

# Practical Course: GPU Programming in Computer Vision

## CUDA Memories

Björn Häfner, Benedikt Löwenhauser, Thomas Möllenhoff

Technische Universität München  
Department of Informatics  
Computer Vision Group

Summer Semester 2017  
September 11 - October 8





## Outline

- 1 Overview of Memory Spaces
- 2 Shared Memory
- 3 Texture Memory
- 4 Constant Memory
- 5 Common Strategy for Memory Accesses



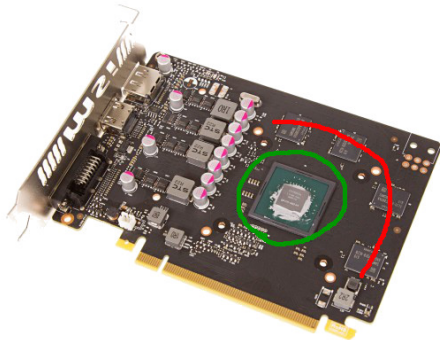


## Outline

- 1 Overview of Memory Spaces
- 2 Shared Memory
- 3 Texture Memory
- 4 Constant Memory
- 5 Common Strategy for Memory Accesses



# CUDA Memories

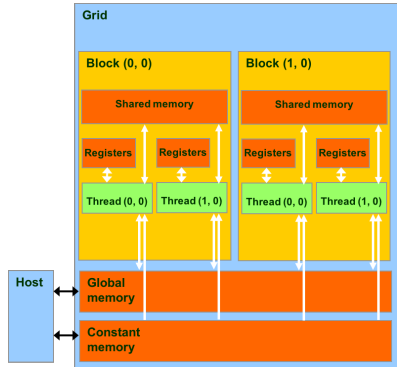


- red line is global memory (off-chip)
- green circle is the chip, contains SMs and on-chip memory



# CUDA Memories

- Each **thread** can:
  - read / write per-**thread** registers
  - read / write per-**block** shared memory
  - read / write per-**grid** global memory
  - read per-**grid** constant memory



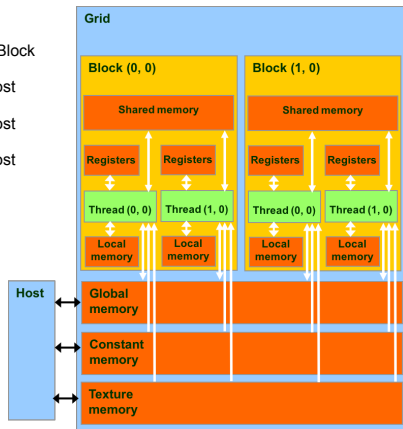


# CUDA Memories

Memory	Location	Access	Scope
Register	On-Chip	Read/Write	1 Thread
Local	Off-Chip	Read/Write	1 Thread
Shared	On-Chip	Read/Write	All Threads in 1 Block
Global	Off-Chip	Read/Write	All Threads + Host
Constant	Off-Chip	Read	All Threads + Host
Texture	Off-Chip	Read/(Write)	All Threads + Host

Other memories:

- local memory
- texture memory
- both are part of global memory





# CUDA Variable Type Qualifiers

Variable declaration	Memory	Scope	Lifetime
<code>int var;</code>	register	thread	thread
<code>int array_var[10];</code>	local memory	thread	thread
<code>__shared__ int shared_var;</code>	shared memory	block	block
<code>__device__ int global_var;</code>	global memory	grid	application
<code>__constant__ int constant_var;</code>	constant memory	grid	application

## Rules of thumb:

- scalar variables without qualifier reside in a register
- (compiler may spill to local memory)
- array variables without qualifier reside in local memory



# CUDA Variable Type Performance

Variable declaration	Memory	Penalty
<code>int var;</code>	register	1x
<code>int array_var[10];</code>	local	100x
<code>__shared__ int shared_var;</code>	shared	1x
<code>__device__ int global_var;</code>	global	100x
<code>__constant__ int constant_var;</code>	constant	1x

- **scalar variables** reside in fast, on-chip registers
- **shared memory** resides in fast, on-chip memories
- **thread local arrays & global variables** reside in off-chip memory (though cached on modern architectures)
- **constant variables** reside in cached off-chip memory



# CUDA Variable Type Scale

Variable declaration	Instances	Visibility
<code>int var;</code>	100,000s	1
<code>int array_var[10];</code>	100,000s	1
<code>__shared__ int shared_var;</code>	100s	100s
<code>__device__ int global_var;</code>	1	100,000s
<code>__constant__ int constant_var;</code>	1	100,000s

