



Information Coding / Computer Graphics, ISY, LiTH

## **Lecture 11**

# **More CUDA**



## **In this episode...**

- **Error checking**
- **Query device capabilities**
  - **CUDA events**
- **More on CUDA memory:**

**Coalescing, Constant memory, Texture memory...**

- **OpenGL integration**
- **Reduction (intro)**

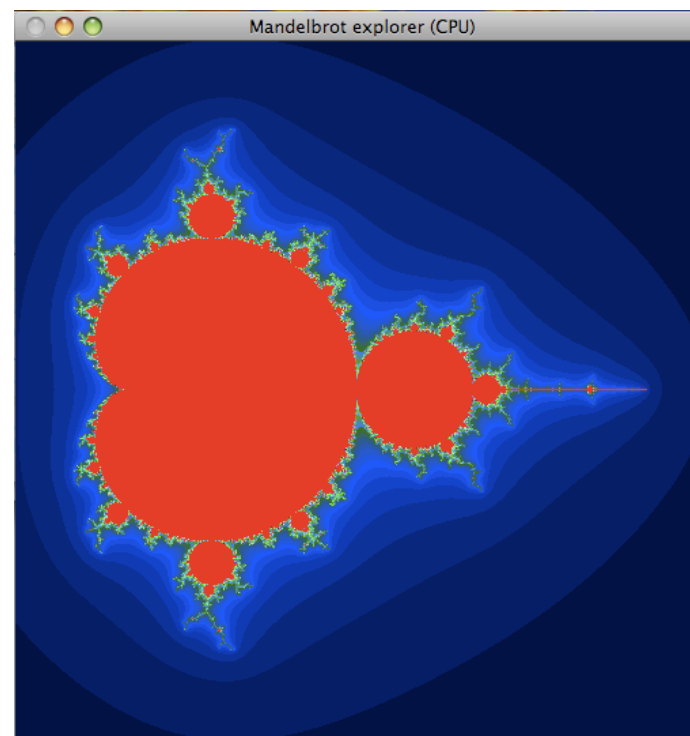


## Lab 4

**Happened earlier than usual. Everybody done or almost done?**

**Last year major change: "Mandelbrot revisited" part, to follow up lab 1.**

**Everything OK so far?**



**(Except for the drivers in the Multicore lab.)**



## **The story so far...**

- **CUDA and its language extensions**
  - **The CUDA architecture**
    - **Intro to memory**
- **Matrix multiplication example, using shared memory**



