# CME213/ME339
# Lecture 4

Eric Darve
Erich Elsen
Austin Gibbons

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

Spring 2012

# Vector Types

- CUDA has built-in vector types
- char1, char2, char3, char4, uchar1, uchar2, uchar3, uchar4
- int1, int2, int3, int4, uint1, uint2, uint3, uint4
- float1, float2, float3, float4
- double1, double2
- also for short, ushort, long, ulong, longlong, ulonglong

| | |
|---|---|
| 1st element | .x |
| 2nd element | .y |
| 3rd element | .z |
| 4th element | .w |

# Vector Types

- In general, avoid the 3 component versions (they don't coalesce well)
- You're not crazy - the 1 component version isn't very useful
- The 2 and 4 componenet versions are very useful

```
1  //how to initialize a vector type
2  float2 foo = make_float2(3.f, 4.f);
3
4  float sum = foo.x + foo.y; //== 7
```

# Simple Performance Example

```
1   __global__
2   void testCalc1(int N, float *d_out)
3   {
4     const int tid = blockIdx.x * blockDim.x + threadIdx.x;
5     const float val = 10000.f;
6
7     if (tid < N)
8       d_out[tid] = powf(lgamma(val) * tgamma(val) + erfcinv(val),
9                         sinf(3.f));
10  }
```

# Timing

| # Threads | MegaThreads / sec |
|-----------|-------------------|
| 1024      | 24.8              |
| 2048      | 30.9              |
| 4096      | 54.3              |
| 8192      | 62.6              |
| 16384     | 68.9              |
| 32768     | 71.7              |
| 65536     | 73.0              |
| 131072    | 73.8              |
| 262144    | 74.2              |

- Peak performance is not reached until a certain threshold is crossed.
- Here the threshold is $\sim$ 32,000 threads which is a good general ballpark number.

# Experiment 1

What if half the threads do nothing?

```
1   __global__
2   void testCalc1(int N, float *d_out)
3   {
4     const int tid = blockIdx.x * blockDim.x + threadIdx.x;
5     const float val = 10000.f;
6
7     if (tid >= N)
8         return;
9
10    if (tid % 2)
11      d_out[tid] = powf(lgamma(val) * tgamma(val) + erfcinv(val),
12                        sinf(3.f));
13  }
```

# Experiment 1

| # Threads | MegaThreads / sec |
|-----------|-------------------|
| 1024      | 24.8              |
| 2048      | 29.9              |
| 4096      | 52.9              |
| 8192      | 62.2              |
| 16384     | 68.0              |
| 32768     | 70.9              |
| 65536     | 72.3              |
| 131072    | 73.1              |
| 262144    | 73.8              |

- No change?! But we did half the work!
- Is there some stride for which it does get faster?

# Experiment 1

| Stride | MegaThreads / sec |
|--------|-------------------|
| 10 | 82.2 |
| 20 | 82.2 |
| 31 | 82.2 |
| 32 | 82.2 |
| 33 | 84.2 |
| 34 | 86.1 |
| 40 | 97.4 |
| 50 | 113.4 |
| 64 | 129.9 |

- N = 1048576
- Any stride < 32 gives the same performance
- After that we see increasing performance with increased stride. Why?

# Warps

- This experiment suggests that the smallest unit of execution is 32 threads
- Each group of 32 threads is known as a warp
- tid 0-31 - 0th warp; tid 32-63 - 1st warp; etc.
- All threads in a warp execute the same instruction at the same time

```
1   if (tid % 2)
2       out[tid] = foo(tid);
3   else
4       out[tid] = bar(tid);
```

- When threads within a warp take different code paths, this is called warp divergence

# Warps

- Warp divergence is handled by *serialization*
- The divergent code path is executed as many times as there are different branches
- Divergence is determined at runtime

```
1  if (blockIdx.x == 3)
2    foo(...);
3  else
4    bar(...);
```

- is NOT divergent - all threads in each warp will take the same path through code

# Quiz

Which of these code snippets diverge? How many times will the divergent code be executed?

```
1   for (int i = 0; i < threadIdx.x; ++i)
2       foo(i);
```

```
1   if (threadIdx.x / 32 == 0)
2       foo(threadIdx.x);
3   else if (threadIdx.x / 32 == 1)
4       bar(threadIdx.x);
5   else
6       foobar(threadIdx.x);
```

```
1   if (threadIdx.x % 32)
2       foo(threadIdx.x);
3   else
4       bar(threadIdx.x);
```

# Answers

- First snippet diverges. The loop will execute 32 times for the first warp of a block, 64 times for the second warp and so on.
- Second snippet *doesn't* diverge. Each warp takes the same path through the conditional statement.
- Third snippet diverges. The divergent code is executed twice.

