Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
CUDA architecture

Beginning CUDA
C
Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD
paradigm

# Introduction to programming in CUDA C

Will Landau

Iowa State University

September 30, 2013

# Outline

A review: GPU parallelism and CUDA architecture

Beginning CUDA C
  Hello world
  Skeleton program
  Simple program
  Vector addition
  Pairwise summation
  Respecting the SIMD paradigm

Introduction to programming in CUDA C

Will Landau

A review: GPU parallelism and CUDA architecture

Beginning CUDA C
Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD paradigm

# Outline

## A review: GPU parallelism and CUDA architecture

## Beginning CUDA C
Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD paradigm

Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
CUDA architecture

Beginning CUDA
C
Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD
paradigm

# The single instruction, multiple data (SIMD) paradigm

▶ SIMD: apply the same command to multiple places in a dataset.

```
1  for ( i = 0;  i < 1e6;  ++i )
2    a[i] = b[i] + c[i];
```
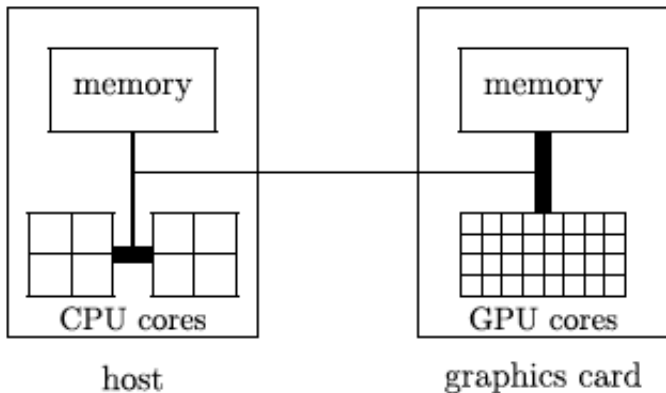
▶ On CPUs, the iterations of the loop run sequentially.

▶ With GPUs, we can easily run all 1,000,000 iterations simultaneously.

```
1  i = threadIdx.x;
2  a[i] = b[i] + c[i];
```

▶ We can similarly *parallelize* a lot more than just loops.

# CPU / GPU cooperation

- ▶ The CPU ("host") is in charge.
- ▶ The CPU sends computationally intensive instruction sets to the GPU ("device") just like a human uses a pocket calculator.
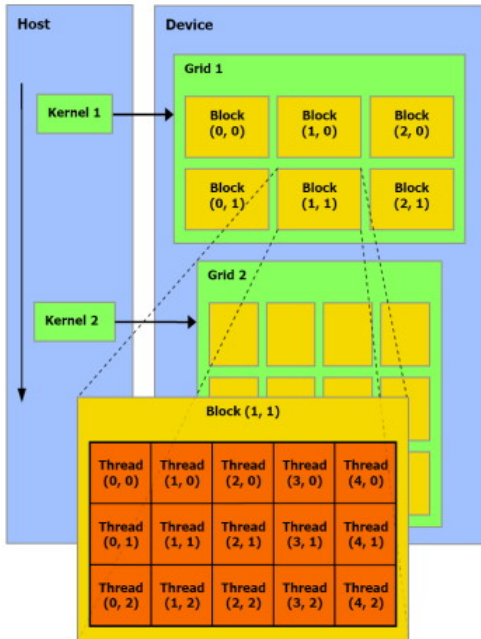
# How GPU parallelism works

1. The CPU sends a command called a **kernel** to a GPU.
2. The GPU executes several duplicate realizations of this command, called **threads**.
   - These threads are grouped into bunches called **blocks**.
   - The sum total of all threads in a kernel is called a **grid**.

