

GPU Computing

Lesson 3: Architecture & Strategy

Jan Lemeire

2024-2025

http://parallel.vub.ac.be/education/gpu

The modern CPU

Sequential' processor: super-scalar out-of-order pipeline





GPU strategy for massive computations

GPU architecture strategy

Light-weight threads, supported by the hardware

- Thread processors, more than 1000 active threads per core
- Switching between threads can happen in 1 cycle!

No caching mechanism, branch prediction, ...

- GPU does not try to be efficient for every program, does not spend transistors on optimization
- Simple straight-forward sequential programming should be abandoned...
- Less higher-level memory:
 - GPU: 16KB shared memory per SIMD multiprocessor
 - CPU: L2 cache contains several MB's
- Massively floating-point computation power
- RISC instruction set instead of CISC
- Transparent system organization

Modern (sequential) CPUs based on simple Von Neumann architecture



GPU processor pipeline

- ▶ 6-24 stages
- in-order execution!!
- no branch prediction!!
- no forwarding!!
- no register renaming!!
- Memory system:
 - relatively small
 - Until recently no caching
 - On the other hand: much more registers (see later)
- No program call stack and no memory stack!
 - All functions inlined
 - No recursion possible

Challenges of GPU computing





Fill the pipelines

GPUs have several pipelines which will be filled with instructions from different kernel threads through:

- 1. Running thread blocks on the different multiprocessors
- 2. Simultaneous multithreading: several hardware threads active at the same time
 - Discussed next
- 3. Single Instruction Multiple Threads (SIMT)
 - Discussed later

Architecture



GPU Architecture





VRIJE UNIVERSITEIT 1 Streaming Multiprocessor = a pipeline





Properties of different architectures

 $\begin{array}{l} \mathsf{N}(\Pi) = \# multiprocessors \\ |\omega| = warp \ size \\ \text{Group & Warp \ slots: \ maximum \ } \\ \text{thread \ blocks \ or \ warps } \end{array}$

	Fermi	Kepler	Maxwell	\mathbf{Pascal}	Tonga
$N(\Pi)$	14	4	3	10	28
f_{clock} (MHz)	1147	1032	1058	1506	1010
issue limit	1	4	4	4	1
ω	32	32	32	32	64
Resources					
Group slots	8	16	32	32	16
Warp slots	48	64	64	64	40
Local memory (KB)	48	48	64	96	64
Registers (KB)	128	256	256	256	

GPUs of our lab and architecture

Full name	Abbreviated name
NIVDIA Tesla C2050	Fermi
NIVDIA GeForce GTX 650 Ti	Kepler
NIVDIA Quadro K620	Maxwell
NIVDIA GeForce GTX 1060 6GB	Pascal
AMD Radeon R9 380	Tonga

	Fermi	Kepler	Maxwell	\mathbf{Pascal}	Tonga
$max(\gamma)$	1024	1024	1024	1024	256
max(local memory) (KB)	48	48	48	48	32
ALU count	32	192	128	128	64
SFU count	8	32	32	32	-
RAM Bandwidth (GB/s)	144	86.4	29	192	176
L2 Cache size (KB)	768	256	2048	1536	512
L2 Cache line size (B)	128	128	128	128	64
L1 Cache size (KB)	16	16	64	48	16
max(global memory) (MB)	1024	672	512	1536	2880
RAM Size (MB)	2688	2048	2048	6144	4096
#LD/STO units =	= 16	32	32	32	

