

# CME213/ME339

## Lecture 11

Erich Elsen

Department of Mechanical Engineering  
Institute for Computational and Mathematical Engineering  
Stanford University

Spring 2013



# What is Thrust?

---

- Template Header library similar to the C++ STL but targeted to parallel execution
- Performance portable way to express parallel algorithms
- Take advantage of expertly tuned implementations targeted for each GPU architecture
- Multiple Backends - run code on GPUs with CUDA or on multiple cores with OMP or TBB
- Highly productive way to program with CUDA



# First Example

---

```
1  #include <thrust/host_vector.h>
2  #include <thrust/device_vector.h>
3  #include <thrust/sequence.h>
4  #include <thrust/sort.h>
5
6  int main(void) {
7      thrust::device_vector<int> d_data(1000);
8
9      //generate a sequence counting down from 1000
10     thrust::sequence(d_data.begin(), d_data.end(), 1000, -1);
11
12     d_data[30] = -6;
13
14     thrust::sort(d_data.begin(), d_data.end());
15
16     thrust::host_vector<int> h_data = d_data;
17 }
```



# Advantages

---

- Automatic Memory Management with containers (no `cudaMalloc` / `cudaFree`)
- Cleaner code for memory movement (no verbose `cudaMemcpy`)
- Read / Write to specific array locations very useful for debugging
- Invoking an algorithm with a device vector will run it on the device
- Invoking an algorithm with a host vector will run it on the host
- Meaning of "device" and "host" can be changed when compiling
- `-DTHRUST_DEVICE_SYSTEM=THRUST_HOST_SYSTEM_OMP`



# STL and Functional Programming

---

Say we want to double every number in a list. The procedural way is:

```
1   for (int i = 0; i < N; ++i)
2       myList[i] *= 2;
```

We declare not only what we want done (doubling) but also HOW it is done (with a for loop).



# Procedural vs. Functional Programming

---

Say we want to double every number in a list. The functional way is:

```
1 //declare a functor
2 struct doubler {
3     void operator()(int& x) {
4         x *= 2;
5     }
6 };
7 std::for_each(mylist.begin(), mylist.end(), doubler());
```

We declare only what we want done (doubling). How this happens is not specified.



# Functors

Can have internal state...

```
1  template<typename T>
2  struct xer {
3      T multiple;
4      //constructor to initialize multiple
5      xer(const T& m) : multiple(m) {}
6
7      void operator()(T& x) {
8          x *= multiple;
9      }
10 };
11
12 //declare our functors
13 xer<float> mf(5.4f);
14 xer<int>    mi(3);
15
16 std::for_each(floatlist.begin(), floatlist.end(), mf);
17 std::for_each(intlist.begin(), intlist.end(), mi);
```



# Iterators

---

- A powerful concept that allows for applying algorithms such as `for_each` to a variety of containers
  - vectors, lists, sets, maps, etc...
- On the GPU the only container that makes much sense is the vector
- Vector iterators are nice in that they are essentially just pointers
- `.begin()` points to the first element
- `.end()` points to one past the last element
- Algorithms also work with raw pointers

```
1  int a[10];  
2  std::for_each(a, a + 10, doubler());
```





# Iterators and Thrust

---

To use a device vector in a normal kernel, you need to do:

```
1 myKernel<<<blocks, threads>>>(  
2   thrust::raw_pointer_cast(&d_vector[0]));
```

To use a raw device pointer with a thrust algorithm, you need to let thrust know this pointer exists on the device:

```
1 thrust::device_ptr<float> my_thrust_ptr(my_device_ptr);  
2 thrust::sort(my_thrust_ptr, my_thrust_ptr + N);
```