- Toy example:
  - CPU says: "Hey, GPU. Sum pairs of adjacent numbers. Use the array, (1, 2, 3, 4, 5, 6, 7, 8)."
  - GPU thinks: "Sum pairs of adjacent numbers" is a kernel.
  - The GPU spawns 2 blocks, each with 2 threads:

| Block | 0 | | 1 | |
|---|---|---|---|---|
| Thread | 0 | 1 | 0 | 1 |
| Action | $1 + 2$ | $3 + 4$ | $5 + 6$ | $7 + 8$ |

- I could have also used 1 block with 4 threads and given the threads different pairs of numbers.

Introduction to programming in CUDA C

Will Landau

A review: GPU parallelism and CUDA architecture

Beginning CUDA C

Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD paradigm

# CUDA: making a gaming toy do science

Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
CUDA architecture

Beginning CUDA
C
Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD
paradigm

- ▶ **CUDA**: Compute Unified Device Architecture.
- ▶ Before CUDA, programmers could only do GPU programming in graphics languages, which are appropriate for video games but clumsy for science.
- ▶ CUDA devices support CUDA C, an extension of C for programs that use GPUs.
- ▶ CUDA-enabled servers at Iowa State:
    - ▶ impact1.stat.iastate.edu
    - ▶ impact2.stat.iastate.edu
    - ▶ impact3.stat.iastate.edu
    - ▶ impact4.stat.iastate.edu (in the works...)

# Outline

A review: GPU parallelism and CUDA architecture

## Beginning CUDA C

Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD paradigm

Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
CUDA architecture

Beginning CUDA
C

# Hello world

▶ A beginner C program:

```
1  #include <stdio.h>
2
3  int main(){
4    printf("Hello, World!\n");
5    return 0;
6  }
```

▶ A beginner CUDA C program:

```
1  #include <stdio.h>
2
3  __global__ void myKernel(){
4  }
5
6  int main(){
7    myKernel<<<1, 1>>>();
8    printf("Hello, World!\n");
9    return 0;
10 }
```

# Hello world

Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
CUDA architecture

Beginning CUDA
C

Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD
paradigm

```
 1  #include <stdio.h>
 2
 3  __global__ void myKernel(){
 4  }
 5
 6  int main(){
 7    myKernel<<<2, 4>>>();
 8    printf("Hello, World!\n");
 9    return 0;
10  }
```

- ▶ `__global__` says that the function is a kernel, which
  - ▶ will be executed on the GPU by one or more simultaneous threads when called.
  - ▶ must return void
- ▶ `<<<2, 4>>>` specifies
  - ▶ number of blocks (first number)
  - ▶ number of threads per block (second number).

# Prefixes in CUDA C

- ► __host__

    - ► Runs once per call on the CPU.
    - ► Only callable from the CPU (i.e., from another host function).
    - ► All functions without explicit prefixes are host functions.

- ► __global__

    - ► Used to specify a kernel.
    - ► Runs multiple times per call on the GPU (that's what <<<#, #>>> is for).
    - ► Only callable from the CPU (i.e., from a host function).

- ► __device__

    - ► Runs once per call on the GPU.
    - ► Only callable from the GPU (i.e., from either a kernel or another device function).

# Prefix example: 2 blocks and 5 threads per block

```
1  #include <stdio.h>
2
3  __device__ int dev1(){
4  }
5
6  __device__ int dev2(){
7  }
8
9  __global__ void pleaseRunThis10Times(){
10    dev1();
11    dev2();
12 }
13
14 int main(){
15    pleaseRunThis10Times<<<2, 5>>>();
16    printf("Hello, World!\n");
17    return 0;
18 }
```