https://en.wikipedia.org/wiki/CUDA

The different Nvidia architectures



CUDA Compute Capability can be queried, also in GPU Caps Viewer

VRIJE UNIVERSITEIT

1st generation: Tesla

Compute Capability = 1.x



2nd generation: Fermi

Compute Capability = 2.x







3rd generation: Kepler

Compute Capability = 3.x

мм											
	Verte	x Fetch			olyMorph Tess	ellator	0	v	iewport T	ransform	
			Attrib	ute Setup			Stream O	lutput			_
					Instructi	on Cache					
		nstructio	on Buffe	r			1	nstructi	on Buffe	ar 🛛	
		Warp Sc	heduler			Warp Scheduler					
Disparch Unit Unsparch Unit							Dispatch Unit Dispatc				ut.
Register File (16,384 x 32-bit)							Regist	er File (1	16,384 x	32-bit)	
Core	Core	Core	Core		SFU	Core	Core	Core	Core		SFU
Core	Core	Core	Core		SFU	Core	Core	Core	Core		SFU
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LDIST	SFU
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LDIST	SFU
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LDIST	SFU
Core	Core	Core	Core		SFU	Core	Core	Core	Core		SFU
Core	Core	Core	Core		SFU	Core	Core	Core	Core		SFU
Core	Core	Core	Core		SFU	Core	Core	Core	Core		SFU
					Texture /	L1 Cache					
						Tex Tex					
				Tex							
		nstructi	on Buffe	i ex			Tex	nstructi	on Buffe	Tex er	
D	lispatch Un	nstructi Warp Si	on Buffe	ir Dispatch U	nit		Tex I	nstructi Warp Se	on Buffe	Tex er Dispatch Ur	sit
D	ispatch Un	Natructi Warp St	on Buffe	Nispatch U	nit		Tex I ispatch Un	nstructi Warp Si	on Buffe	Tex er Dispatch Ur	nit
	ispatch Un F Regist	nstructi Warp Si it er File ('	on Buffe cheduler (16,384 x	Hr Dispatch U 32-bit)	nit		Tex Ispatch Un Regist	nstructi Warp Se it ter File (on Buffe cheduler 1 16,384 x	Tex er Dispatch Ur E : 32-bit)	48
Core	ispatch Un Regist Core	NSTRUCTI Warp St It er File (Core	on Buffe theduler 16,384 x Core	Nispatch U 32-bit)	nit	Core	Tex Ispatch Un Regist	nstructi Warp Se at ter File (Core	on Buffi cheduler 16,384 x Core	Tex Dispatch Ur 32-bit)	sFU
D Core Core	ispatch Un Tegist Core Core	nstructi Warp Se it er File (' Core Core	on Buffr theduler 16,384 x Core Core	I UX Pr Dispatch UX 32-bit) LD/ST	nit SFU SFU	Core	Tex ispatch Un Rogist Core Core	nstructi Warp St it it cer File (Core Core	on Buffi cheduler 16,384 x Core Core	Tex Dispatch Ur 32-bit) LDIST LDIST	sFU SFU
D Core Core Core	Ispatch Un Regist Core Core Core	Narp Si Warp Si It Core Core Core	on Buffe cheduler 16,384 x Core Core Core	Itox Dispatch U 32-bit) LDIST LDIST LDIST	SFU SFU SFU	Core Core	tex tespatch Un Rogist Core Core Core	Instructi Warp So it ier File (Core Core Core	on Buffe cheduler 16,384 x Core Core Core	Tex Pr Dispatch Ur 32-bit) LD/ST LD/ST LD/ST	SFU SFU SFU
Core Core Core Core	Ispetch Un Regist Core Core Core	Natructi Warp So It Core Core Core Core	on Buffe cheduler 16,384 x Core Core Core	I Dispatch U 32-bit) LD/ST LD/ST LD/ST	SFU SFU SFU SFU	Core Core Core	ispatch Un Regist Core Core Core	nstructi Warp St it Core Core Core Core	on Bufficheduler 16,384 x Core Core Core	Tex Pr Dispatch Ur 32-bit) LD/ST LD/ST LD/ST LD/ST	Nit SFU SFU SFU SFU
Core Core Core Core Core	Rogist Core Core Core Core Core	Warp St it Core Core Core Core Core	on Buffi theduler 16,384 x Core Core Core Core	I Dispatch U 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST	NK SFU SFU SFU SFU SFU	Core Core Core Core	Tex ispatch Un Rogist Core Core Core Core Core	nstructi Warp So it Core Core Core Core Core	on Buffe cheduler 16,384 x Core Core Core Core	Tex Dispatch Ur 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST	SFU SFU SFU SFU SFU
Core Core Core Core Core Core	ispatch Un Regist Core Core Core Core Core	It Warp St It Core Core Core Core Core Core Core Core	on Buffi cheduler 16,384 x Core Core Core Core Core Core	I UX I US I Spatch UX I DIST I DIST I DIST I DIST I DIST	SFU SFU SFU SFU SFU SFU	Core Core Core Core Core	tex ispatch Un Rogist Core Core Core Core Core Core	Warp So it Core Core Core Core Core Core Core	on Buffi cheduler 16,384 x Core Core Core Core Core Core	Tex Dispatch Ur 32-bit) LDIST LDIST LDIST LDIST LDIST	SFU SFU SFU SFU SFU SFU
Core Core Core Core Core Core Core	ispetch Un Rogist Core Core Core Core Core Core	Astruction Warp Struction for File (* Core Core Core Core Core Core Core	on Buffe cheduler 16,384 x Core Core Core Core Core Core Core	Nepatch U Standard U Standard U Standard U Standard Standard S	SFU SFU SFU SFU SFU SFU	Core Core Core Core Core Core	Tex ispatch Un Rogist Core Core Core Core Core Core Core	Warp Stite Warp Stite Core Core Core Core Core Core Core Cor	on Buffrecheduler theduler 16,384 x Core Core Core Core Core Core Core	Tex Pr Sispatch Ur 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST	SFU SFU SFU SFU SFU SFU SFU
D Core Core Core Core Core Core	Rogist Core Core Core Core Core Core Core Core	Warp Si Warp Si t Core Core Core Core Core Core Core	Core Core Core Core Core Core Core Core	er ar 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST	SFU SFU SFU SFU SFU SFU SFU SFU	Core Core Core Core Core Core Core Core	Fex Ferrit Tex Ferrit	Instruction Warp Sr Warp Sr Warp Sr Instruction Core Core Core Core Core Core	Core Core Core Core Core Core Core Core	Tex ar bispatch Ur 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST	SFU SFU SFU SFU SFU SFU SFU
D D Core Core Core Core Core Core	Rogist Core Core Core Core Core Core Core Core	Warp Sri terr File (* Core Core Core Core Core Core Core Core	Core Core Core Core Core Core Core Core	III IIII IIIIIIIIIIIIIIIIIIIIIIIIIIIII	SFU SFU SFU SFU SFU SFU SFU SFU Texture	Core Core Core Core Core Core Core	Tex Ispatch Un Rogist Core Core Core Core Core Core Core Core	Instruction Warp State File (Core Core Core Core Core Core Core Core	Core Core Core Core Core Core Core Core	Tex sr liquate Ur 2 liquate Ur	SFU SFU SFU SFU SFU SFU SFU SFU
Core Core Core Core Core Core Core	Ispatch Unit Regist Core Core Core Core Core Core Core	Warp Sri terror File (Core Core Core Core Core Core Core Core	Core Core Core Core Core Core Core	rex Pispatch U. Vispatch U. Vispatch U. Vispatch U. Vispatch U. Vispatch U. Vispatch Vispatc	SFU SFU SFU SFU SFU SFU SFU SFU Texture	Core Core Core Core Core Core Core Core	Tex Ispatch Un Rogist Core Core Core Core Core Core Core Core	Instruction Warp State Terr File (Core Core Core Core Core Core Core	Core Core Core Core Core Core Core	Tex ar ar ar ar ar ar ar ar ar ar	NT SFU SFU SFU SFU SFU SFU SFU

