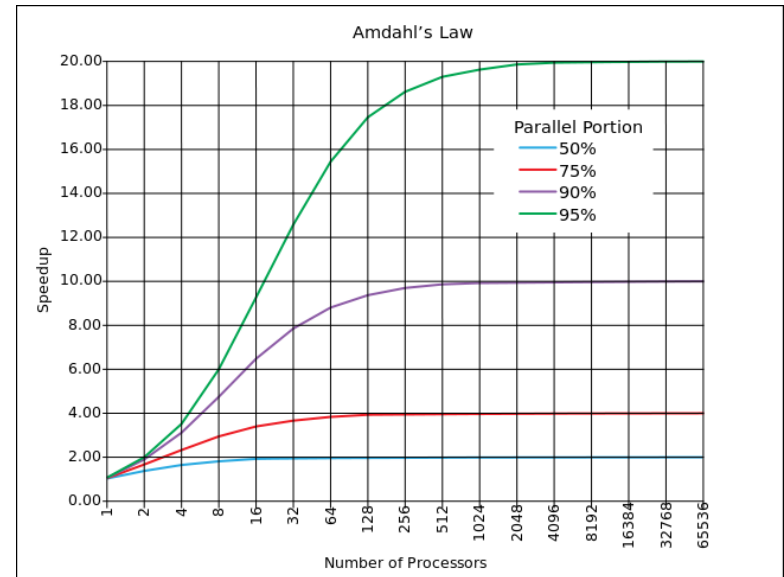


(neoseeker.com)

GPU limitations and performance considerations

17

- Not all code is highly parallelizable – Amdahl's law
- Severe penalty for branching (if-then-else)
- Limited GPU memory – increases code complexity
- Complex operations not available or much slower
- Communication between GPU and CPU is relatively slow
- Writing efficient code for the GPU is not easy



Homework assignment

Due date TBD

18

Code the Vector Addition program
(<http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>)

1. For both the CPU and GPU
2. `cudaSetDevice(1)` – first statement in program
3. Generate random values for input vectors
4. Output total of all C_i , $i = 1$ to n , for verification
5. Graphically plot execution times for different vector lengths, for CPU and GPU

GTA: Yao Zhang, yaozhang@vt.edu

McBryde office hours: Fri 9:00-11:00 am

Procedures for working on kuprin

19

1. `ssh -X CS4414S13@bio.cs.vt.edu`
2. `ssh -X kuprin`
3. `mkdir groupxx`
4. `cd groupxx`
5. Use `vi`, `gvim`, `ed`, or `gedit`, to create/edit your program
6. `cudaSetDevice(1)` – first statement in program

Ethical use of shared computer resources

20

- Do not access accounts and directories that are not yours
- Use resources only for the authorized purpose – assignments and projects
- Do not allow anyone else to use your account
- Do not try anything that might crash the machine
- Do not overload the machine so that others can not use it (denial of service)
- Follow all licensing and copyright laws
- Report any violations of these rules

If you're interested in learning more

21

- Yong Cao
 - CS4204 "Computer Graphics"
 - CS4644 "Creative Computing Studio: Video Game Design"
 - CS5984 "Advanced Computer Graphics: Parallel Computing and Visualization on GPU"
 - CS6204 "Character Animation"
- Wu Feng – has several GPU project, led the HokieSpeed project
- Alexey Onufriev – we have some projects
- CS summer seminar
- <http://developer.amd.com/pages/default.aspx>
- <http://developer.nvidia.com/category/zone/cuda-zone>
- HokieSpeed