

# Mergesort with CUDA

25.02.2020

Architecture

Mapping Cuda calls to Hardware

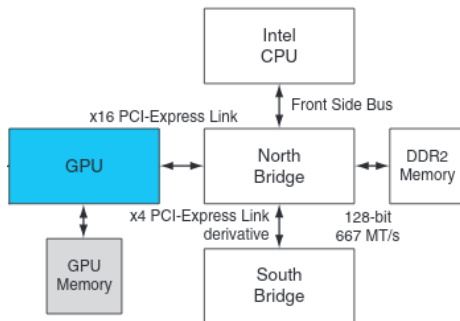
Merge

Memory Hirachy of CUDA

Merging with local memory

# Architecture

# GPU and CPU uses different memories



Init:

```
1 T* gpup;  
2 int sz = 10*sizeof(int);  
3 cudaMalloc((void**)&gpup, sz);
```

Memory Transfer:

```
1 cudaMemcpy(gpup, cpup, nBytes, cudaMemcpyHostToDevice);
```

# Ideas for Memory Management

```
1  class Storage {
2      public:
3          explicit Storage(const std::vector<T>&);
4
5      private:
6          std::vector<T> _data;
7          T* _cpu_pointer;
8          T* _gpu_pointer;
9          void initialize_gpu_memory();
10 };
```

Memory pool, takes ownership

Initializes the gpu memory as copy

Pointers for cpu/gpu locations

# Lazy Memory Sync

```
1  class Storage {
2      public:
3          T* cpu_pointer();
4          T* gpu_pointer();
5          const T* cpu_pointer_const();
6          const T* gpu_pointer_const();
7
8      private:
9          std::string head; \\\ head = {CPU, GPU, SYNC}
10         void sync_to_cpu();
11         void sync_to_gpu();
12     };
```

## Example:

After Initialization: head = SYNC

Access, gpu\_pointer\_const(): head remains

Access cpu\_pointer(): head = CPU

Acess gpu\_pointer(): sync\_to\_cpu(); head = SYNC

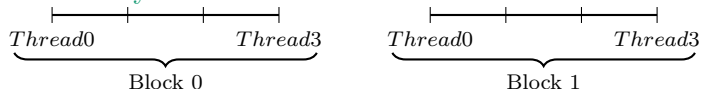
# Mapping Cuda calls to Hardware

# Launching CUDA threads

## Cuda Program

```
1 dim3 Grid(2)
2 dim3 Block(4)
3 add_kernel<<<Grid, Block>>>(...)
```

## Thread Layout:



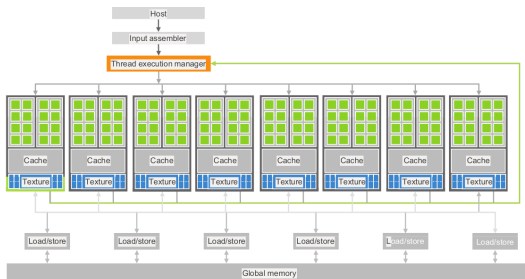
## Addition:

```
1 __global__
2 add_kernel(float* A, float* B, float* C, int n) {
3     int i = blockDim.x * blockIdx.x + threadIdx.x;
4     if (i < n) {
5         C[i] = A[i] + B[i];
6     }
7 }
```



# Software/Hardware connection

**Programmer:** Specifies number of blocks and threads

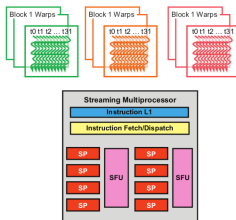


**Dispatch:**

1. **Blocks to SM:** Block 1  $\rightarrow$  SM1, Block 2  $\rightarrow$  SM2 ...