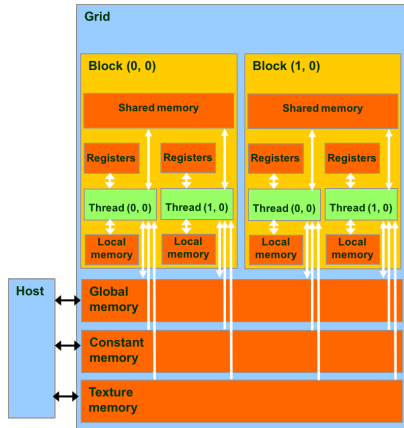
- 100,000s per-thread variables, read/write by 1 thread
- 100s shared variables, each read/write by 100s of threads
- 1 global variable, is read/write by 100,000s of threads
- 1 constant variable, is read by 100,000s of threads



## Local Memory

Compiler might place variables in **local memory**:

- too many register variables
- a structure consumes too much register space
- an array is not indexed with constant quantities, i.e., when the addressing of the array is not known at compile time







## Outline

- 1 Overview of Memory Spaces
- 2 Shared Memory**
- 3 Texture Memory
- 4 Constant Memory
- 5 Common Strategy for Memory Accesses





## Global and Shared Memory

- **Global memory is located off-chip**
  - high latency (often the bottleneck of computation)
  - important to minimize accesses
  - cached (L1 and L2) for reasonably modern GPUs (not cached for CC 1.x GPUs)
  - difficulty: try to coalesce accesses (more later)
- **Shared memory is on-chip**
  - low latency
  - like a user-managed per-SM cache
  - GPUs in lab: 48kb per multiprocessor
  - minor difficulty: try to minimize or avoid bank conflicts (more tomorrow)





## Take Advantage of Shared Memory

- Hundreds of times faster than global memory
- Threads can cooperate via shared memory
- Avoid multiple loads of same data by different threads of the block
- Use one/a few threads to load/compute data shared by all threads in the block



## Shared Memory: Example

```
1  // forward differences discretization of derivative
2  __global__ void diff_global(float *result, float *input, int n)
3  {
4      int i = threadIdx.x + blockDim.x*blockIdx.x;
5
6      float res = 0;
7      if (i+1 < n)
8      {
9          // each thread loads two elements from global memory
10         float xplus1 = input[i+1];
11         float x0 = input[i];
12         res = xplus1 - x0;
13     }
14     if(i<n) result[i] = res;
15 }
```

- `input[i]` is read by thread `i-1` and by thread `i`
- Idea: eliminate redundancy by sharing data



## Shared Memory: Example

```
1  #define BLOCK_SIZE 32
2
3  // forward differences discretization of derivative
4  __global__ void diff_shared(float *result, float *input, int n)
5  {
6      int i = threadIdx.x + blockDim.x*blockIdx.x;
7      int iblock = threadIdx.x; // local "block" version of i
8
9      // allocate shared array, of constant size BLOCK_SIZE
10     __shared__ float sh_data[BLOCK_SIZE];
11
12     // each thread reads one element and writes into sh_data
13     if (i<n) sh_data[iblock] = input[i];
14
15     // ensure all threads finish writing before continuing
16     __syncthreads();
17     ...
```



## Shared Memory: Example

```
1    ...
2
3    float res = 0;
4    if (i+1 < n)
5    {
6        // handle thread block boundary
7        int xplus1 = (iblock+1<blockDim.x)? sh_data[iblock+1] : input[i+1];
8        int x0 = sh_data[iblock];
9        res = xplus1 -x0;
10   }
11   if(i<n) result[i] = res;
12 }
```