# skeleton.cu: outlining a CUDA C workflow

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <cuda.h>
4  #include <cuda_runtime.h>
5
6  __global__ void some_kernel(...) {...}
7
8  int main (void){
9     // Declare all variables.
10    ...
11    // Allocate host memory.
12    ...
13    // Dynamically allocate device memory for GPU results.
14    ...
15    // Write to host memory.
16    ...
17    // Copy host memory to device memory.
18    ...
19
20    // Execute kernel on the device.
21    some_kernel<<< num_blocks, num_theads_per_block >>>(...);
22
23    // Write GPU results in device memory back to host memory.
24    ...
25    // Free dynamically-allocated host memory
26    ...
27    // Free dynamically-allocated device memory
28    ...
29  }
```

# simple.cu: a program that actually does something

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <cuda.h>
4  #include <cuda_runtime.h>
5
6  __global__ void colonel(int *a_d){
7    *a_d = 2;
8  }
9  int main(){
10   int a = 0, *a_d;
11
12   cudaMalloc((void**) &a_d, sizeof(int));
13   cudaMemcpy(a_d, &a, sizeof(int), cudaMemcpyHostToDevice);
14
15   colonel<<<1,1>>>(a_d);
16
17   cudaMemcpy(&a, a_d, sizeof(int), cudaMemcpyDeviceToHost);
18
19   printf("a = %d\n", a);
20   cudaFree(a_d);
21
22  }
```

# Compiling and running `simple.cu`

Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
CUDA architecture

Beginning CUDA
C
Hello world
Skeleton program
**Simple program**
Vector addition
Pairwise summation
Respecting the SIMD
paradigm

```
1 > nvcc simple.cu -o simple
2 > ./simple
3 a = 2
```

- ▶ Notes:
    - ▶ `nvcc` is the NVIDIA CUDA C compiler,
    - ▶ CUDA C source files usually have the `*.cu` extension, though they sometimes have have `*.c` and `*.cpp` extensions.
    - ▶ This code is available at `http://will-landau.com/gpu/Code/CUDA_C/simple/simple.cu`.
    - ▶ Most of the example code I present will be linked from pages at `will-landau.com/gpu/talks`.

# Builtin CUDA C variables

- ► `maxThreadsPerBlock`: exactly that: 1024 on impact1.
- ► For a kernel call with $B$ blocks and $T$ threads per block,
  - ► `blockIdx.x`
    - ► ID of the current block (in the $x$ direction).
    - ► Integer from 0 to $B - 1$ inclusive.
  - ► `threadIdx.x`
    - ► within the current block, ID of the current thread (in the $x$ direction).
    - ► Integer from 0 to $T - 1$ inclusive.
  - ► `gridDim.x`: number of blocks in the current grid (in the $x$ direction).
  - ► `blockDim.x`: number of threads per block (in the $x$ direction).
- ► With some modifications that I will describe in later lectures, you can use the $y$ and $z$ directions with variables like `threadIdx.y`, `threadIdx.z` etc.

# Vector addition: `vectorsums.cu`

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <cuda.h>
4   #include <cuda_runtime.h>
5
6   #define N 10
7
8   __global__ void add(int *a, int *b, int *c){
9       int bid = blockIdx.x;
10      if(bid < N)
11          c[bid] = a[bid] + b[bid];
12  }
13
14  int main(void) {
15      int i, a[N], b[N], c[N];
16      int *dev_a, *dev_b, *dev_c;
17
18      cudaMalloc((void**) &dev_a, N*sizeof(int));
19      cudaMalloc((void**) &dev_b, N*sizeof(int));
20      cudaMalloc((void**) &dev_c, N*sizeof(int));
21
22      for(i=0; i<N; i++){
23          a[i] = -i;
24          b[i] = i*i;
25      }
26
27      cudaMemcpy(dev_a, a, N*sizeof(int), cudaMemcpyHostToDevice);
28      cudaMemcpy(dev_b, b, N*sizeof(int), cudaMemcpyHostToDevice);
```

# Vector addition: `vectorsums.cu`

