

SIGGRAPH Asia 2009 GPU Computing master Class Future Direction in GPU Computing

Toru Baji / NVIDIA Solution Architect

16th Dec., 2009

Content



• Why GPU Computing?

- Single-threaded CPU performance is no longer scaling
- CUDA can exploit "Performance = Parallelism"
- Evolution of NVIDIA GPU
- Key Milestone in GPU Computing 'Fermi'
 - Processor / Memory Architecture
 - Enhanced Thread Management
 - CUDA3.0
 - Integrated Development Environment 'Nexus'
- An NVIDIA ExaScale Machine in 2017
- Conclusion

Single-threaded CPU performance is no longer scaling (1) Moore's Law





Moore, Electronics 38(8) April, 1965

- In 1965, Gordon Moore predicted the number of Tr will grow at the rate of:
 x 2 / year
 later revised to
 x 2 / 18-months
- No prediction of CPU performance

Single-threaded CPU performance is no longer scaling (2) The End of ILP (Instruction Level Parallelism) Scaling



The increase in Tr counts was used to increasing the parallelism of the CPU instruction execution (e.g. superscalar, pipeline) and other speed-ups



Performance is measured as the time required to execute one instruction (ps/inst)

Bill Dally et al., the Last Classical Computer, ISAT Study, 2001

© 2009 NVIDIA Corporation

Single-threaded CPU performance is no longer scaling (3) Explicit Parallelism is Now Attractive



The increase in Tr counts is used to increasing the number of processors



Bill Dally et al., the Last Classical Computer, ISAT Study, 2001

© 2009 NVIDIA Corporation

CUDA can exploit Performance = Parallelism



Now it is clear that "Performance = Parallelism"

- Then what is required to exploit this Parallelism?
 - Many efficient processors
 - A programming system that abstract the parallelism
- "CUDA (Compute Unified Device Architecture"
 - GPU has hundreds of full-featured processors
 - Multi-threading architecture use the best of those massive parallel cores
 - CUDA starts from extension to the familiar C-language, and now expanded to Fortran, OpenCL and direct Compute
 - CUDA can isolate the programmer from the details of parallel programming

CUDA GPU Computing Visual Demo separate animation tool used



Hierarchical thread definition





Hide memory Latency



Ease of Programming and Performance Increase



- Ease of programming represented as "Development time"
- Performance achievement measured in GFLOPS





Programmable Shader

Unified Shader

CUDA, OpenCL, DirectX Compute etc



Key Milestone in GPU Computing 'Fermi'

- Improved peak performance
- Improved throughput efficiency
- Broader applicability
- Full integration within modern software development environment







SM Multiprocessor Architecture

- 32 CUDA cores per SM (512 total)
- 8x peak FP64 performance
 - 50% of peak FP32 performance
- Dual Thread Scheduler
- 64 KB of RAM for shared memory and L1 cache (configurable)





IEEE 754-2008 Floating Point



• IEEE 754-2008 results

- 64-bit double precision
- 32-bit single precision
- full-speed denormal operands & results
- NaNs, +/- Infinity

IEEE 754-2008 rounding

nearest even, zero, +inf, -inf

IEEE 754-2008 Fused Multiply-Add (FMA)

- D = A*B + C;
- No loss of precision
- IEEE divide & sqrt use FMA









Fused Multiply-Add (FMA): $D = A^*B + C$;

Cached Memory Hierarchy

- Configurable L1 cache per SM
 - 16KB L1\$ / 48KB Shared Memory
 - 48KB L1\$ / 16KB Shared Memory
- Shared 768KB L2 cache
- Motivation:
 - Caching captures locality, amplifies bandwidth
 - Caching more effective than Shared Memory RAM for irregular or unpredictable access
 - Ray tracing, sparse matrix multiply, physics kernels ...
 - Caching helps latency sensitive cases

Larger, Faster Memory Interface

GDDR5 memory interface

• 2x improvement in peak speed over GDDR3

• Up to 1 Terabyte of memory attached to GPU

Operate on large data sets

Unified Load/Store Addressing

Non-unified Address Space

Unified Address Space

ECC Memory Protection

• All major internal memories protected by ECC

- Register file
- L1 cache
- L2 cache

• External DRAM protected by ECC

• ECC is a must have for many computing applications

Clear customer feedback

GigaThread[™] Hardware Thread Scheduler

• Hierarchically manages tens of thousands of simultaneously active threads

• 10x faster context switching on Fermi

Overlapping kernel execution

Overlapping Kernel Execution

© 2009 NVIDIA Corporation

GigaThread Streaming Data Transfer Engine

- Dual DMA engines
- Simultaneous CPU→GPU and GPU→CPU data transfer

Fully overlapped with CPU/GPU processing

	G80	GT200	Fermi
Transistors	681 million		3.0 billion
CUDA Cores	128	240	512
Double Precision Floating Point	-	30 FMA ops/clock	256 FMA ops/clock
Single Precision Floating Point	128 MAD ops/clock	240 MAD ops/clock	512 FMA ops/clock
Special Function Units (per SM)	2	2	4
Warp schedulers (per SM)	1	1	2
Shared Memory (per SM)	16 KB	16 KB	Configurable 48/16 KB
L1 Cache (per SM)	-	-	Configurable 16/48 KB
L2 Cache	-	-	768 KB
ECC Memory Support	-	-	Yes
Concurrent Kernels	-	-	Up to 16
Load/Store Address Width	32-bit	32-bit	64-bit

