CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

# CUDA C: performance measurement and memory

Will Landau

Iowa State University

October 14, 2013

# Outline

Timing kernels on the GPU

Memory

# Outline

Timing kernels on the GPU

Memory

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

# Measuring CPU time

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
1  #include <stdio.h>
2  #include <time.h>
3
4  int main(){
5    float elapsedTime;
6    clock_t start = clock();
7
8    // SOME CPU CODE YOU WANT TO TIME
9
10   elapsedTime = ((double) clock() - start) /
         CLOCKS_PER_SEC;
11
12   pritnf("CPU time elapsed: %f seconds \n",
         elapsedTime);
13   return 0;
14 }
```

# Events

- ▶ **Event**: a time stamp on the GPU
- ▶ Use events to measure GPU execution time.
- ▶ `time.cu`:

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
1  #include <stdlib.h>
2  #include <stdio.h>
3  #include <cuda.h>
4  #include <cuda_runtime.h>
5
6  int main(){
7    float    elapsedTime;
8    cudaEvent_t start, stop;
9    cudaEventCreate(&start);
10   cudaEventCreate(&stop);
11   cudaEventRecord( start, 0 );
12
13   // SOME GPU WORK YOU WANT TIMED HERE
14
15   cudaEventRecord( stop, 0 );
16   cudaEventSynchronize( stop );
17   cudaEventElapsedTime( &elapsedTime, start, stop );
18   cudaEventDestroy( start );
19   cudaEventDestroy( stop );
20   printf("GPU Time elapsed: %f milliseconds\n", elapsedTime);
21 }
```

- ▶ GPU time and CPU time must be measured separately.

# Example: `pairwise_sum_timed.cu`

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <math.h>
4   #include <time.h>
5   #include <unistd.h>
6   #include <cuda.h>
7   #include <cuda_runtime.h>
8
9   /* This program computes the sum of the elements of
10   * vector v using the pairwise (cascading) sum algorithm. */
11
12  #define N 1024 // length of vector v. MUST BE A POWER OF 2!!!
13
14  // Fill the vector v with n random floating point numbers.
15  void vfill(float* v, int n){
16    int i;
17    for(i = 0; i < n; i++){
18      v[i] = (float) rand() / RAND_MAX;
19    }
20  }
21
22  // Print the vector v.
23  void vprint(float* v, int n){
24    int i;
25    printf("v = \n");
26    for(i = 0; i < n; i++){
27      printf("%7.3f\n", v[i]);
28    }
29    printf("\n");
30  }
```

# Example: pairwise_sum_timed.cu

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
31  // Pairwise−sum the elements of vector v and store the result in v
        [ 0 ] .
32  __global__ void psum(float *v){
33    int t = threadIdx.x;  // Thread index.
34    int n = blockDim.x;  // Should be half the length of v.
35
36    while (n != 0) {
37      if (t < n)
38        v[t] += v[t + n];
39      __syncthreads();
40      n /= 2;
41    }
42  }
43
44  // Linear sum the elements of vector v and return the result
45  float lsum(float *v, int len){
46    float s = 0;
47    int i;
48    for(i = 0; i < len; i++){
49      s += v[i];
50    }
51    return s;
52  }
```

# Example: pairwise_sum_timed.cu

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
54  int main (void){
55    float *v_h, *v_d; // host and device copies of our vector,
             respectively
56
57    // dynamically allocate memory on the host for v_h
58    v_h = (float*) malloc(N * sizeof(*v_h));
59
60    // dynamically allocate memory on the device for v_d
61    cudaMalloc ((float**) &v_d, N *sizeof(*v_d));
62
63    // Fill v_h with N random floating point numbers.
64    vfill(v_h, N);
65
66    // Print v_h to the console
67    // vprint(v_h, N);
68
69    // Write the contents of v_h to v_d
70    cudaMemcpy( v_d, v_h, N * sizeof(float), cudaMemcpyHostToDevice );
71
72    // compute the linear sum of the elements of v_h on the CPU and
             return the result
73    // also, time the result.
74    clock_t start = clock();
75    float s = lsum(v_h, N);
```

# Example: pairwise_sum_timed.cu

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
76    float elapsedTime = ((float) clock() - start) / CLOCKS_PER_SEC;
77    printf("Linear Sum = %7.3f, CPU Time elapsed: %f seconds\n", s,
          elapsedTime);
78
79    // Compute the pairwise sum of the elements of v_d and store the
          result in v_d[0].
80    // Also, time the computation.
81
82    float    gpuElapsedTime;
83    cudaEvent_t gpuStart, gpuStop;
84    cudaEventCreate(&gpuStart);
85    cudaEventCreate(&gpuStop);
86    cudaEventRecord( gpuStart, 0 );
87
88    psum<<< 1, N/2 >>>(v_d);
89
90    cudaEventRecord( gpuStop, 0 );
91    cudaEventSynchronize( gpuStop );
92    cudaEventElapsedTime( &gpuElapsedTime, gpuStart, gpuStop ); // time
          in milliseconds
93    cudaEventDestroy( gpuStart );
94    cudaEventDestroy( gpuStop );
95
96    // Write the pairwise sum, v_d[0], to v_h[0].
97    cudaMemcpy( v_h, v_d, sizeof(float), cudaMemcpyDeviceToHost );
```