```
29
30    add<<<N,1>>>(dev_a, dev_b, dev_c);
31
32    cudaMemcpy(c, dev_c, N*sizeof(int), cudaMemcpyDeviceToHost);
33
34    printf("\na + b = c\n");
35    for(i = 0; i<N; i++){
36      printf("%5d + %5d = %5d\n", a[i], b[i], c[i]);
37    }
38
39    cudaFree(dev_a);
40    cudaFree(dev_b);
41    cudaFree(dev_c);
42  }
```

# Compiling and running `vectorsums.cu`

Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
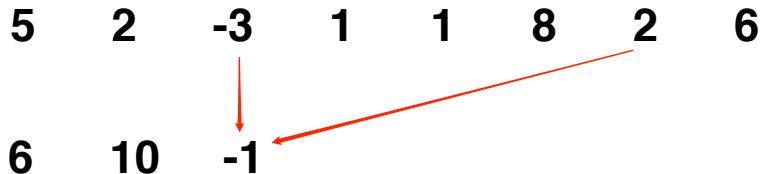CUDA architecture

Beginning CUDA
C

Hello world
Skeleton program
Simple program
**Vector addition**
Pairwise summation
Respecting the SIMD
paradigm

```
 1  > nvcc vectorsums.cu −o vectorsums
 2  > ./vectorsums
 3  a + b = c
 4       0 +     0 =     0
 5      −1 +     1 =     0
 6      −2 +     4 =     2
 7      −3 +     9 =     6
 8      −4 +    16 =    12
 9      −5 +    25 =    20
10      −6 +    36 =    30
11      −7 +    49 =    42
12      −8 +    64 =    56
13      −9 +    81 =    72
```

# Synchronizing threads within blocks: the pairwise sum revisited

- Example: pairwise sum of the vector (5, 2, -3, 1, 1, 8, 2, 6)

**5    2    -3    1    1    8    2    6**

**6**

# Thread 0

# Synchronizing threads within blocks: the pairwise sum revisited

**5        2        -3        1        1        8        2        6**

**6        10**

**Thread  1**

# Synchronizing threads within blocks: the pairwise sum revisited

**5      2      -3      1      1      8      2      6**

**6      10      -1**

## Thread 2

# Synchronizing threads within blocks: the pairwise sum revisited

**5    2    -3    1    1    8    2    6**

**6    10    -1    7**

## Thread 3

# Synchronizing threads within blocks: the pairwise sum revisited

**5      2      -3      1      1      8      2      6**

**6      10      -1      7**

## Synchronize threads

# Synchronizing threads within blocks: the pairwise sum revisited

**5    2    -3    1    1    8    2    6**

**6    10    -1    7**

**5**

**Thread 0**

# Synchronizing threads within blocks: the pairwise sum revisited

| 5 | 2 | -3 | 1 | 1 | 8 | 2 | 6 |

| 6 | 10 | -1 | 7 |

| 5 | 17 |

**Thread 1**

# Synchronizing threads within blocks: the pairwise sum revisited

**5    2    -3    1    1    8    2    6**

**6    10    -1    7**

**5    17**

## Synchronize Threads

# Synchronizing threads within blocks: the pairwise sum revisited

**5      2      -3      1      1      8      2      6**

**6      10      -1      7**

**5      17**

**Thread 0**

**22**

# Pairwise sum in pseudocode

▶ Let $n = 2^m$ be the length of the vector.

▶ Denote the vector by $(x_{(0,\ 0)}, \ldots, x_{(0,\ n-1)})$

▶ Spawn 1 grid with a single block of $n/2$ threads.

▶ Do:

1. Set offset $= n/2$.
2. For parallel threads $j = 0, \ldots,$ offset $- 1$, compute:

$$x_{(i,\ j)} = x_{(i-1,\ j)} + x_{(i-1,\ j+\text{offset})}$$

3. Synchronize threads.
4. Integer divide offset by 2.
5. Return to step 2 if offset $> 0$.

# pairwise_sum.cu

