

Mergesort with CUDA

25.02.2020

GPUs is fast and easily programmable

Fast: See number of flops **Easy:** CUDA: Extension to C-language, called from C/C++

Simple Program and Big Architecture Ideas

Merge

Memory Hirachy of CUDA

Merging with local memory

Simple Program and Big Architecture Ideas

A Simple Cuda Programm

1: Call Cuda Program from C/C++:

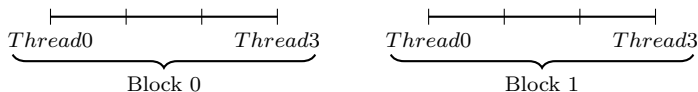
```
1  dim3 Grid(2)
2  dim3 Block(4)
3  add_arrays<<<Grid, Block>>>(...)
```

A Simple Cuda Programm

1: Call Cuda Program from C/C++:

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

2: Blocks & Grid spawn many threads on GPU:

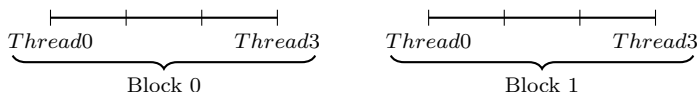


A Simple Cuda Programm

1: Call Cuda Program from C/C++:

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

2: Blocks & Grid spawn many threads on GPU:



3: Cuda Code is a C-Extension:

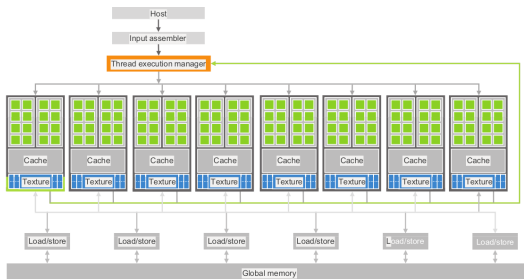
```
1 __global__
2 add_arrays(float* A, float* B, float* C, int n) {
3     int i = blockDim.x * blockIdx.x + threadIdx.x;
4     if (i < n) {
5         C[i] = A[i] + B[i];
6     }
```

Software/Hardware interface

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

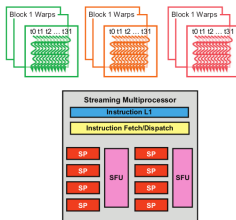
1. Software: Specify Blocks and Threads

Software/Hardware interface



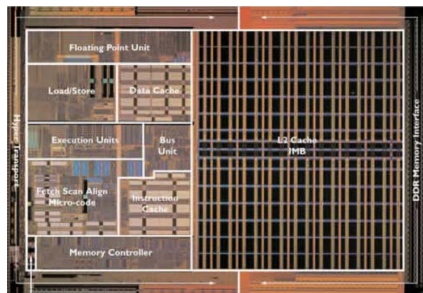
1. **Software:** Specify Blocks and Threads
2. **Hardware:**
 1. **Blocks to SM:** Block 1 \rightarrow SM1, Block 2 \rightarrow SM2 ...

Software/Hardware interface



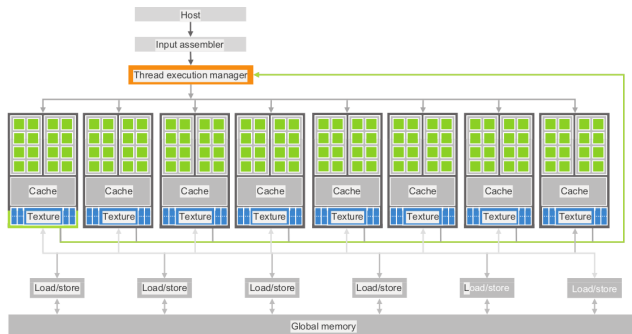
1. **Software:** Specify Blocks and Threads
2. **Hardware:**
 1. **Blocks to SM:** Block 1 \rightarrow SM1, Block 2 \rightarrow SM2 ...
 2. **Threads to Processors:** Wrappend threads are allocated

Reasons for strong GPU performance



CPU: Chip full with fancy Branch Prediction, Caches, etc.