CUDA Parallel Computing Architecture

© 2009 NVIDIA Corporation

OpenCL is trademark of Apple Inc. used under license to the Khronos Group Inc.

CUDA C 3.0

- Unified addressing for C and C++ pointers
- Global, shared, local addresses
- Enables 3rd party GPU callable libraries, dynamic linking
- One 40-bit address space for load/store instructions

Compiling for native 64-bit addressing

IEEE 754-2008 single & double precision

• C99 math.h support

© 2009 NVIDIA Corporation

NVIDIA Integrated Development Environment Code-named 'Nexus'

Industry's 1st IDE for massively parallel applications

 Accelerate development of CPU + GPU co-processing applications

Complete Visual Studio-integrated development environment

- C, C++, OpenCL, & DirectCompute platforms
- DirectX and OpenGL graphics APIs

🖓 Nexus CUDA Sam	nples.90 (Debugging) - Microsoft Visual Stud	lio (Administrator)							
File Edit View	Project Build Debug Nexus Tools	Test Window Help							
i 🗊 + 🛅 + 💕 🔓	a 🗿 X 🗈 🛍 🕫 - 🔍 - 🚚 - 🖳	🕨 🕨 Debug 🕞 Win32 🚽 🏄 🖉	_vbo_bu	ffer	- 🞝 🚰	* 🐋 🛠 💽 🖂 - 📜 🗊		CUDA (0,0,0), (0,0,0)	
• • • • • • •	🔶 🗐 🗊 Hex 👒 🏹 - 📜 🗄 🖏	No N	<u>ا چا</u>						
Process: [4560] GP	U - matrixMul.e - Thread: [2772376] <no n<="" td=""><td>lame> 🛛 👻 😵 Stack Frame: Module: 60956632 - [0] _Z</td><td>n - 📮</td><td></td><td></td><td></td><td></td><td></td><td></td></no>	lame> 🛛 👻 😵 Stack Frame: Module: 60956632 - [0] _Z	n - 📮						
Solution Explorer - m	natrixMul 🗸 🕂 🗙	matrixMul_kernel.cu			- ×	Nexus CUDA Device Su	mmary		- ×
🗎 🚯 눩 🖧 -		// Loop over all the sub-matrice	s of A	and B		Name		Detaile	*
Solution 'Nexus CUDA Samples.90' (3 projects) // required to compute the block sub			sub-n	atrix	^			Details	
for (int a = aBegin, b = bBegin;									
🔠 🗀 inc		a <= aEnd;				Device 0			
🖃 🗁 src		a += aStep, b += bStep) {				Device 1			
- matri	rixMul.cu					Context 277237	6	Device 0	
		// Declaration of the shared	memor	y array As	used to	Module 6095	6632	c:/ProgramData/NVIDIA Nexus 1.0/Samples/CUI	DA/Debug
WIDIA Nexus -	- CUDA Focus Picker	// store the sub-matrix of A	751 (DT	OCK ST2E1.		🖃 🎟 Grid		_Z9matrixMulPfS_S_ii<<<(8,5),(16,16,1), 0>>>	
		SharedIIOat_AS[bLOCK_SI	CE] [DI	JOCK_512E];		- 🗑 Block 0		Warp Mask: 0x000000FF	
	Dimensions	// Declaration of the shared	memor	w array Be	used to	≣ Warp 0		Active Mask: 0xFFFFFFFF, PC: 0x000703E8, ma	atrixMuL k
Block:	0,0,0 8,5,1 // store the sub-matrix of B			used to	= Wam 1		Active Mask: 0xEEEEEEE PC: 0x000703E8		
		shared float Bs(BLOCK SIZE)(BLOCK SIZE):						Active Made: 0x FFFFFFFF PC: 0x000702E9	
Thread:	read: 0,0,0 16,16,1							Active Mask, 0XFFFFFFFF, FC, 0X000703E6,	