## Shared Memory: Example

```
1  __global__ void diff_global(float *result,  
2  float *input, int n) {  
3  int i = threadIdx.x + blockDim.x*blockIdx.x;  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  float res = 0;  
19  if (i+1 < n) {  
20      // each thread loads  
21      // two elements from global memory  
22      float xplus1 = input[i+1];  
23      float x0 = input[i];  
24      res = xplus1 - x0;  
25  }  
26  if(i<n) result[i] = res;  
27  }
```

```
1  __global__ void diff_shared(float *result,  
2  float *input, int n) {  
3  int i = threadIdx.x + blockDim.x*blockIdx.x;  
4  // local "block" version of i  
5  int iblock = threadIdx.x;  
6  
7  // allocate shared array of size BLOCK_SIZE  
8  __shared__ float sh_data[BLOCK_SIZE];  
9  
10 // each thread reads one element  
11 // and writes into sh_data  
12 if (i<n) sh_data[iblock] = input[i];  
13  
14 // ensure all threads finish  
15 // writing before continuing  
16 __syncthreads();  
17  
18 float res = 0;  
19 if (i+1 < n) {  
20     // handle thread block boundary  
21     int xplus1 = (iblock+1<blockDim.x)?  
22         sh_data[iblock+1] : input[i+1];  
23     int x0 = sh_data[iblock];  
24     res = xplus1 - x0;  
25 }  
26 if (i<n) result[i] = res;  
27 }
```



# Shared Memory: Dynamic Allocation

Size known at compile time:

```
1  __global__ void kernel (...)  
2  {  
3      ...  
4      __shared__ float s_data[BLOCK_SIZE];  
5      ...  
6  }  
7  
8  int main(void)  
9  {  
10     ...  
11  
12  
13  
14     kernel <<<grid,block>>> (...);  
15     ...  
16 }
```

Size known at kernel launch:

```
1  __global__ void kernel (...)  
2  {  
3      ...  
4      extern __shared__ float s_data[];  
5      ...  
6  }  
7  
8  int main(void)  
9  {  
10     ...  
11     // allocate enough shared memory  
12     size_t smBytes = block.x * block.y * block.z  
13                     * sizeof(float);  
14     kernel <<<grid,block,smBytes>>> (...);  
15     ...  
16 }
```

## ■ Always use dynamic allocation

- flexibility w.r.t. maximal block size: can specify at run time
- no waste of resources: more blocks can run in parallel





## Shared Memory: Synchronization

- `__syncthreads()` ;
- Synchronizes all threads in a block
  - generates a barrier synchronization instruction
  - no thread can pass this barrier until all threads in the block reach it
  - used to avoid Read-After-Write / Write-After-Read / Write-After-Write hazards for shared memory accesses
- Allowed in conditional code („if“, „while“, etc.) only if the conditional is uniform across the block
  - e.g. every thread follows the same „if“- or „else“-path



## Shared Memory: Synchronization

- Always use `__syncthreads()`; after writing to the shared memory to ensure that data is ready for accessing
- Don't synchronize or serialize unnecessarily

```
1  __global__ void share_data(int *input)
2  {
3      extern __shared__ int data[];
4      data[threadIdx.x] = input[threadIdx.x];
5      __syncthreads();
6      // the state of the entire data array
7      // is now well-defined for all threads in the block
8  }
```





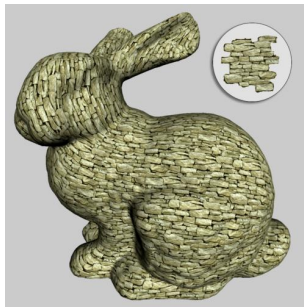
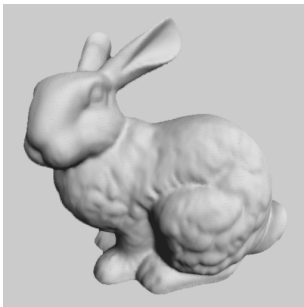
## Outline

- 1 Overview of Memory Spaces
- 2 Shared Memory
- 3 Texture Memory**
- 4 Constant Memory
- 5 Common Strategy for Memory Accesses



# Texture Memory

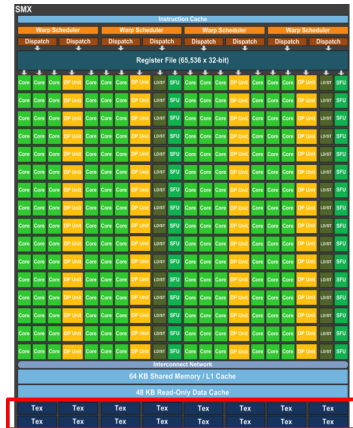
- GPUs were originally intended to do computer graphics
- Still contain specialized hardware for frequent operations such as texture mapping





# Textures

- Texture memory is part of global memory
- (Read-only), cached
- Global memory reads are performed through **extra hardware** for texture manipulation