# Software/Hardware connection

**Programmer:** Specifies number of blocks and threads

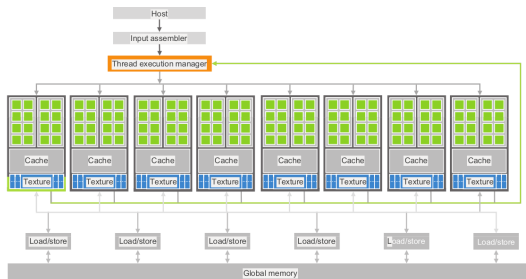


**Dispatch:**

1. **Blocks to SM:** Block 1  $\rightarrow$  SM1, Block 2  $\rightarrow$  SM2 ...
2. **Grouping:** Wrappend threads are allocated

# Software/Hardware connection

**Programmer:** Specifies number of blocks and threads



**Dispatch:**

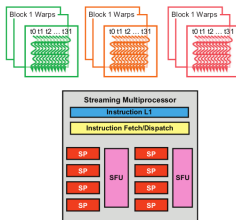
1. **Blocks to SM:** Block 1  $\rightarrow$  SM1, Block 2  $\rightarrow$  SM2 ...
2. **Grouping:** Wrappend threads are allocated

**Stengths:**

**Quantity:** No branch prediction/ chaches, just cores

# Software/Hardware connection

**Programmer:** Specifies number of blocks and threads



**Dispatch:**

1. **Blocks to SM:** Block 1  $\rightarrow$  SM1, Block 2  $\rightarrow$  SM2 ...
2. **Grouping:** Wrappend threads are allocated

**Stengths:**

**Quantity:** No branch prediction/ chaches, just cores

**Latency Hiding:** If wrap 1 stalls, scheduler starts wrap 2

Merge

# Basic Merge Operation

A = 5 7 8 9 12 14 15 16

B = 1 2 3 4 6 10 11 13

C = ? ? ? ? ? ? ? ... ?

```
1 void merge(T* a, T* b, T* c, int sz_a, int sz_b) {
2     int i = 0, j = 0, k = 0;
3     while (k < sz_a + sz_b) {
4         if (i == sz_a)
5             c[k++] = b[j++];
6         else if (j == sz_b)
7             c[k++] = a[i++];
8         else if (a[i] <= b[j])
9             c[k++] = a[i++];
10        else
11            c[k++] = b[j++];
12    }
13 }
```

Merge

# How to split A and B to spawn many threads?

Example:

$$A = 0 \quad 0 \quad 0 \quad 0$$

$$B = 1 \quad 1 \quad 1 \quad 1$$

$$C = ? \quad ? \quad ? \quad ? \quad ? \quad ? \quad ? \quad ?$$

# How to split A and B to spawn many threads?

Example:

$$A = 0 \quad 0 \quad 0 \quad 0$$

$$B = 1 \quad 1 \quad 1 \quad 1$$

$$C = ? \quad ? \quad ? \quad ? \quad ? \quad ? \quad ? \quad ?$$

Naive: 2 Threads, half A and B



# How to split A and B to spawn many threads?

Example:

$$A = 0 \quad 0 \quad 0 \quad 0$$

$$B = 1 \quad 1 \quad 1 \quad 1$$

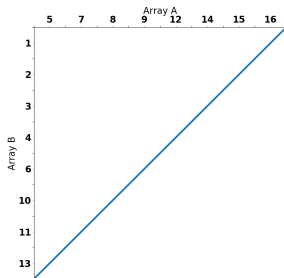
$$C = ? \quad ? \quad ? \quad ? \quad ? \quad ? \quad ? \quad ?$$

Naive: 2 Threads, half A and B    Result:

## Instead: How to split arrays

NEED TO CITE PAPER!

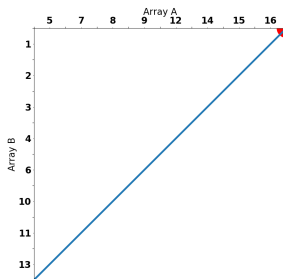
**How to split two arrays in equal chunks?**



## Instead: How to split arrays

NEED TO CITE PAPER!