This is only half of an SM



5th generation: Pascal

Compute Capability = 6.x

4rd generation: Maxwell

Compute Capability = 4.x or 5.x



Compute Capability = 8.x if < 8.5: Volta

5th generation: Pascal

6th generation: Volta & Turing

Pascal SM



Turing SM



Without double precision (DP) units

7th generation: Ampere

Compute Capability = 9.x

Simultaneous multithreading



Multithreading

- Performing multiple threads of execution in parallel
 - Replicate registers, PC, etc.
 - Fast switching between threads
- Fine-grain multithreading
 - Switch threads after each cycle
 - Interleave instruction execution
 - If one thread stalls, others are executed
- Coarse-grain multithreading
 - Only switch on long stall (e.g., L2-cache miss)
 - Simplifies hardware, but doesn't hide short stalls (eg, data hazards)



verhead

Multithreading on CPU

- I process/thread active per core
- When activating another thread: *context switch*
 - Stop program execution: flush pipeline (let all instructions finish)
 - Save state of process/thread into Process Control Block : registers, program counter and operating system-specific data
 - Restore state of activated thread
 - Restart program execution and refill the pipeline
- Processor 'sees' only 1 thread
- Called Software threads



Running threads on same CPU core



- Executed one by one
- Context switch
- Thread's state in core: instruction fetch buffer, return address stack, register file, control logic/state, ...
- Supported by hardware
- Takes time!

Coarse-grain multithreading



Fine multi-threading: Hardware threads

- In several modern CPUs
 - typically 2 HW threads (Intel: hyperthreading)
- Devote extra hardware for keeping process state
- Thread switching by hardware
 - (almost) no overhead
 - within 1 cycle!
 - Instructions in flight from different threads



Simultaneous Multithreading

- In multiple-issue dynamically scheduled processor
 - Schedule instructions from multiple threads
 - Instructions from independent threads execute when function units are available
 - Within threads, dependencies handled by scheduling and register renaming
- Example: Intel Pentium-4 HyperThreading
 - Two threads: duplicated registers, shared function units and caches



Multi-Threading (MT) possibilities





Benefits of fine-grained multithreading

- Independent instructions (no bubbles)
- More time between instructions: possibility for *latency hiding*
 - Hide memory accesses
- If pipeline full
 - Forwarding not necessary
 - Branch prediction not necessary



Running a simple addition kernel



Runtime increases only when all pipelines are full (8000 threads)



The execution on a GPU



- Thread blocks are scheduled on MultiProcessors.
- Warps of active threads are scheduled on the multiprocessor



Concurrency

- Keep all processing units busy!
 - Enough threads
- All Multiprocessors (MPs)
- All Scalar Processors (SPs)
- Full pipeline of scalar processor
 - Pipeline of up to 24 stages

What determines the occupancy



Occupancy

- **Occupancy** = #*concurrent warps running on a* Multiprocessor
- A higher occupancy means that more work can be scheduled and in general a higher performance
- Limited resources will limit the number of work groups that can be simultaneously active and run concurrently:
 - 1. Registers (private memory) needed per work group
 - Each kernel's local variables are stored in register memory
 - 2. Local memory needed per work group
 - 3. Maximum number of concurrent work groups
 - 4. Maximum number of threads (work items)
- The most constrained resource determines the occupancy
 - Each Multiprocessor has resources (depends on architecture, can be queried)
 - For Pascal architecture: 256KB registers, 96KB local memory,
 - max. 32 work groups, max. 2048 threads(=64 warps)

The effect of occupancy will be studied with the **Pipeline Model**