## Why use Textures?

- Data is cached, cache is optimized for 2D spatial locality
- Filtering (interpolation) with no additional costs
  - linear / bilinear / trilinear
- Wrap modes with no additional costs
  - for „out-of-bounds“ addresses
- Addressable in 1D, 2D, or 3D
  - using integer or normalized  $[0,1)$  coordinates



## Texture Usage: Overview

- Host (CPU) code:
  - allocate global memory
  - create a texture reference object
  - bind the texture reference to the allocated memory
  - use texture reference in kernels
  - when done: unbind texture reference
- Device (GPU) code:
  - fetch (read) using texture reference
  - `tex1D(texRef, x)`, `tex2D(texRef, x, y)`,  
`tex3D(texRef, x, y, z)`
- Work best together with `cudaArray` (more later)



## Texture Usage: Texture Reference

- Define a texture reference at file scope:

```
texture <Type, Dim, ReadMode> texRef;
```

- Type: int, float, float2, float4, ...
- Dim: 1, 2, or 3, data dimension
- ReadMode:
  - cudaReadModeElementType for integer-valued textures:  
return value as is
  - cudaReadModeNormalizedFloat for integer-valued textures:  
normalize value to [0,1)



## Texture Usage: Set Parameters

- Set boundary conditions for x and y
  - `texRef.addressMode[0] = cudaAddressModeClamp;`
  - `texRef.addressMode[1] = cudaAddressModeClamp;`
  - `cudaAddressModeClamp, cudaAddressModeWrap`
- Enable/disable filtering
  - `texRef.filterMode = cudaFilterModePoint;`
  - `cudaFilterModePoint, cudaFilterModeLinear`
- Set whether coordinates are normalized to  $[0, 1)$ 
  - `texRef.normalized = false;`





## Texture Usage: Bind and Unbind

### ■ Bind texture to array:

```
cudaBindTexture2D(NULL, &texRef, ptr, &desc,  
width, height, pitch)
```

- ptr: pointer to allocated array in global memory
- width: width of array
- height: height of array
- pitch: pitch of array in bytes, if ptr was allocated using `cudaMalloc`, this is `width*sizeof(ptr[0])`
- desc: number of bits for each texture channel, e.g., `cudaCreateChannelDesc<float>()`

### ■ Unbind texture:

```
cudaUnbindTexture(texRef);
```



## Textures: Example

```
1  texture<float,2,cudaReadModeElementType> texRef; // at file scope
2
3  __global__ void kernel (...)
4  {
5      int x = threadIdx.x + blockDim.x*blockIdx.x;
6      int y = threadIdx.y + blockDim.y*blockIdx.y;
7      float val = tex2D(texRef, x+0.5f, y+0.5f); // add 0.5f to get center of pixel
8      ...
9  }
10
11 int main()
12 {
13     ...
14     texRef.addressMode[0] = cudaAddressModeClamp; // clamp x to border
15     texRef.addressMode[1] = cudaAddressModeClamp; // clamp y to border
16     texRef.filterMode = cudaFilterModeLinear; // linear interpolation
17     texRef.normalized = false; // access as (x+0.5f,y+0.5f), not as ((x+0.5f)/w,(y+0.5f)/h)
18     cudaChannelFormatDesc desc = cudaCreateChannelDesc<float>();
19     cudaBindTexture2D(NULL, &texRef, d_ptr, &desc, w, h, w*sizeof(d_ptr[0]));
20     kernel <<<grid,block>>> (...);
21     cudaUnbindTexture(texRef);
22     ...
23 }
```





## Textures: cudaArray

- Textures can use memory stored in a space filling curve



- Better texture cache hit rate due to improved 2D locality
- Copying data to a cudaArray will cause it to be formatted to such a curve
- { cudaArray, cudaMallocArray, cudaMemcpyToArray, cudaBindTextureToArray, cudaFreeArray }

```
1  cudaArray* cuArray;  
2  cudaMallocArray(&cuArray, &channelDesc, width, height);  
3  cudaMemcpyToArray(cuArray, 0, 0, h_data, size, cudaMemcpyHostToDevice);  
4  cudaBindTextureToArray(texRef, cuArray, channelDesc);  
5  ...  
6  cuFreeArray(cuArray); // free device memory
```



## Surface Memory

- Device code ( $CC \geq 2.0$ ) can read/write to `cudaArray` by using **surfaces**
- See CUDA SDK example `simpleSurfaceWrite`
- Surface operations have
  - no interpolation or data conversion
  - but some boundary handling
- Some caveats:
  - texture cache is not notified of `cudaArray` modifications
  - similarly, texture cache is also not notified of global memory modifications
  - start new kernel to pick up modifications
  - surface write/reads take x coordinates in byte size