# Example: `pairwise_sum_timed.cu`

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
98      // Print the pairwise sum.
99      printf("Pairwise Sum = %7.3f, GPU Time elapsed: %f seconds\n", v_h
           [0], gpuElapsedTime/1000.0);
100
101     // Free dynamically−allocated host memory
102     free(v_h);
103
104     // Free dynamically−allocated device memory
105     cudaFree(&v_d);
106 }
```

▶ Output:

```
1 > nvcc pairwise_sum_timed.cu −o pairwise_sum_timed
2 > ./pairwise_sum_timed
3 Linear Sum = 518.913, CPU Time elapsed: 0.000000 seconds
4 Pairwise Sum = 518.913, GPU Time elapsed: 0.000037 seconds
```

# Outline

CUDA C:
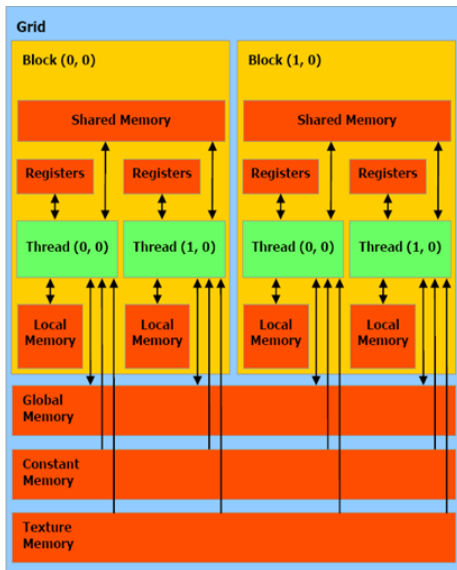performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

Timing kernels on the GPU

Memory

# Types of memory

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

# What happens in `myKernel<<<2, 2>>>(b, t)`?

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
1  __global__ void myKernel(int *b_global, int *
       t_global){
2
3    __shared__ int t;
4    __shared__ int b;
5
6    int b_local, t_local;
7
8    *t_global = threadIdx.x;
9    *b_global = blockIdx.x;
10
11     t_shared = threadIdx.x;
12     b_shared = blockIdx.x;
13
14     t_local = threadIdx.x;
15     b_local = blockIdx.x;
16 }
```

# At the end of myKernel<<<4, 7>>>(b, t)...

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

▶ b_local and t_local are in local memory (or registers), so each thread gets a copy.

| (block, thread) | (0, 0) | (0, 1) | (1, 0) | (1, 1) |
|---|---|---|---|---|
| b_local | 0 | 0 | 1 | 1 |
| t_local | 0 | 1 | 0 | 1 |

▶ b_shared and t_shared are in shared memory, so each block gets a copy.

| (block, thread) | (0, 0) | (0, 1) | (1, 0) | (1, 1) |
|---|---|---|---|---|
| b_shared | 0 | 0 | 1 | 1 |
| t_shared | ? | ? | ? | ? |

  ▶ ? = last thread in its block to write to t_shared.

# At the end of myKernel<<<4, 7>>>(b, t)...

▶ b_global and t_global point to global memory, so there is only one copy.

| (block, thread) | (0, 0) | (0, 1) | (1, 0) | (1, 1) |
|-----------------|--------|--------|--------|--------|
| *b_global       | ??     | ??     | ??     | ??     |
| *t_global       | ?      | ?      | ?      | ?      |

  ▶ ? = last thread in its block to write to *t_global.
  ▶ ?? = block of the last thread to write to *b_global.