## CUDA and its language extensions

Kernel invocation `myKernel<<<>>>()`

`__global__ __device__ __host__`

`cudaMalloc(), cudaMemcpy()`

`threadIdx, blockIdx, blockDim, gridDim`

Using `nvcc`



Information Coding / Computer Graphics, ISY, LiTH

# **The CUDA architecture**

**Blocks and threads**

**Grid-block-thread hierarchy**

**Indexing data with thread/block numbers**



# **Intro to memory**

**global memory**

**shared memory**

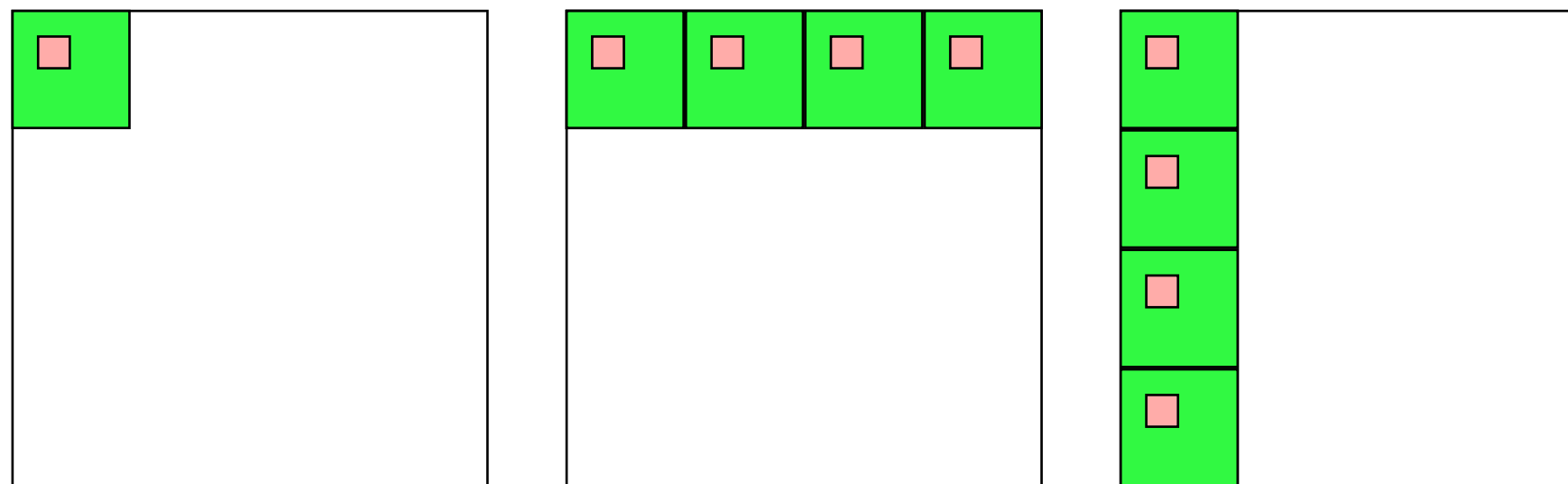
**constant memory**

**local memory**

**texture memory/texture units**



## Matrix multiplication example, using shared memory



**Huge speedup - my measly 9400M went from obvious loser to clearly faster than CPU!**





Information Coding / Computer Graphics, ISY, LiTH

**Over to today's episode:**



## **Lecture questions:**

- 1. Why can using constant memory improve performance?**
- 2. What is CUDA Events used for?**
- 3. What does coalescing mean and what should we do to get a speedup from coalescing?**
- 4. How can you efficiently calculate the maximum of a dataset in parallel?**



## **Error checking**

- **Functions returns error codes (but kernel launch does not)**
  - **cudaGetLastError()**
  - **cudaPeekLastError()**



# **Asynchronous error checking**

**Asynchronous errors can not be returned  
by the function call!**

**Call `cudaDeviceSynchronize()` and check  
its returned error code.**



## **Query devices**

**You can't trust all devices to have the same - or even similar - data.**

**New boards may have totally different data.**

**Query CUDA for a list of features using `cudaGetDeviceProperties()`**



## Example query result

```
----- Information for GeForce 9400M -----  
    Compute capability:  1.1  
Total global memory (VRAM):  259712 kB  
    Total constant Mem:  64 kB  
        Number of SMs:  2  
    Shared mem per SM:  16 kB  
    Registers per SM:  8192  
    Threads in warp:  32  
    Max threads per block:  512  
Max thread dimensions:  (512, 512, 64)  
Max grid dimensions:  (65535, 65535, 1)
```



## What is important?

Compute capability - can this board at all work with our program?

Amount of shared memory - make sure we fit.

Max threads, max dimensions - make sure we fit.

Threads in warp: A lower bound for performance.

Number of SMs: Lower bound for blocks



# Compute capability

**Essentially CUDA/architecture version number.**

- 1.0: Original release.**
- 1.1: Mapped memory, atomic operations.**
- 1.3: Double support.**
- 2.0: Fermi.**
- 3.0: Kepler.**
- 5.0: Maxwell.**





Feature Support	Compute Capability					
	1.0	1.1	1.2	1.3	2.x, 3.0	3.5
<b>(Unlisted features are supported for all compute capabilities)</b>						
Atomic functions operating on 32-bit integer values in global memory (Atomic Functions)	No	Yes				
atomicExch() operating on 32-bit floating point values in global memory (atomicExch())						
Atomic functions operating on 32-bit integer values in shared memory (Atomic Functions)	No	Yes				
atomicExch() operating on 32-bit floating point values in shared memory (atomicExch())						
Atomic functions operating on 64-bit integer values in global memory (Atomic Functions)						
Warp vote functions (Warp Vote Functions)						
Double-precision floating-point numbers	No			Yes		
Atomic functions operating on 64-bit integer values in shared memory (Atomic Functions)	No				Yes	
Atomic addition operating on 32-bit floating point values in global and shared memory (atomicAdd())						
__ballot() (Warp Vote Functions)						
__threadfence_system() (Memory Fence Functions)						
__syncthreads_count(), __syncthreads_and(), __syncthreads_or() (Synchronization Functions)						
Surface functions (Surface Functions)						
3D grid of thread blocks						
Funnel shift (see reference manual)	No				Yes	

