Code Execution

Memory Types

Advanced Features

Parallel Algorithms

Compute Shaders

# Realtime Computer Graphics on GPUs GPGPU I

### Jan Kolomazník

Department of Software and Computer Science Education Faculty of Mathematics and Physics Charles University in Prague



Computer Graphics Charles University

|--|

# Introduction

Introduction

0000000

## VLSI technology evolves very quickly

Code Execution

Frequency scaling slowed after 2005 (2-3GHz)

Memory Types

Advanced Features

Parallel Algorithms

- modern desktop parts now boost up to 6 GHz
- Current trend multi-core approach
  - parallelizable tasks
  - changes in programming techniques (traditionally all the algorithms used to be sequential)
  - High-performance computing (HPC) community practices are becoming common in consumer segment

Introduction CUDA Code Execution Memory Types Advanced Features Parallel Algorithms

# **GPU** ADVANCES

#### graphic hardware evolution

- 1. single-purpose
- 2. configurable
- 3. programmable
- currently almost arbitrary algorithm can run on a GPU
  - code size and memory are the only limits
- but GPU is efficient for special class of algorithms only

Code Execution

Introduction

00000000

- running many sequential algorithms
- two- to many-core CPUs
- full x86-64 instruction set, SSE, AVX extensions

Memory Types

Advanced Features

Parallel Algorithms

- example 1: Intel Core i9 (Sky Lake X, Sep 2017)
  - 14 nm, 24MB L3 cache, TDP: 165W
  - 18 cores + hyperthreading = 36 computing threads
  - out-of-order execution
- example 2: Intel Xeon 6 "Sierra Forest" 6700E (2024)
  - Up to 288 Crestmont E-cores, 96MB L3 cahe, TDP: 330W

# MANY-CORE COMPUTING

Introduction

00000000

emphasis on parallel algorithms

Code Execution

- graphic accelerators programmable GPUs
- more simple instruction sets
- very high brute computing power
- example 1: NVIDIA Quadro GP100 (Pascal, 2017)

Memory Types

Advanced Features

- 16 nm, 3584 cores, 1476MHz, TDP: 235W
- 16 GB HBM2 RAM (4096-bit bus)
- example 2: NVIDIA H100 (Hopper, 2023)
  - 132 SMs / 16 896 CUDA cores, 1.9GHz, TDP: 700W
  - 80 GB HBM3 RAM 3 TB/s

Compute Shaders

Parallel Algorithms

Introduction CUDA Code Execution Memory Types Advanced Features Parallel Algorithms Compute Shaders

# **FLOPS**





# CPU vs. GPU



_												
_												
-												
-												
-				_								
-				_								
-												
DRAM												

GPU

Introduction CUDA Code Execution Memory Types Advanced Features Par

Parallel Algorithms

Compute Shaders

# MEMORY BANDWIDTH



Introduction	CUDA •000000000	Code Execution	Memory Types	Advanced Features	Parallel Algorithms	Compute Shaders

# CUDA

# Introduction CUDA Code Execution Memory Types Advanced Features Parallel Algorithms Compute Shaders

- ► GPGPU (history) computing in shaders, data in textures
  - more elegant access to GPU resources was needed
- 2007: NVIDIA Tesla architecture (G80, GeForce 8800)
  - general model for parallel programming
  - computing thread hierarchy
  - barriers for synchronization
  - atomic operations
- CUDA: Compute Unified Device Architecture
  - C for CUDA: programming language

# **CUDA ARCHITECTURE**

00000000 0000000

Introduction

CUDA

 CUDA GPU: group of highly threaded streaming multiprocessors (SM)

Memory Types

Advanced Features

Parallel Algorithms

**Compute Shaders** 

- number of SM per GPU depends on GPU generation
- each SM has a set of streaming processors (SP) CUDA cores
  - G80: 8 SP per SM
  - GF100: 32 SP per SM

Code Execution

- GP100: 64 SP per SM
- GA100: 64 SP per SM
- GH100: 128 SP per SM
- SP within one SM share control circuits, instruction decoder and instruction cache
- SIMT (Single Instruction Multiple Threads) warp