```c
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <math.h>
4  #include <cuda.h>
5  #include <cuda_runtime.h>
6
7  /*
8   * This program computes the sum of the elements of
9   * vector v using the pairwise (cascading) sum algorithm.
10  */
11
12 #define N 8 // 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 }
```

# pairwise_sum.cu

```
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  int main (void){
45    float *v_h, *v_d; // host and device copies of our vector,
          respectively
46
47    // dynamically allocate memory on the host for v_h
48    v_h = (float*) malloc(N * sizeof(*v_h));
49
50    // dynamically allocate memory on the device for v_d
51    cudaMalloc ((float**) &v_d, N *sizeof(*v_d));
52
53    // Fill v_h with N random floating point numbers.
54    vfill(v_h, N);
55
56    // Print v_h to the console
57    vprint(v_h, N);
```

# pairwise_sum.cu

```
58      // Write the contents of v_h to v_d
59      cudaMemcpy( v_d, v_h, N * sizeof(float), cudaMemcpyHostToDevice );
60
61      // Compute the pairwise sum of the elements of v_d and store the
               result in v_d[0].
62      psum<<< 1, N/2 >>>(v_d);
63
64      // Write the pairwise sum, v_d[0], to v_h[0].
65      cudaMemcpy(v_h, v_d, sizeof(float), cudaMemcpyDeviceToHost );
66
67      // Print the pairwise sum.
68      printf("Pairwise sum = %7.3f\n", v_h[0]);
69
70      // Free dynamically-allocated host memory
71      free(v_h);
72
73      // Free dynamically-allocated device memory
74      cudaFree(v_d);
75  }
```

# Compiling and running `pairwise_sum.cu`

Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
CUDA architecture

Beginning CUDA
C

Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD
paradigm

```
 1 > nvcc pairwise_sum.cu -o pairwise_sum
 2 > ./pairwise_sum
 3 v =
 4     0.840
 5     0.394
 6     0.783
 7     0.798
 8     0.912
 9     0.198
10     0.335
11     0.768
12
13 Pairwise sum =     5.029
```

# Best practices: respect the SIMD paradigm

- ► SIMD: "Single Instruction, Multiple Data"
- ► Under this paradigm, the thread in a kernel call write to different memory spaces.
- ► When threads write to the same memory (SISD), problems can arise.

# sisd.cu: violating the SIMD paradigm

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <cuda.h>
4  #include <cuda_runtime.h>
5
6  __global__ void colonel(int *a_d){
7    *a_d = blockDim.x * blockIdx.x + threadIdx.x;
8  }
9
10 int main(){
11
12   int a = 0, *a_d;
13
14   cudaMalloc((void**) &a_d, sizeof(int));
15   cudaMemcpy(a_d, &a, sizeof(int), cudaMemcpyHostToDevice);
16
17   colonel<<<4,5>>>(a_d);
18
19   cudaMemcpy(&a, a_d, sizeof(int), cudaMemcpyDeviceToHost);
20
21   printf("a = %d\n", a);
22   cudaFree(a_d);
23
24 }
```

▶ What is the output?

# sisd.cu: violating the SIMD paradigm

```
1 > nvcc sisd.cu −o sisd
2 > ./sisd
3 a = 14
```

▶ The output is unpredictable because the threads modify
  the same variable in an unpredictable order.

# Outline

A review: GPU parallelism and CUDA architecture

Beginning CUDA C
  Hello world
  Skeleton program
  Simple program
  Vector addition
  Pairwise summation
  Respecting the SIMD paradigm

Introduction to
programming in
CUDA C

Will Landau

A review: GPU
parallelism and
CUDA architecture

Beginning CUDA
C

Hello world
Skeleton program
Simple program
Vector addition
Pairwise summation
Respecting the SIMD
paradigm

# Resources

▶ Texts:

　　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.

▶ Code:

　　▶ skeleton.cu
　　▶ simple.cu
　　▶ vectorsums.cu
　　▶ pairwise_sum.cu

# That's all for today.

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