# Example: dot product

$$a \bullet b = (a_0, \ldots, a_{15}) \bullet (b_0, \ldots, b_{15}) = a_0 \cdot b_0 + \cdots + a_{15} \cdot b_{15}$$

1. In this example, spawn 2 blocks and 4 threads per block.
2. Give each block a subvector of $a$ and an analogous subvector of $b$.
   ▶ Block 0:

   $$(a_0, a_1, a_2, a_3, \ a_8, a_9, a_{10}, a_{11})$$
   $$(b_0, b_1, b_2, b_3, \ b_8, b_9, b_{10}, b_{11})$$

   ▶ Block 1:

   $$(a_4, a_5, a_6, a_7, \ a_{12}, a_{13}, a_{14}, a_{15})$$
   $$(b_4, b_5, b_6, b_7, \ b_{12}, b_{13}, b_{14}, b_{15})$$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

## Example: dot product

3. Create an array, `cache`, in shared memory:

▶ Block 0:

$$\text{cache}[0] = a_0 \cdot b_0 + a_8 \cdot b_8$$
$$\text{cache}[1] = a_1 \cdot b_1 + a_9 \cdot b_9$$
$$\text{cache}[2] = a_2 \cdot b_2 + a_{10} \cdot b_{10}$$
$$\text{cache}[3] = a_3 \cdot b_3 + a_{11} \cdot b_{11}$$

▶ Block 1:

$$\text{cache}[0] = a_4 \cdot b_4 + a_{12} \cdot b_{12}$$
$$\text{cache}[1] = a_5 \cdot b_5 + a_{13} \cdot b_{13}$$
$$\text{cache}[2] = a_6 \cdot b_6 + a_{14} \cdot b_{14}$$
$$\text{cache}[3] = a_7 \cdot b_7 + a_{15} \cdot b_{15}$$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

## Example: dot product

4. Compute the pairwise sum of cache in each block and write it to cache[0]

   ▶ Block 0:

   $$
   \begin{aligned}
   \text{cache}[0] = {} & a_0 \cdot b_0 + a_8 \cdot b_8 \\
   & + a_1 \cdot b_1 + a_9 \cdot b_9 \\
   & + a_2 \cdot b_2 + a_{10} \cdot b_{10} \\
   & + a_3 \cdot b_3 + a_{11} \cdot b_{11}
   \end{aligned}
   $$

   ▶ Block 1:

   $$
   \begin{aligned}
   \text{cache}[0] = {} & a_4 \cdot b_4 + a_{12} \cdot b_{12} \\
   & + a_5 \cdot b_5 + a_{13} \cdot b_{13} \\
   & + a_6 \cdot b_6 + a_{14} \cdot b_{14} \\
   & + a_7 \cdot b_7 + a_{15} \cdot b_{15}
   \end{aligned}
   $$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