# Now with Thrust

```
1  template<typename T>
2  struct xer {
3      T multiple;
4      //constructor to initialize multiple
5      __host__ __device__
6      xer(const T& m) : multiple(m) {}
7
8      __host__ __device__
9      void operator()(T& x) {
10         x *= multiple;
11     }
12 };
13
14 //declare our functors
15 xer<int>    mi(3);
16
17 thrust::host_vector<int> h_intlist(10);
18 thrust::device_vector<int> d_intlist(10);
19
20 //runs on the CPU
21 thrust::for_each(h_intlist.begin(), h_intlist.end(), mi);
22 //runs on the GPU
23 thrust::for_each(d_intlist.begin(), d_intlist.end(), mi);
```



# Radix Sort with Thrust

---

Sort this sequence: 0 1 1 0 1 0 1

- Idea: Each 0 needs to know how many 0s are before it
- Equivalent to knowing the current position and how many 1s before us
- Each 1 needs to know many 1s are before it and also how many 0s
- We need Scan, Transform and Scatter

Scan it: 0 0 1 2 2 3 3 4

Now each position "knows" how many ones come before it and also how many total ones there are (the last entry)



# Radix Sorting with Thrust

---

## Transform

- If we are a 0 at position 5, then we need to know how many 0s are before us, but we know how many 1s are before us.
- That is equal to our position (5) minus the number of 1s before us. This is our output position.
- If we are a 1, our final position is total number of 0s plus the number of 1s before us.

```
1  if (digit == 0)
2      outputPos = ourPosition - ScanVal;
3  else //digit == 1
4      outputPos = numZeros + ScanVal;
```

Transform: 0 3 4 1 5 2 6



# Radix Sort With Thrust

---

## Scatter

- $\text{output}[\text{outputLoc}[i]] = \text{input}[i]$
- $\text{outputLoc}$  = transformed values
- $\text{input}$  = original input sequence

Input:            0 1 1 0 1 0 1

OutputLocs: 0 3 4 1 5 2 6

Output:           0 0 0 1 1 1 1



# Radix Sort With Thrust

---

- We would perform this process for total number of bits that we need to sort
- Thrust has a better way of sorting that we've seen earlier
- The built-in radix sort is *extremely* fast
- This example is for pedagogical purposes, don't actually write your own sort

```
1 thrust::sort(input.begin(), input.end());
```

Done.



# Kernel Fusion

Image we want to compute  $\sum_{i=0}^N x_i^2$

Approach 1:

```
1  struct square {
2      __host__ __device__
3      float operator()(float x) {
4          return x * x;
5      }
6  };
7
8  thrust::device_vector<float> input(100);
9  thrust::device_vector<float> squares(100);
10
11  //transform(InputStart, InputEnd, OutputStart, Function);
12  thrust::transform(input.begin(), input.end(),
13                  squares.begin(), square());
14  float sum = thrust::reduce(squares.begin(), squares.end());
```

- Slow - we write the squares back to memory and then read them in again
- Wastes space - we don't need to store the temporary values



# Transform Iterators and Kernel Fusion

---

Transform Iterators Apply a transform in-place allowing the new value to be used locally.

```
1  struct square {
2      __host__ __device__
3      float operator()(float x) {
4          return x * x;
5      }
6  };
7
8  thrust::device_vector<float> input(100);
9
10 float sum = thrust::reduce(
11     thrust::make_transform_iterator(
12         input.begin(), square()),
13     thrust::make_transform_iterator(
14         input.end(), square()));
```





# Constant and Counting Iterators

---

We often need to represent a constant number or an increasing sequence. Explicitly representing them would waste both storage space and memory resources moving them around.

```
1 //sum = 100
2 int sum = thrust::reduce(thrust::make_constant_iterator(1),
3                          thrust::make_constant_iterator(1) + 100);
4
5 //sum = 5050
6 sum = thrust::reduce(thrust::make_counting_iterator(1),
7                      thrust::make_counting_iterator(101));
```



# Histogram Example

---

Given the sequence:

[2 1 0 0 2 2 1 1 1 1 4]

the dense histogram would be:

[2 5 3 0 1]

the sparse histogram would be:

[(0,2), (1,5), (2,3), (4,1)]

- How might we implement these operations with thrust algorithms?
- The first step - sort!
- Even when a serial algorithm might not involve sorting it is often a useful primitive when using thrust
- Let's examine the sparse case - we can do it with one thrust call!



