GPU Computing



Napping Programming model to hardware - the SMs	Zoom on the SM	The SIMT Execution Model on CUDA Cores	The SIMT Execution Model on CUDA Cores				
 Any block of any grid can be mapped to any SM Image: Solution of the solution of the	 SM organizes blocks into warps 1 warp = group of 32 threads 3 warp = group of 32 threads 4 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes blocks into warps 1 warp = group of 32 threads 5 Worganizes block of 96 threads maps to: 3 (physical) warps of 32 threads 5 Worganizes blocks into warps 5 Worganizes blocks 5 Worganizes	• Divergent code paths (branching) pile up: Image: Divergent code paths (branching) pile up:	<pre>What is the latency (in term of inst) of this code in the better and worst case ? if a > 0: inst-a if b > 0: inst-b; else inst-c else: inst-d</pre>				
Soom on the CUDA cores	The SIMT Execution Model on CUDA Cores	The SIMT Execution Model on CUDA Cores	The SIMT Execution Model on CUDA Cores				
1 core = 1 thread	• Divergent code paths (branching) pile up! The second	<pre>What is the latency (in term of inst) of this code in the better and worst case ? if a > 0: inst-a if b > 0: inst-b; else inst-c else: inst-d · Best case: a > 0 is false for every thread. For all threads: inst-d · Worst case: a > 0 and b > 0 is true for some but not all threads. For all threads: inst-a</pre>	Loops i = 0 while i < thread_id: i += 1				
Warning: SIMT allows to specify the execution and branching behavior of a single thread but maps to SIMD processors!		inst-c					

inst-d



External memory: embedded GPU	GPU: on-chip memory	Memory model hierarchy (software)	Non-aligned and strided load in L1 & L2
Most GPUs today: • Same memory • May support memory coherence (GPU can read directly from CPU caches) • More contention on external memory GPU CPU Cache 16 GB/s Main memory 8GB	Cache area in CPU vs GPU:	eru sidors to Data	Q: how to make aligned-load with 2D-arrays ? ima(x,y) = y * width + x Lifetime Thread Thread Sk block Is host Is host Host Sk host Host Sk host Sk host
GPU: on-chip memory	Memory model hierarchy (hardware)	Non-aligned and strided load in L1 & L2	Shared Memory & conflicts
Cache area in CPU vs GPU:	Store State	O: how to make aligned-load with 2D-arrays ? int is the control of the control block into	Used a lot for local copy and fast-access · Organized in memory bank (accessed concurrently) · Every bank can provide 64 bits every cycle · Only two modes: · Only two mode

						1441				
						128	2008	2010	2012	2014
GPU	Registers files + caches		CPU	GPU	(mv)	64 32	-		,	\sim
NVidia Maxwell	8.3 MB	Register / Core	256	65K	loues	16	·~	~	÷	-
AMD Hawaii GPU	15.8 MB				lenal	4	-	T.	MADAG	U
Core i7 CPU	9.3 MB				2	2	-		AMD CI Intel CI Intel V	
							72 0000 77 0002 775 0000 7754	DENAME STORAGE	gtx/88	NUCC BIT

 \rightarrow GPU has many more registers but made of simpler memory

 Reduce throughput demand on main memory L1 Managed by hardware (L1, L2) or software (shared memory)

Keep frequently-accessed data Core

On CPU, caches are designed to avoid memory latency

· On GPU, multi-threading deals with memory latency

ima(x,y) = y * stride + x



2-way conflict

No conflict Conflict = Serialized access (bad!) No conflict (broadcast)

CUDA crash course