# Example: dot product

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

5. Compute an array, partial_c in global memory:

partial_c[0] = cache[0] from block 0

partial_c[1] = cache[0] from block 1

6. The pairwise sum of partial_c is the final answer.

# dot_product.cu

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
1   #include "../common/book.h"
2   #include <stdio.h>
3   #include <stdlib.h>
4   #define imin(a,b) (a<b?a:b)
5
6   const int N = 32 * 1024;
7   const int threadsPerBlock = 256;
8   const int blocksPerGrid = imin( 32, (N+threadsPerBlock-1) /
        threadsPerBlock );
9
10  __global__ void dot( float *a, float *b, float *partial_c ) {
11
12    __shared__ float cache[threadsPerBlock];
13    int tid = threadIdx.x + blockIdx.x * blockDim.x;
14    int cacheIndex = threadIdx.x;
15    float temp = 0;
16
17    while (tid < N) {
18      temp += a[tid] * b[tid];
19      tid += blockDim.x * gridDim.x;
20    }
21
22    // set the cache values
23    cache[cacheIndex] = temp;
```

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

dot<<2,4>>(a, b, c)

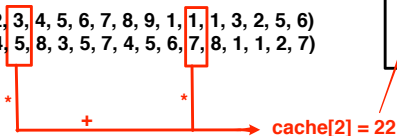blockDim.x = 4
gridDim.x = 2

**a = (1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 1, 1, 3, 2, 5, 6)**
**b = (2, 4, 5, 8, 3, 5, 7, 4, 5, 6, 7, 8, 1, 1, 2, 7)**

| Block 0 | Block 1 |
|---|---|
| cache[0] = | cache[0] = |
| cache[1] = | cache[1] = |
| cache[2] = | cache[2] = |
| cache[3] = | cache[3] = |

# dot<<<2, 4>>>(a, b, c) with $N = 16$

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau
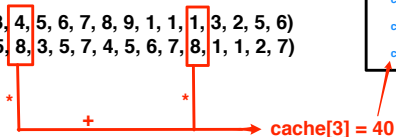
Timing kernels on
the GPU

Memory



dot<<2,4>>(a, b, c)

blockDim.x = 4
gridDim.x = 2

**threadIdx.x = 1**
**blockIdx.x = 0**

a = (1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 1, 1, 3, 2, 5, 6)
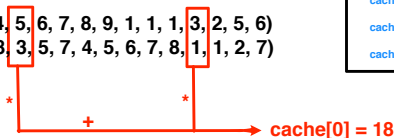b = (2, 4, 5, 8, 3, 5, 7, 4, 5, 6, 7, 8, 1, 1, 2, 7)

\*                      \*

+

**cache[1] = 14**

| Block 0 | Block 1 |
|---------|---------|
| cache[0] = 47 | cache[0] = |
| cache[1] = 14 | cache[1] = |
| cache[2] = | cache[2] = |
| cache[3] = | cache[3] = |

# dot<<<2, 4>>>(a, b, c) with $N = 16$

dot<<2,4>>(a, b, c)

blockDim.x = 4
gridDim.x = 2

**threadIdx.x = 2**
**blockIdx.x = 0**

a = (1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 1, 1, 3, 2, 5, 6)
b = (2, 4, 5, 8, 3, 5, 7, 4, 5, 6, 7, 8, 1, 1, 2, 7)

*          *

+

**cache[2] = 22**

| **Block 0** | **Block 1** |
|---|---|
| cache[0] = 47 | cache[0] = |
| cache[1] = 14 | cache[1] = |
| cache[2] = 22 | cache[2] = |
| cache[3] = | cache[3] = |

# dot<<<2, 4>>>(a, b, c) with $N = 16$

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

dot<<2,4>>(a, b, c)

blockDim.x = 4
gridDim.x = 2

threadIdx.x = 0
blockIdx.x = 1

a = (1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 1, 1, 3, 2, 5, 6)
b = (2, 4, 5, 8, 3, 5, 7, 4, 5, 6, 7, 8, 1, 1, 2, 7)

*                    *

+           cache[0] = 18

| Block 0 | Block 1 |
|---|---|
| cache[0] = 47 | cache[0] = 18 |
| cache[1] = 14 | cache[1] = |
| cache[2] = 22 | cache[2] = |
| cache[3] = 40 | cache[3] = |

# dot<<<2, 4>>>(a, b, c) with $N = 16$

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
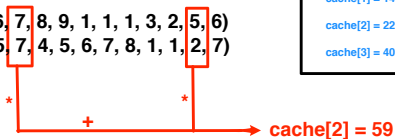the GPU

Memory
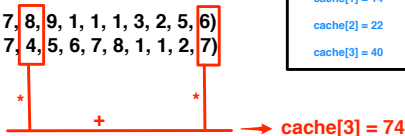
dot<<2,4>>(a, b, c)

blockDim.x = 4
gridDim.x = 2

**threadIdx.x = 2**
**blockIdx.x = 1**

a = (1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 1, 1, 3, 2, 5, 6)
b = (2, 4, 5, 8, 3, 5, 7, 4, 5, 6, 7, 8, 1, 1, 2, 7)

| Block 0 | Block 1 |
|---|---|
| cache[0] = 47 | cache[0] = 18 |
| cache[1] = 14 | cache[1] = 32 |
| cache[2] = 22 | cache[2] = 59 |
| cache[3] = 40 | cache[3] = |

\*      \*

+

**cache[2] = 59**

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

dot<<2,4>>(a, b, c)

blockDim.x = 4
gridDim.x = 2

threadIdx.x = 3
blockIdx.x = 1

| Block 0 | Block 1 |
|---------|---------|
| cache[0] = 47 | cache[0] = 18 |
| cache[1] = 14 | cache[1] = 32 |
| cache[2] = 22 | cache[2] = 59 |
| cache[3] = 40 | cache[3] = 74 |

a = (1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 1, 1, 3, 2, 5, 6)
b = (2, 4, 5, 8, 3, 5, 7, 4, 5, 6, 7, 8, 1, 1, 2, 7)

\*        +        \*        → cache[3] = 74

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

▶ Make sure `cache` is full before continuing.

```
24  // synchronize threads in this block
25  __syncthreads();
```

▶ Execute a pairwise sum of `cache` for each block.