LiTH



## Information Coding / Computer Graphics, ISY, LiTH

	FERMI GF100	FERMI GF104	KEPLER GK104	KEPLER GK110
<b>Compute Capability</b>	2.0	2.1	3.0	3.5
<b>Threads / Warp</b>	32	32	32	32
<b>Max Warps / Multiprocessor</b>	48	48	64	64
<b>Max Threads / Multiprocessor</b>	1536	1536	2048	2048
<b>Max Thread Blocks / Multiprocessor</b>	8	8	16	16
<b>32-bit Registers / Multiprocessor</b>	32768	32768	65536	65536
<b>Max Registers / Thread</b>	63	63	63	255
<b>Max Threads / Thread Block</b>	1024	1024	1024	1024
<b>Shared Memory Size Configurations (bytes)</b>	16K 48K	16K 48K	16K 32K 48K	16K 32K 48K
<b>Max X Grid Dimension</b>	$2^{16}-1$	$2^{16}-1$	$2^{32}-1$	$2^{32}-1$
<b>Hyper-Q</b>	No	No	No	Yes
<b>Dynamic Parallelism</b>	No	No	No	Yes

Compute Capability of Fermi and Kepler GPUs



# Information Coding / Computer Graphics, ISY, LiTH

<b>Compute Capability</b>	<b>1.0</b>	<b>1.1</b>	<b>1.2</b>	<b>1.3</b>	<b>2.0</b>	<b>2.1</b>	<b>3.0</b>	<b>3.5</b>
<i>SM Version</i>	sm_10	sm_11	sm_12	sm_13	sm_20	sm_21	sm_30	sm_35
<i>Threads / Warp</i>	32	32	32	32	32	32	32	32
<i>Warps / Multiprocessor</i>	24	24	32	32	48	48	64	64
<i>Threads / Multiprocessor</i>	768	768	1024	1024	1536	1536	2048	2048
<i>Thread Blocks / Multiprocessor</i>	8	8	8	8	8	8	16	16
<i>Max Shared Memory / Multiprocessor (bytes)</i>	16384	16384	16384	16384	49152	49152	49152	49152
<i>Register File Size</i>	8192	8192	16384	16384	32768	32768	65536	65536
<i>Register Allocation Unit Size</i>	256	256	512	512	64	64	256	256
<i>Allocation Granularity</i>	block	block	block	block	warp	warp	warp	warp
<i>Max Registers / Thread</i>	124	124	124	124	63	63	63	255
<i>Shared Memory Allocation Unit Size</i>	512	512	512	512	128	128	256	256
<i>Warp allocation granularity</i>	2	2	2	2	2	2	4	4
<i>Max Thread Block Size</i>	512	512	512	512	1024	1024	1024	1024
<i>Shared Memory Size Configurations (bytes)</i>	16384	16384	16384	16384	49152	49152	49152	49152
<i>[note: default at top of list]</i>					16384	16384	16384	16384
							32768	32768
<i>Warp register allocation granularities</i>					64	64	256	256
<i>[note: default at top of list]</i>					128	128		



Information Coding / Computer Graphics, ISY, LiTH

# **Do I care about Compute capability?**

**While learning CUDA - not much. Stick to the basics, it works on all.**

**But if you write professional CUDA code, of course.**



# CUDA Events

## Timing!

**Two ways of timing CUDA programs:**

- **CPU timer. Synchronize at start and end.**
- **CUDA Events. Synchronize at end.**

**Synchronize? Because CUDA runs asynchronously.**



## **CUDA Events API**

**cudaEventCreate** - initialize an event variable

**cudaEventRecord** - place a marker in the queue

**cudaEventSynchronize** - wait until all markers  
have received values

**cudaEventElapsedTime** - get the time difference  
between two events