Reasons for strong GPU performance

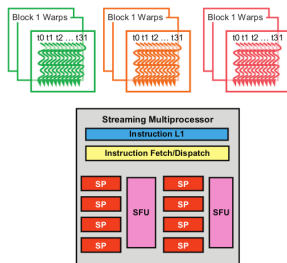


CPU: Chip full with fancy Branch Prediction, Caches, etc.

GPU:

1. **Simplicity:** Just Cores on chip

Reasons for strong GPU performance



CPU: Chip full with fancy Branch Prediction, Caches, etc.

GPU:

1. **Simplicity:** Just Cores on chip
2. **Multithreading:** Wrap 1 in a long-latency op: Start Wrap 2

Merge

Serial Merge

$$A = \begin{bmatrix} 5 & 7 & 8 & 12 \end{bmatrix}$$

$$B = \begin{bmatrix} 3 & 4 & 6 & 10 \end{bmatrix}$$

$$C = \begin{bmatrix} ? & ? & ? & ? & ? & ? & ? & ? \end{bmatrix}$$

```
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 }
```

Problem: Merge complexity $\mathcal{O}(n)$

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:

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

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

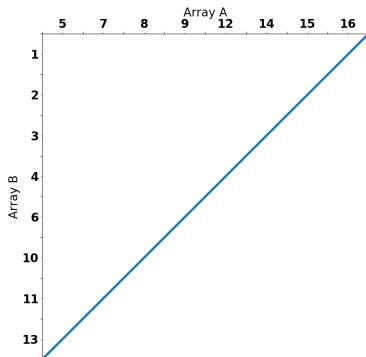
$$C = ? \quad ? \quad ? \quad ? \quad ? \quad ? \quad ? \quad ?$$

Naive: 2 Threads, half A and B Result:

Instead: How to split arrays

Odeh et al. (2012)

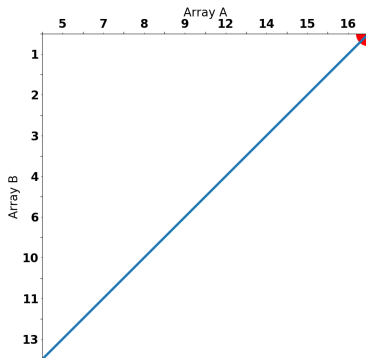
How to split two arrays in equal chunks?



Instead: How to split arrays

Odeh et al. (2012)

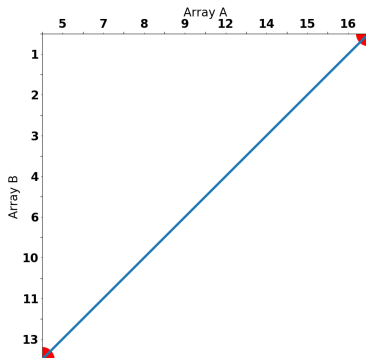
Feasible split: Array A to Thread 1, B to Thread 2



Instead: How to split arrays

Odeh et al. ([2012](#))

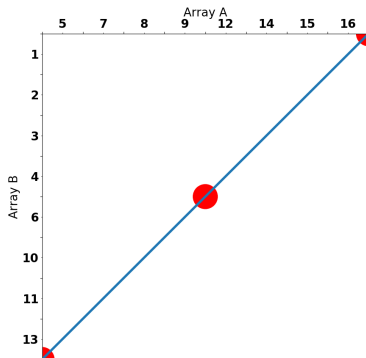
Another split: Array B to Thread 1, A to Thread 2



Instead: How to split arrays

Odeh et al. (2012)

Another (as before): Thread 1 gets half of A and B

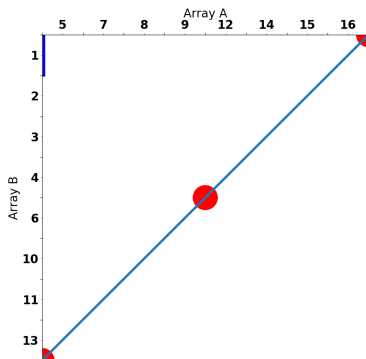


Summary: All allocations along vertical line split work equally!