# Segmented Reduction

```
1 thrust::pair<OutputIterator1,OutputIterator2>  
2 thrust::reduce_by_key(InputIterator1 keys_first,  
3                       InputIterator1 keys_last,  
4                       InputIterator2 values_first,  
5                       OutputIterator1 keys_output,  
6                       OutputIterator2 values_output);
```

We perform multiple reductions using adjacent identical keys to determine which reductions to perform

Keys : 1 3 3 3 2 2 1  
Vals : 9 8 7 6 5 4 3

Out Keys : 1 3 2 1  
Out Vals : 9 21 9 3

So how do we use this to do the histogram?



# Sparse Histogram

After sorting...

Sequence:           0 0 1 1 1 1 1 2 2 2 4 (keys)

Const. Iterator: 1 1 1 1 1 1 1 1 1 1 1 (vals)

Out Keys: 0 1 2 4 E X X X X X X

Out Vals: 2 5 3 1 E X X X X X X

```
1 thrust::device_vector<int> data(11); // = [2 1 0 0 2 2 1 1 1 1 4]
2 thrust::device_vector<int> histogram_values(11);
3 thrust::device_vector<int> histogram_counts(11);
4
5 thrust::sort(data.begin(), data.end());
6
7 typedef thrust::device_vector<int>::iterator devIt;
8 thrust::pair<devIt, devIt> endIterators =
9     thrust::reduce_by_key(data.begin(), data.end(),
10                          thrust::make_constant_iterator(1),
11                          histogram_values.begin(),
12                          histogram_counts.begin());
13
14 int num_values = endIterators.first - histogram_values.begin();
```



# Sparse Histogram

---

- What if we don't want to over-allocate space for `histogram_values` and `histogram_counts`?
- Must count how many distinct values in keys
- Can do this with a clever use of `inner_product`
- Inner product:  $(a_0 \otimes b_0) \oplus (a_1 \otimes b_1) \oplus \dots$
- Dot product is  $\otimes = \times$  and  $\oplus = +$

Sequence A:            0 0 1 1 1 1 1 2 2 2 4 (keys)

Sequence B:            0 0 1 1 1 1 1 2 2 2 4 (keys)

What should our operations be?



# Sparse Histogram

```
1 int num_bins = thrust::inner_product(data.begin(), data.end() - 1,  
2                                     data.begin() + 1,  
3                                     (int)1, thrust::plus<int>(),  
4                                     thrust::not_equal_to<int>());
```

Sequence A:            0 0 1 1 1 1 1 2 2 2 4 (keys)

Sequence B:            0 0 1 1 1 1 1 2 2 2 4 (keys)

A != B            :            0 1 0 0 0 0 1 0 0 1

num\_bins        :            1 + 0+1+0+0+0+0+1+0+0+1 = 4



## Second Example - Point Binning

---

- First generate a random collection of 2d points
- Then bin these points into a 2d grid of cells
- Finally extract which cells have only one point in them

|       |   |     |     |
|-------|---|-----|-----|
| •     |   | • • |     |
| • • • | • |     |     |
|       |   | •   | • • |
| • • • | • |     |     |



# Random Number Generation

## With Thrust

---

- Can do RNG on both host and device
- We will focus on host generation
- We need both a generator AND a distribution
- A Generator produces random bits
  - `thrust::default_random_engine`
- The Distribution turns these bits into something useful - uniformly distributed floats between  $[-3, 10]$  for example
  - `thrust::uniform_real_distribution`
  - `thrust::uniform_int_distribution`