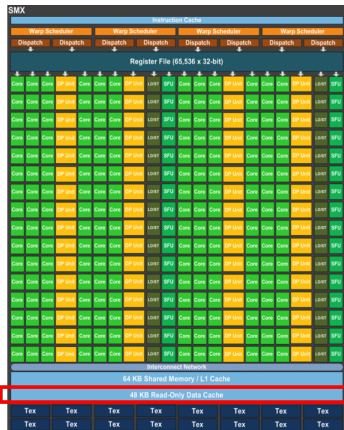
## Outline

- 1 Overview of Memory Spaces
- 2 Shared Memory
- 3 Texture Memory
- 4 Constant Memory**
- 5 Common Strategy for Memory Accesses



# Constant Memory

- Part of global memory
- Read-only, cached
  - cache is dedicated
  - will not be overwritten by other global memory reads
- Fast!
- Limited size, use it to store a few crucial parameters (convolution kernel,  $4 \times 4$  camera matrix, ...)





## Constant Memory

- Defined at file scope
- Qualifier: `__constant__`
- Examples:
  - `__constant__ float myparam;`
  - `__constant__ float constKernel[KERNEL_SIZE];`
  - array size must be known, no dynamic allocation possible
- Reading only on device:  
`float val = myparam; val = constKernel[0];`
- Writing only on host:  
`cudaMemcpyToSymbol(constKernel, h_ptr, szBytes);`





## Outline

- 1 Overview of Memory Spaces
- 2 Shared Memory
- 3 Texture Memory
- 4 Constant Memory
- 5 Common Strategy for Memory Accesses**





## Global Memory: Coalescing

- Global memory access is **slow** (400-800 clock cycles)
- Hardware coalesces (combines) memory accesses
  - chunks of size 32 B, 64 B, 128 B
  - aligned to multiples of 32 B, 64 B, 128 B, respectively
- Coalescing is per warp
  - each thread reads a char:  $1\text{B} \cdot 32 = 32\text{ B chunk}$
  - each thread reads a float:  $4\text{B} \cdot 32 = 128\text{ B chunk}$
  - each thread reads a int2:  $8\text{B} \cdot 32 = 2 \cdot 128\text{ B chunks}$





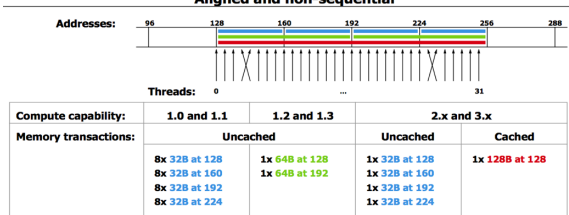
## Global Memory: Coalescing

- Make sure threads within a warp access
  - a contiguous memory region
  - as few 128 B segments as possible
- Huge performance hit for non-coalesced accesses
  - memory accesses per warp will be **serialized**
  - worst case: reading chars from random locations

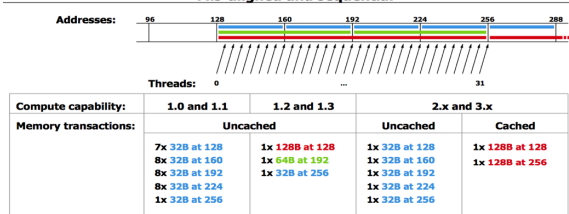


# Global Memory: Coalescing

## Aligned and non-sequential



## Mis-aligned and sequential





# The Most Important CUDA Optimization

Minimize the number of global memory accesses

- they are the slowest operations
- essentially the only reason for slow kernel run time
- if you access global memory, do it coalesced

Rules of thumb:

- neighboring threads must access neighboring elements
  - `array[threadId.x + blockDim.x * blockIdx.x]`
- (two `float` arrays are better than one `float2` array)
  - use layered memory layout for multi-channel images
- value is used a lot in same thread: load in local variable
  - even if used just more than once
- if one value is used by lots of threads: shared memory
  - but if used only by 2 or so threads, global mem is still OK





## Recommended Further Reading

CUDA Programming Guide (linked on course page)

- Appendix B.1 – B.4
- Chapter 3, sections 3.2.1 – 3.2.3

Best Practices Guide (linked on course page)

- Chapter 9, section 9.2