# Warps And Memory Transactions

- Threads do not individually request memory
- Memory transactions are always 128 bytes
- We want as many bytes as possible of each request to useful
- Memory transactions are issued per warp for sizes $\leq 4$ bytes per thread
- For 8 byte transactions each half warp issues a transaction
- And for 16 byte transactions each quarter warp issues a transaction

# Examples

Imagine a warp consists of 8 threads instead of 32.

| Thread Id | Memory Location | Size |
|:---------:|:---------------:|:----:|
| 0 | 0x80 | 4 |
| 1 | 0x84 | 4 |
| 2 | 0x88 | 4 |
| 3 | 0x8C | 4 |
| 4 | 0x90 | 4 |
| 5 | 0x94 | 4 |
| 6 | 0x98 | 4 |
| 7 | 0x9C | 4 |

For a warp of 32 threads, this access pattern would generate one memory request of $32 * 4 = 128$ bytes.

# Examples

Imagine a warp consists of 8 threads instead of 32.

| Thread Id | Memory Location | Size |
|:---------:|:---------------:|:----:|
| 0 | 0x80 | 4 |
| 1 | 0x88 | 4 |
| 2 | 0x84 | 4 |
| 3 | 0x9C | 4 |
| 4 | 0x90 | 4 |
| 5 | 0x94 | 4 |
| 6 | 0x98 | 4 |
| 7 | 0x8C | 4 |

A permutation of the memory locations still results in one transaction of 128 bytes.

# Examples

| Thread Id | Memory Location | Size |
|:---------:|:---------------:|:----:|
| 0 | 0x80 | 4 |
| 1 | 0x24 | 4 |
| 2 | 0x88 | 4 |
| 3 | 0x8C | 4 |
| 4 | 0x90 | 4 |
| 5 | 0x94 | 4 |
| 6 | 0x98 | 4 |
| 7 | 0x9C | 4 |

Because the 1st thread is now reading memory not in the transaction starting from warp 0, another transaction of size 128 bytes is issued for this remaining memory. Most of the bytes read in this second transaction will not be used.

# Quiz

| Thread Id | Memory Location | Size |
|:---:|:---:|:---:|
| 0 | 0x00 | 8 |
| 1 | 0x08 | 8 |
| 2 | 0x10 | 8 |
| 3 | 0x18 | 8 |
| 4 | 0x40 | 8 |
| 5 | 0x48 | 8 |
| 6 | 0x50 | 8 |
| 7 | 0x58 | 8 |

Which memory transactions are issued? Imagine this pattern continues for a 32 thread warp.

# Answer

Four 128 byte memory transactions will be issued for a 32 thread warp. They start at:

- 0x00
- 0x80
- 0x100
- 0x180

# Pointers about Pointers

We can control the size of each memory transaction by changing the type of the pointer.

```
1    char *data;
2    data[0]; //will fetch 1 byte
```

```
1    float *data;
2    data[0]; //will fetch 4 bytes
```

```
1    float2 *data;
2    data[0]; //will fetch 8 bytes
```

# Pointers about Pointers

By casting a pointer of one type to another, we can change the size of the transaction.

```
1   char *data;
2   char val1 = data[tid + 0]; //fetches 1 byte
3   char val2 = data[tid + 1]; //fetches 1 byte
4   char val3 = data[tid + 2]; //fetches 1 byte
5   char val4 = data[tid + 3]; //fetches 1 byte
6
7   typedef unsigned int uint;
8   uint iVal = ((uint *)data)[tid]; //fetches 4 bytes
9                                    //in one transaction
10
11  (iVal & 0xFF000000) >> 24 == val1
12  (iVal & 0x00FF0000) >> 16 == val2
13  (iVal & 0x0000FF00) >> 8  == val3
14   iVal & 0x000000FF        == val4
```