# Random Point Generation

---

```
1  // return a random vec2 in  $[0,1]^2$ 
2  vec2 make_random_float2(void)
3  {
4      //The static is important!
5      static thrust::default_random_engine rng;
6      static thrust::uniform_real_distribution<float> u01(0.0f, 1.0f);
7      float x = u01(rng);
8      float y = u01(rng);
9      return vec2(x,y);
10 }
11
12 thrust::host_vector<float2> h_points(N);
13 thrust::generate(h_points.begin(), h_points.end(), make_random_float2);
14
15 thrust::device_vector<float2> points = h_points;
```



# Grid Structure

---

```
1 // allocate storage for a 2D grid
2 // of dimensions w x h
3 unsigned int w = 200, h = 100;
4
5 // the grid data structure keeps a range per grid bucket:
6 // each bucket_begin[i] indexes the first element of
7 // bucket i's list of points
8 thrust::device_vector<unsigned int> bucket_begin(w*h);
9
10 // each bucket_end[i] indexes one past the last element of
11 // bucket i's list of points
12 thrust::device_vector<unsigned int> bucket_end(w*h);
13
14 // allocate storage for each point's bucket index
15 thrust::device_vector<unsigned int> bucket_indices(N);
```



# Point to Bucket Functor

```
1 // hash a point in the unit square to the index of
2 // the grid bucket that contains it
3 struct point_to_bucket_index :
4     thrust::unary_function<float2,unsigned int>
5 {
6     float width; // buckets in the x dimension (grid spacing = 1/width)
7     float height; // buckets in the y dimension (grid spacing = 1/height)
8
9     __host__ __device__
10    point_to_bucket_index(unsigned int width, unsigned int height)
11        : width(width), height(height) {}
12
13    __host__ __device__
14    unsigned int operator()(const float2& v) const
15    {
16        // find the raster indices of p's bucket
17        unsigned int x = static_cast<unsigned int>(v.x * width);
18        unsigned int y = static_cast<unsigned int>(v.y * height);
19
20        // return the bucket's linear index
21        return y * width + x;
22    }
23
24 };
```



# Points → Cells

```
1  // transform the points to their bucket indices
2  thrust::transform(points.begin(),
3                    points.end(),
4                    bucket_indices.begin(),
5                    point_to_bucket_index(w,h));
6
7  // sort the points by their bucket index
8  thrust::sort_by_key(bucket_indices.begin(),
9                     bucket_indices.end(),
10                     points.begin());
```

Transform (assume w=10 h=10 here):

|         |         |         |         |         |            |
|---------|---------|---------|---------|---------|------------|
| Points: | (.7,.8) | (.4,.1) | (.6,.4) | (.2,.3) | (.62, .43) |
| Index : | 87      | 14      | 64      | 32      | 64         |

Sort:

|         |         |         |         |            |          |
|---------|---------|---------|---------|------------|----------|
| Index : | 14      | 32      | 64      | 64         | 87       |
| Points: | (.4,.1) | (.2,.3) | (.6,.4) | (.62, .43) | (.7, .8) |



# Determine Contents of Each Cell

---

- For each cell  $[0, w \times h)$  we need to figure out its bounds
- Hmm...makes me think of a `counting_iterator`
- The algorithms we need are `lower_bound` and `upper_bound`
- `lower_bound` takes a sequence and a list of search values
- For each search value, it finds the first place in the sequence it could be inserted without changing the ordering

Seq: 0 0 1 2 4 4 5 6 6 6 7

3 6

`lower_bound`: 3 6

Output: 4 7



# Determine Contents of Each Cell

---

- `upper_bound` is similar except it finds the *last* place in the sequence the value can be inserted without changing the ordering

Seq: 0 0 1 2 4 4 5 6 6 6 7  
          3                  6

`upper_bound`: 3 6

Output: 4 10

Now we can determine the number of points in each cell by subtracting the output of `lower_bound` from `upper_bound`