Code Execution 00000000 0000000

Memory Types

Advanced Features

Parallel Algorithms

**Compute Shaders** 

# **GPU SCHEME**

CUDA



Code Execution 00000000 0000000

CUDA

Memory Types

Advanced Features

Parallel Algorithms

**Compute Shaders** 

# STREAMING MULTIPROCESSOR

SM	M														
	Instruction Cache														
	Instruction Buffer Instruction Buffer														
			Warp So	heduler							Warp So	heduler			
Dispatch Unit Dispatch Unit Dispatch Unit Dispatch Unit															
Register File (32,768 x 32-bit) Register File (32,768 x 32-bit)															
Core	Core	DP Unit	Core	Core	DP Unit		SFU	Core	Core	DP Unit	Core	Core	DP Unit		SFU
Core	Core	DP Unit	Core	Core	DP Unit		SFU	Core	Core	DP Unit	Core	Core	DP Unit		SFU
Core	Core	DP Unit	Core	Core	DP Unit		SFU	Core	Core	DP Unit	Core	Core	DP Unit		SFU
Core	Core	DP Unit	Core	Core	DP Unit	LD/ST	SFU	Core	Core	DP Unit	Core	Core	DP Unit	LD/ST	SFU
Core	Core	DP Unit	Core	Core	DP Unit	LD/ST	SFU	Core	Core	DP Unit	Core	Core	DP Unit	LD/ST	SFU
Core	Core	DP Unit	Core	Core	DP Unit	LD/ST	SFU	Core	Core	DP Unit	Core	Core	DP Unit	LD/ST	SFU
Core	Core	DP Unit	Core	Core	DP Unit	LD/ST	SFU	Core	Core	DP Unit	Core	Core	DP Unit	LD/ST	SFU
Core	Core	Unit	Core	Core	Unit	LD/ST	SFU	Core	Core	Unit	Core	Core	Unit	LD/ST	SFU
	Texture / L1 Cache														
						6	4KB Sha	red Memo	ry						

CUDA

000000000

## Specification of set of HW features

Code Execution

Memory Types

Advanced Features

Parallel Algorithms

- Instruction availability
- Memory sizes
- Synchronization options
- Mostly backward compatible
- Compilation must target specific feature set
- Each device supports specific Compute Capability version

# **PROGRAM STRUCTURE I**

00000000 0000000

Introduction

CUDA

- one or more phases (one phase = one kernel)
- deployment on CPU or GPU

Code Execution

CUDA source code can contain both parts

Memory Types

Advanced Features

Parallel Algorithms

- nvcc (NVIDIA C compiler) separates these parts during compilation
- Ianguage: from ANSI C to complete C++ support
- GPU code = kernel
  - deployed on thousands of threads
  - GPU threads are much more light-weighted (thread creation, ctxsw a couple of machine cycles)

**PROGRAM STRUCTURE II** 

Code Execution

Memory Types

CUDA

000000000



Advanced Features

Parallel Algorithms

00000000 0000000

CUDA

## Keywords \_\_global\_\_, \_\_host\_\_, \_\_device\_\_\_

Memory Types

Advanced Features

Parallel Algorithms

Compute Shaders

where the code can run

Code Execution

from where it can be called

Keyword	Runs on	Executable from
global float KernelFunc()	GPU	host, GPU (CC 3.5)
device float DeviceFunc()	GPU	GPU
host float HostFunc()	host	host
hostdevice		
float HostDeviceFunc()	host, GPU	host, GPU

	Introduction 00000000	CUDA 000000000	Code Execution	Memory Types	Advanced Features	Parallel Algorithms	Compute Shaders
--	--------------------------	-------------------	----------------	--------------	-------------------	---------------------	-----------------

# Code Execution

Introduction CUDA Code Execution Memory Types Advanced Features Parallel Algorithms Compute Shaders

```
// Kernel definition
___qlobal___ void VecAdd(float *A, float *B, float *C)
{
   int i = threadIdx.x:
   C[i] = A[i] + B[i];
int main()
   . . .
   // Run kernel in N threads
   VecAdd<<<1, N>>>(A, B, C);
   . . .
```

# KERNEL EXECUTION

Code Execution

0000000

MyKernel<<<pre>CgridSize, blockSize, dynamicSharedMemorySize, streamID>>>(arg1, arg2, ...);

Advanced Features

Parallel Algorithms

**Compute Shaders** 

Memory Types

- ► gridSize
  - int or dim3
  - Specifies 3D structure of thread blocks
- blockSize
  - int or dim3
  - Specifies 3D structure of threads in block
- dynamicSharedMemorySize
  - amount of dynamically allocated shared memory
  - 0 is default
- streamID
  - Which stream is used for execution
  - 0 is default

 Usual approach – map grid and blocks on input data

Code Execution

0000000

Memory Types

Advanced Features

- Block:
  - Executed on single SM
  - Cannot be removed until finished
- Threads in warp:
  - Set of threads that all share the same code
  - Follow the same execution path (masking execute on branching)