**One feasible split: Array A to Thread 1, Array B to**

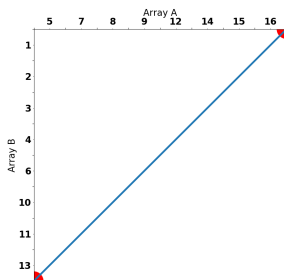


**Thread 2**

## Instead: How to split arrays

NEED TO CITE PAPER!

**Another feasible split: Array B to Thread 1, Array A**

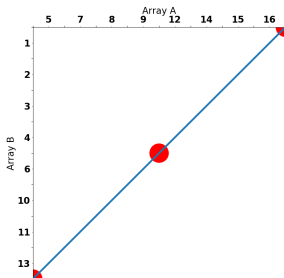


**to Thread 2**

## Instead: How to split arrays

NEED TO CITE PAPER!

**Another split (seen before): Thread 1 gets half of A**



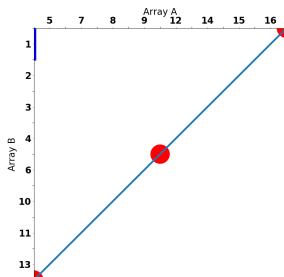
**and B**

**Summary:** All allocations along vertical line split work equally!

## Instead: How to split arrays

NEED TO CITE PAPER!

**Mergepath: Define the optimal split: One Element of B**



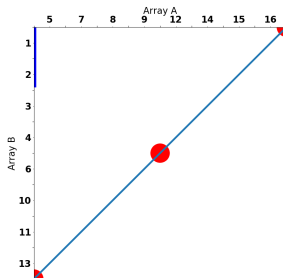
**Summary:** All allocations along vertical line split work equally!

**Mergepath:** Vertical move: pick array B; horizontal move: Array

## Instead: How to split arrays

NEED TO CITE PAPER!

**Mergepath: Define the optimal split: Another Element**

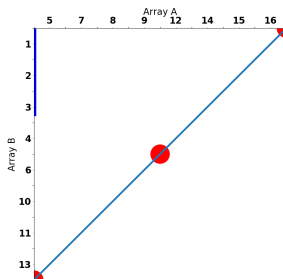


of B

# Instead: How to split arrays

NEED TO CITE PAPER!

**Mergepath: Define the optimal split: Another Element**



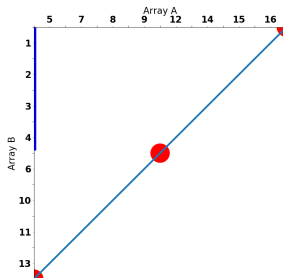
of B



## Instead: How to split arrays

NEED TO CITE PAPER!

**Mergepath: Define the optimal split: Another Element**

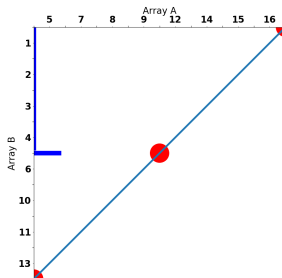


of B

## Instead: How to split arrays

NEED TO CITE PAPER!

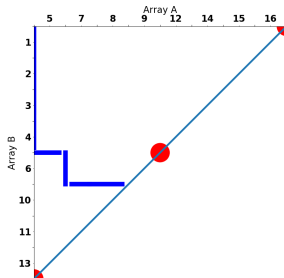
**Mergepath: Define the optimal split: Pick first element**



of A

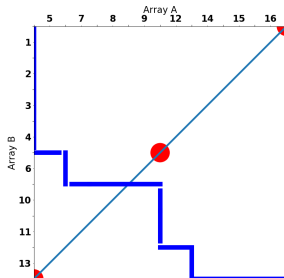
# Instead: How to split arrays

NEED TO CITE PAPER!



# Instead: How to split arrays

NEED TO CITE PAPER!