# Code for Cell Determination

---

- The sequence we're searching *in* comes first
- The values we're searching *for* come next
- The output goes last

```
1  // find the beginning of each bucket's list of points
2  thrust::counting_iterator<unsigned int> search_begin(0);
3
4  thrust::lower_bound(bucket_indices.begin(),
5                      bucket_indices.end(),
6                      search_begin,
7                      search_begin + w*h,
8                      bucket_begin.begin());
9
10 // find the end of each bucket's list of points
11 thrust::upper_bound(bucket_indices.begin(),
12                    bucket_indices.end(),
13                    search_begin,
14                    search_begin + w*h,
15                    bucket_end.begin());
```



# Extracting Lonely Points

---

The function we need for this is called `remove_copy_if`

Cell:        0 1 2 3 4 5 6 7 8 9

# Points: 0 1 3 2 1 4 0 4 1 0

Out: 1 4 8

We need a predicate to determine which cells to remove based on the value of points

```
1  struct is_equal_to_one : thrust::unary_function<int, int>
2  {
3      __host__ __device__
4      int operator()(const int& v)
5      {
6          return v == 1;
7      }
8  };
```





# Extract Cells

---

```
1 thrust::device_vector<int> bucket_sizes(N);
2 thrust::transform(bucket_end.begin(), bucket_end.end(),
3                   bucket_begin.begin(), bucket_sizes.begin(),
4                   thrust::minus<int>());
5
6 int num_lonely_cells = thrust::count_if(bucket_sizes.begin(),
7                                         bucket_sizes.end(),
8                                         is_equal_to_one());
9
10 thrust::device_vector<int> lonely_cells(num_lonely_cells);
11 thrust::remove_copy_if(make_counting_iterator(0),
12                       make_counting_iterator(w*h),
13                       bucket_sizes.begin(),
14                       lonely_cells.begin(),
15                       is_equal_to_one() );
```



# Maximum Number of Points

---

To determine the cell with the most points, we could use `max_element` which returns an *iterator* to the largest element

```
1 thrust::device_vector<int>::iterator maxIt;  
2 maxIt = thrust::max_element(bucket_sizes.begin(), bucket_sizes.end());  
3  
4 int maxNum = *maxIt;  
5 int maxPos = maxIt - bucket_sizes.begin();
```

# Points: 0 1 3 2 1 4 0 4 1 0

maxNum: 4

maxPos: 5 (returns the first if there are duplicates)



# Additional Thrust Information

---

<https://github.com/thrust/thrust/wiki/Quick-Start-Guide>

<https://github.com/thrust/thrust/wiki/Documentation>

<http://developer.download.nvidia.com/CUDA/training/GTCthrust.mp4>



# Vigenère Cipher

---

Plain text: ILIKEMYTEACHER

KEY: NOTNOTNOTNOTNO

=====

Cipher text: VZBXSFLHXNQARF

- By using multiple shifts (or permutations) instead of just 1 as in a substitution cipher the frequency distribution of the cipher text is obscured
- If we knew the key length, then we could solve multiple substitution ciphers

Plain text: ILIKEMYTEACHER

KEY: NOTNOTNOTNOTNO

=====

Cipher text: VZBXSFLHXNQARF



# Determine Key Length

First define an Index of Coincidence (IOC) between two texts as:

$$\frac{26 * \sum_{i=0}^N A_i == B_i}{N}$$

```
MYTEACHERISAWESOME
ILOVETHRUSTCODING
=====
000000100000000000
```

$$\text{IOC} = 26 * 1 / 17 = 1.53$$

- Defined so that the IOC between two *randomly* chosen texts is 1
- In English the IOC between two different texts (Moby Dick and The Great Gatsby) is  $\sim 1.73$
- Which is because letters are not chosen randomly, but have a non-uniform frequency distribution



# Breaking the Vigenère

---

Shifts don't line up and the IOC is  $\sim 1$ , the texts appear random

VZBXSFLHXNQARF  
VZBXSFLHXNQARF

VZBXSFLHXNQARF  
VZBXSFLHXNQARF

Shifts line up and suddenly the IOC jumps to  $\sim 1.73$ , because now at each position the "alphabet" though jumbled, is the same

VZBXSFLHXNQARF  
VZBXSFLHXNQARF