Instead: How to split arrays

Odeh et al. (2012)

Mergepath: Optimal split: One Element of B



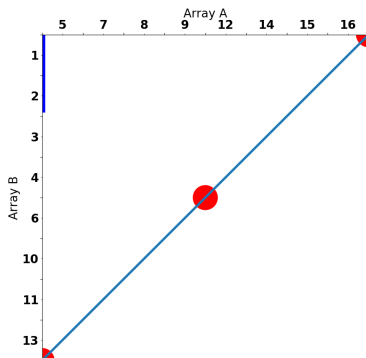
Summary: All allocations along vertical line split work equally!

Mergpath: Vertical move: pick array B; horizontal move: Array
Merge

Instead: How to split arrays

Odeh et al. (2012)

Mergepath: Optimal split: Another Element of B



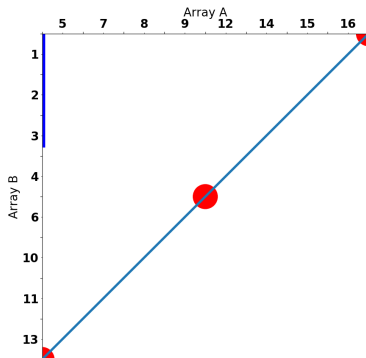
Summary: All allocations along vertical line split work equally!

Mergpath: Vertical move: pick array B; horizontal move: Array
Merge

Instead: How to split arrays

Odeh et al. (2012)

Mergepath: Optimal split: Another Element of B



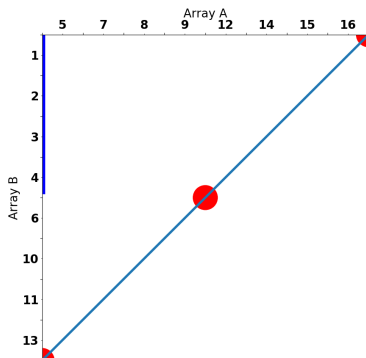
Summary: All allocations along vertical line split work equally!

Mergpath: Vertical move: pick array B; horizontal move: Array
Merge

Instead: How to split arrays

Odeh et al. (2012)

Mergepath: Optimal split: Another Element of B



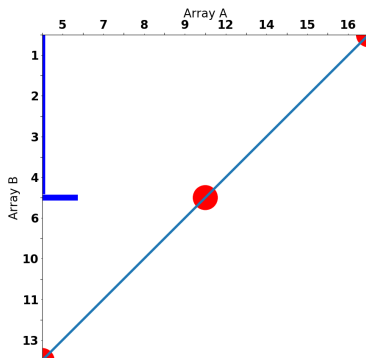
Summary: All allocations along vertical line split work equally!

Mergpath: Vertical move: pick array B; horizontal move: Array
Merge

Instead: How to split arrays

Odeh et al. (2012)

Mergepath: Optimal split: First element of A



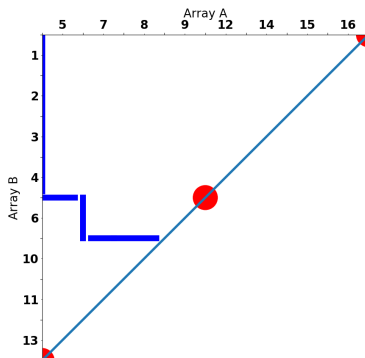
Summary: All allocations along vertical line split work equally!

Mergpath: Vertical move: pick array B; horizontal move: Array
Merge

Instead: How to split arrays

Odeh et al. (2012)

Mergepath: Optimal split!