```
26      // threadsPerBlock must be a power of 2
27      int i = blockDim.x/2;
28      while (i != 0) {
29        if (cacheIndex < i)
30          cache[cacheIndex] += cache[cacheIndex + i
              ];
31        __syncthreads();
32        i /= 2;
33      }
```
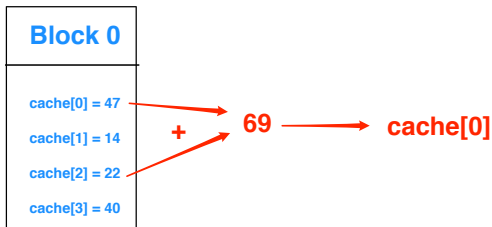
▶ Record the result in `partial_c`.

```
34      if (cacheIndex == 0)
35        partial_c[blockIdx.x] = cache[0];
36  }
```

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory



**dot<<2,4>>(a, b, c)**

**blockDim.x = 4**
**gridDim.x = 2**

**cacheIndex = threadIdx.x = 0**
**blockIdx.x = 0**
**i = 2**

**Block 0**

**cache[0] = 47**

**cache[1] = 14**    **+**    **69** ⟶ **cache[0]**

**cache[2] = 22**

**cache[3] = 40**

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

dot<<2,4>>(a, b, c)

blockDim.x = 4
gridDim.x = 2

cacheIndex = threadIdx.x = 1
blockIdx.x = 0
i = 2

**Block 0**

cache[0] = 69

cache[1] = 14

cache[2] = 22

cache[3] = 40

**+** **54** **cache[1]**

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

dot<<2,4>>(a, b, c)

blockDim.x = 4
gridDim.x = 2

**cacheIndex = threadIdx.x = 1**
**blockIdx.x = 0**
**i = 2**

### Block 0

cache[0] = 69

cache[1] = 54

cache[2] = 22

cache[3] = 40

# __synchthreads();

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

**dot<<2,4>>(a, b, c)**

**blockDim.x = 4**
**gridDim.x = 2**

**cacheIndex = threadIdx.x = 0**
**blockIdx.x = 0**
**i = 1**

**Block 0**

cache[0] = 69
cache[1] = 54      **+**      **123**  ⟶  **cache[0]**
cache[2] = 22
cache[3] = 40

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

**dot<<2,4>>(a, b, c)**

**blockDim.x = 4**
**gridDim.x = 2**

**Block 0**

cache[0] = 123

cache[1] = 54

cache[2] = 22

cache[3] = 40

**cacheIndex = threadIdx.x = 0**
**blockIdx.x = 0**
**i = 1**

# __synchthreads();

# dot<<<2, 4>>>(a, b, c) with $N = 16$

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

dot<<2,4>>(a, b, c)

blockDim.x = 4
gridDim.x = 2

**Block 0**

cache[0] = 123

cache[1] = 54

cache[2] = 22

cache[3] = 40

**cacheIndex = threadIdx.x = 0
blockIdx.x = 0
i = 0**

**i = 0, so end the pairwise sum.**

**The result for block 0 is cache[0] = 123.**

# Sum up partial_c inside int main()

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

```
37    dot<<<blocksPerGrid , threadsPerBlock >>>( dev_a ,
         dev_b , dev_partial_c );
38
39    // copy partial_c to the CPU
40    cudaMemcpy( partial_c , dev_partial_c ,
         blocksPerGrid * sizeof( float ),
         cudaMemcpyDeviceToHost );
41
42    // finish up on the CPU side
43    c = 0;
44    for ( int i=0; i<blocksPerGrid; i++) {
45      c += partial_c [ i ];
46    }
```

# Outline

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

Timing kernels on the GPU

Memory

# Resources

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

- ▶ Guides:

    1. J. Sanders and E. Kandrot. CUDA by Example. Addison-Wesley, 2010.
    2. D. Kirk, W.H. Wen-mei, and W. Hwu. *Programming massively parallel processors: a hands-on approach.* Morgan Kaufmann, 2010.
    3. Michael Romero and Rodrigo Urra. CUDA Programming. Rochester Institute of Technology. http://cuda.ce.rit.edu/cudaoverview/cudaoverview.html.

- ▶ Code:
    - ▶ time.cu
    - ▶ pairwise_sum_timed.cu
    - ▶ dot_product.cu

# That's all for today.

CUDA C:
performance
measurement and
memory

Will Landau

Timing kernels on
the GPU

Memory

▶ Series materials are available at
  http://will-landau.com/gpu.