Parallel Algorithms

AUTOMATIC SCALABILITY

Code Execution

0000000

Memory Types



Advanced Features

Parallel Algorithms

## SIMPLE EXAMPLE IMPROVED

00000000 0000000

Code Execution

```
// Kernel definition
___qlobal___ void VecAdd(float *A, float *B, float *C, int N)
   int i = blockIdx.x * blockDim.x + threadIdx.x;
   if (i < N) {
      C[i] = A[i] + B[i];
}
int main()
   . . .
   constexpr int M = 256;
   // Run kernel in M threads per block
   VecAdd<<<(N + M - 1) / M, M>>>(A, B, C, N);
   . . .
```

Memory Types Advanced Features Parallel Algorithms



- CUDA Stream queue of commands (kernel execution, memory transfers, event)
- Commands in stream serialized
- Different streams possible concurrency
- Default stream 0 always exists (can be per thread)
- cudaStreamCreate()
- Synchronization:
  - cudaStreamSynchronize(stream)
  - Event system

	Introduction 00000000	CUDA 000000000	Code Execution	Memory Types •00000	Advanced Features	Parallel Algorithms	Compute Shaders
--	--------------------------	-------------------	----------------	------------------------	-------------------	---------------------	-----------------

# Memory Types

Code Execution

- Normal memory RAM
- By default cannot be accessed from device

Memory Types

000000

Advanced Features

Parallel Algorithms

**Compute Shaders** 

Must be copied to device memory

GLOBAL MEMORY

Actual GPU memory

Code Execution

Used as normal linear memory – pointer arithmetics

Memory Types

00000

- Management:
  - cudaMalloc(), cudaMallocPitch(), cudaMalloc3D()

Advanced Features

Parallel Algorithms

- cudaMemcpy(), cudaMemcpyToSymbol()
- Read from kernel can take hundreds of cycles

TEXTURE MEMORY

Allow usage of texturing HW:

Code Execution

- Spatial caching
- Filtering
- Limited by predefined element types (colors)

Memory Types

00000

Advanced Features

Parallel Algorithms

**Compute Shaders** 

No custom structures

# SHARED MEMORY

- Same space as L1 cache
  - Division is customizable

Code Execution

- \_\_shared\_\_keyword
- Shared by threads in block
- Use when same value accessed is multiple times in block execution (not necessarily by same thread)

Memory Types

000000

Advanced Features

Parallel Algorithms

Advanced Features

Parallel Algorithms

**Compute Shaders** 

# REGISTERS

- 32-bit registers
- Divided between active warps
- Shared memory + registers limit occupancy:
  - Number of active warps vs. max possible warps on SM
- Limiting number of registers may lower performance
  - May be necessary to run at least 1 block on SM

Introduction 00000000	CUDA 000000000	Code Execution	Memory Types	Advanced Features	Parallel Algorithms	Compute Shaders

# **Advanced Features**

# UNIFIED MEMORY

 Single memory address space accessible from any processor in a system

Advanced Features

0000

Parallel Algorithms

**Compute Shaders** 

Replace malloc(), new with calls to cudaMallocManaged()

Memory Types

- New on Kepler architecture
- Pascal:
  - 49-bit virtual addressing

Code Execution

- On-demand page migration
- Oversubscription to GPU memory

DA Code Execution

Memory Types

Advanced Features

Parallel Algorithms

Compute Shaders

# DYNAMIC PARALLELISM I

#### too coarse



too fine



## DYNAMIC PARALLELISM II

Code Execution

Kernels can be executed from kernels on device

Memory Types

Advanced Features

00000

Parallel Algorithms

- Parent kernel waits until children finishes
- Allows adaptive thread execution

GRAPHICAL APIS CONNECTION

Code Execution

CUDA

Introduction

 Acquire access to memory used by buffer objects in OpenGL, DX, Vulcan

Advanced Features

0000

Parallel Algorithms

**Compute Shaders** 

cudaGraphicsMapResources(...),...

Memory Types

- Read or update the memory
- Unmap the resource

Introduction 00000000	CUDA 000000000	Code Execution	Memory Types	Advanced Features	Parallel Algorithms ●0000000	Compute Shaders

# **Parallel Algorithms**

# Introduction CUDA Code Execution Memory Types Advanced Features Parallel Algorithms Compute Shaders

## INTRODUCTION

### Algorithms for massively parallel architectures:

- Often bottom up design
- Shallow datastructures
- Memory access patterns considered first
- Try to make all operations local only
- Problem reformulation:
  - Search for possible constrains
  - Solve dual problem
  - Cellular automata
  - ▶ ...