Split index of A: 3

Split index of B: 5

Merge

Merge with Cuda

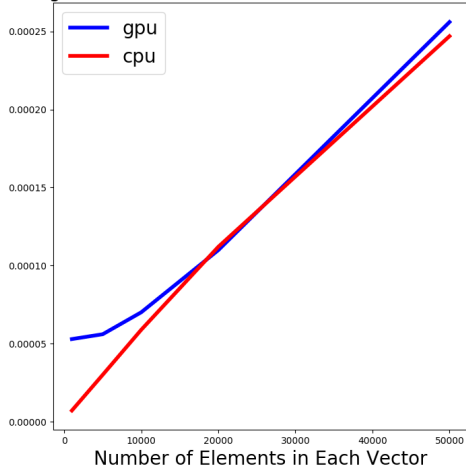
```
1  __global__
2  void parallelMerge(int* a, int sz_a, int* b, int sz_b,
3                    int* c, int length)
4  {
5      int diag = threadIdx.x * length;
6      int a_split = mergepath(a, sz_a, b, sz_b, diag);
7      int b_split = diag - a_split;
8      merge(a, a_split, sz_a, b, b_split, sz_b,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

Each identifies split indices

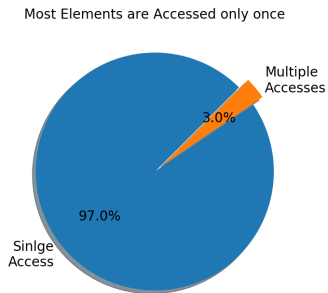
Split indices suffice to merge two sub-arrays into c

Problem: Slow as a Snail

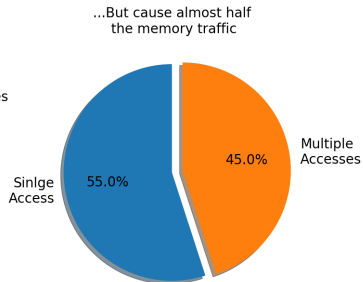
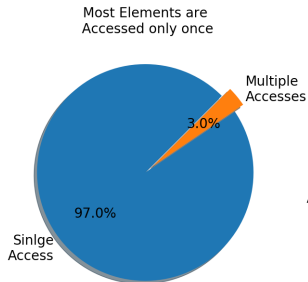
Merge time vs Number of Elements for CPU and GPU



Reason: So much global memory access

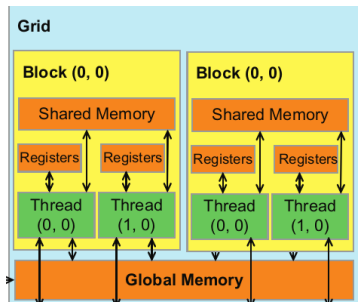


Reason: So much global meory access



Memory Hirachy of CUDA

Different Memories and Size



Visibility:

Global \Leftrightarrow Shared: All threads \Leftrightarrow Private for Block

Size:

Global \Leftrightarrow Shared: 2 GB \Leftrightarrow Total: 192 KB, per block: 48 KB

Latency:

Global \Leftrightarrow Shared: 8 GB/s \Leftrightarrow 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

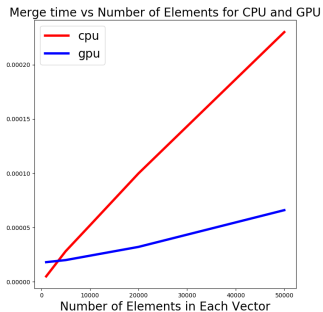
- Each thread merges subset (as before)

```

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





Saher Odeh et al. “Merge Path - Parallel Merging Made Simple”. In: *2012 IEEE 26th International Parallel and Distributed Processing Symposium Workshops PhD Forum*. ISSN: null. May 2012, pp. 1611–1618 (cit. on pp. 19–28).