# Computation Procedure

```
1  __global__
2  void parallelMerge(int* a, int sz_a, int* b, int sz_b,
3                    int* c, int length)
4  {
5      int diag = threadIdx.x * length;
6      int a_start = mergepath(a, sz_a, b, sz_b, diag);
7      int b_start = diag - a_start;
8      merge(a, a_start, sz_a, b, b_start, sz_b, c, diag, length);
9  }
```

Each thread works on one part

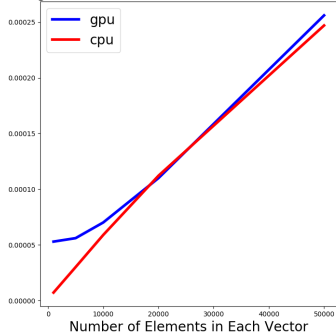
Thread calculates  $A_{lower}$  mergepath

Obtain  $B_{lower}$  as difference

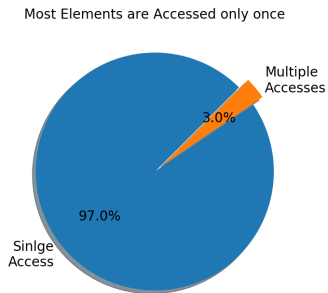
merges two sub-arrays

# Problem: Slow as a Snail

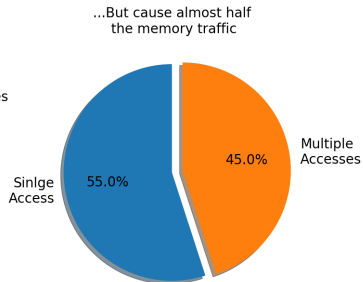
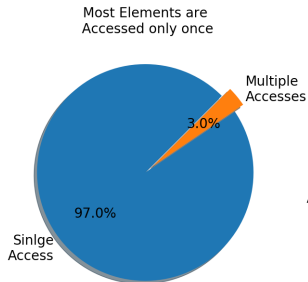
Merge time vs Number of Elements for CPU and GPU



Reason: So much global memory access



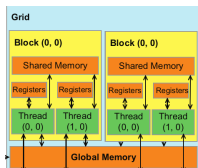
# Reason: So much global memory access





# Memory Hirachy of CUDA

# The different memories and their sizes



## Visibility:

Global Memory: Accessed by all threads

Shared Memory: Private for Block

Registers: Private for Thread

## Size:

Global Memory: 2 GB

Shared Memory: Total: 192 KB, per block: 48 KB

## Latency:

Global Memory: 8 GB/s

Shared Memory: 80 GB/s

# Merging with local memory

# General Idea

**Problem:** Block shared memory: 48 Kb -> need to make blocks small enough

**Solution:**

- Determine large but small enough chunks of arrays to load into blocks

- Load into shared memory

- Each thread determines its range (as before), from shared memory

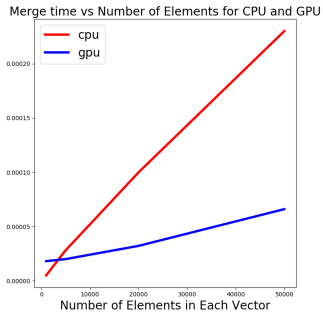
- Each thread merges subset (as before)

```

1  __global__
2  void parallelMerge(const int* a, int sz_a, const int* b, int
3                      int* c, int* boundaries, int length,
4                      int size_shared) {
5      extern __shared__ int shared[];
6      __shared__ int block_ranges[4];
7      ranges(block_ranges, sz_a, sz_b, boundaries);
8      loadtodevice(a, sz_a, b, sz_b, block_ranges, shared);
9      int diag = threadIdx.x * length;
10     if (diag < block_ranges[2] + block_ranges[3]) {
11         int a_start =
12             mergepath(shared, block_ranges[2], &shared[block
13                     block_ranges[3], diag);
14         int b_start = diag - a_start;
15         merge(shared, a_start, block_ranges[2], &shared[blo
16             b_start, block_ranges[3], c, diag + blockIdx
17             length);
18     }
19 }

```

# Show the results



Merging with local memory