BASIC META-ALGORITHMS

Code Execution

Memory Types

Advanced Features

Parallel Algorithms

00000000

**Compute Shaders** 

#### ► Map :

- ForEach (inplace?)
- Transform
- Spatial filters with limited support:
  - Convolution
  - Morphological operations

REDUCE (FOLD)

 Associative binary operation to combine input elements into single value:

Advanced Features

Parallel Algorithms

0000000

**Compute Shaders** 

Memory Types

- Sum
- Multiplication
- Min/Max
- On GPU:
  - 1. Tree-based approach used within each thread block
  - 2. Global sync (multiple kernels)
  - 3. Reduce block results

Code Execution



- Optimizations:
  - Prevent thread idling
  - Shared memory access patterns
  - Ref: Mark Harris Optimizing Parallel Reduction in CUDA (30x speedup)



- Associative binary operation to combine input elements in front of each of the processed elements
- Inclusive vs. exclusive
- Similar implementation and optimizations as reduction

# EXAMPLE: SELECT ELEMENTS BY INDICATOR

Memory Types

Advanced Features

- Task:
  - Output elements selected by predicate
- Naive approach:
  - Adding Elements to output array directly
  - Bottleneck size update

Code Execution

- Better solution:
  - Two level solution
  - Store local selection into shared memory
  - Update global output size once per block
- Advanced solution:
  - Run parallel prefix-sum on indicator set (0/1 for each element)
  - Computes indices for all selected elements and final size
  - Final step store all selected elements on precomputed positions

Compute Shaders

Parallel Algorithms

## EXAMPLE: INTEGRAL IMAGES

Code Execution

Introduction

- Apply prefix-sum to rows and columns
- Sum/average queries for rectangular regions with constant complexity

Advanced Features

Memory Types

**Parallel Algorithms** 

0000000

- Used in feature detectors:
  - Haar features
  - Blur filter approximation



# DATASTRUCTURES

#### Basic categories:

- Dense arrays
- Hash tables
- Sparse structures

Code Execution

- Matrices
- Graphs
- Different criteria:
  - Read/write
  - Incremental changes vs. Rebuild from scratch

Memory Types

Advanced Features

Parallel Algorithms

0000000

- Space waste
- Two level design:
  - Local datastructure living in shared memory
  - Main datastructure living in global memory
  - Write local instances at the end of computation

Introduction 00000000	CUDA 000000000	Code Execution	Memory Types	Advanced Features	Parallel Algorithms	Compute Shaders

# Introduction CUDA Code Execution Memory Types Advanced Features Parallel Algorithms Compute Shaders

# MOTIVATION

## Why not OpenCL or CUDA?

- One API for graphics and GP processing
- Avoid interop
- Avoid context switches
- You already know GLSL
- APIs:
  - Core since OpenGL 4.3
  - Part of OpenGL ES 3.1
  - Vulcan

WHERE IT BELONGS?

Code Execution

# **OpenGL 4.3 with Compute Shaders**

Memory Types

Advanced Features

Parallel Algorithms





# Write compute shader in GLSL

- Define memory resources
- Write main()function
- Initialization
  - Allocate GPU memory (buffers, textures)
  - Compile shader, link program
- Run it
  - Bind buffers, textures, images, uniforms
  - Call glDispatchCompute(...)

# SAMPLE COMPUTE SHADER

Code Execution

Introduction

CUDA

```
#version 430
layout(local_size_x = 1, local_size_y = 1) in;
layout(rgba32f, binding = 0) uniform image2D img_output;
void main() {
 // base pixel colour for image
 vec4 pixel = vec4(0.0, 0.0, 0.0, 1.0);
 // get index in global work group i.e x, y position
  ivec2 pixel_coords = ivec2(ql_GlobalInvocationID.xy);
 // output to a specific pixel in the image
  imageStore(img_output, pixel_coords, pixel);
```

Memory Types Advanced Features

Parallel Algorithms

IDA Code Execution Memory Types

Advanced Features

Parallel Algorithms

Compute Shaders

## **BUILTIN VARIABLES**

in uvec3 gl\_NumWorkGroups;

in uvec3 gl\_WorkGroupID;

in uvec3 gl\_LocalInvocationID;

in uvec3 gl\_GlobalInvocationID;

in uint gl\_LocalInvocationIndex;