(i) Example	*	<pre>// Load the matrices from de</pre>	vice n	nemory	E			Active Mask: UXFFFFFFF, PC: UXU00703E8,	
#129 for bl	() Examples (// to shared memory; each thread)			loads		≣ Warp 4		Active Mask: 0xFFFFFFF, PC: 0x000703E8,	
10 for coor	dinates 10.0	<pre>// one element of each matri</pre>	ĸ			🗄 Warp 5		Active Mask: 0xFFFFFFF, PC: 0x000703E8,	
10. 5 for co	ordinates 10, 5	AS(ty, tx) = A[a + wA * ty + tx];				🗄 Warp 6		Active Mask: 0xFFFFFFF, PC: 0x000703E8,	
	· · · · · · · · · · · · · · · · · · ·	BS(ty, tx) = B[b + wB * ty +	tx];			🗄 Warp 7		Active Mask: 0xFFFFFFF, PC: 0x000703E8,	
	OK Cancel								
		<pre>// Synchronize to make sure</pre>	the ma	trices are	loaded				
<u></u>		<pre>syncthreads();</pre>			-				
Solution Explorer	Class View	•			۱.				-
Locals			ч×	Memory 1					- -
Name	Value	Туре	^	Address: 0x0	0000024			- {\$> Columns: Auto	0 ·
+ OlockDim	$\{x = 16, y = 16, z = 1\}$	const dim3		0x0000024	9c e9 4d 3e	e0 f1 6f 3e 0c f9 85	3e 82 2	1 41 3e 6c e7 35 3f 7f œéM>àño>.ù.>.!A	1>1ç5?.
	$\{X = 8, y = 5, z = 1\}$	Const dim3		0x0000039	63 3f 3f 97	4d 4b 3f 91 59 48 3f	0a e1 0	4 3e 1f 69 0f 3f fc 11 c??-MK?'YH?.á.>	1.?ü.
- V AS	0x00000024 {{0.20108646, 0.23432112, 0.2	16657, 0.18860439,}, {0.8812524, float[16][16]shared	- 1	0x0000004E	fe 3d 1c d5	0d 3f 5a 3d ad 3e e4	27 72 3	f 8a f9 44 3e ca c7 64 þ=.0.?Z=-≻ä'r?S 2 51 21 22 24 11 52 22 2008-20 2008-20 2006	MD>EÇd
	0x00000024 {0.20108040, 0.23432112, 0.20	10057, 0.1000439,} IIO4([10]Shared	E	0x00000078	eb 43 75 3f	8a ed c4 3e ba dd dc	3e 54 2	5 51 21 30 44 11 52 30	?#u`>I
·····································	0x000000a4 {0.55427718 0.1802118 0.766	0004 (0.6812347) 0.21962402, 0.139/1923, 0.139/3963,; III04[10]_Shared_ 0x000008D 7			75 24 3f c5	b1 62 3f 9a 3b 4d 3f	a9 bf 5	4 3f 88 33 44 3f 47 65 u\$?űb?š;M?©;T?	?^3D?Ge
-+ ¢ [3]	0x000000e4 {0.60716575, 0.673513, 0.2610	0x0000			a3 3e 1c e5	Od 3f 71 89 38 3e 89	57 44 3	f 22 d9 10 3f 13 71 09 £>.å.?q.8>.WD?"	'Ù.?.q.
	0x00000424 {{0.80645162, 0.41080967, 0.12955107, 0.26792198,}, {0.179754C float[16][16] shared			0x00000B7	3e ba c1 dc	3d dc e3 6d 3f 9c c1	cd 3e 7	6 c1 3a 3e 70 c9 37 3f >°ÁÜ=Üãm?œÁÍ>vÁ	1:>pÉ7?
a a	0	int		0x000000CC	21 4d 10 3f	d6 37 6b 3f db 7d 6d	3f d9 8	5 6c 3f 42 3d 21 3f 47 !M.?O7k?U}m?U.1	.?B=!?G
🥥 b	0	int		0x000000F6	5b 3e 10 35	88 3e 21 6f 10 3f 04	2d 82 3	e 89 51 c4 3d 5b 99 ad [>.5^>!o.?>.	.QÄ=[™-
🥥 bx	0	int		0x0000010B	3e 2e 15 97	3e cd ab 66 3f 23 91	11 3f c	6 f1 62 3f 20 e9 0f 3e >>Í≪f?‡`.?Æñ	ib? é.≻
🧼 by	0	int		0x00000120	7f 49 3f 3f	dc 11 ee 3d 84 c9 41	3e da 0	1 6d 3c d8 dd 6b 3f ee .I??Ü.î=.ÉA>Ú.m	n≺ØÝk?î
A tv	0	int		0x00000135	01 f7 3b ee	fb 76 3f 2f 81 17 3d	6e f3 3	6 3f 48 1d a4 3e ef b1 .÷;îûv?/=nó6?	/H.¤>ï±
🔜 Autos 👼 Locals	👼 Threads 🛛 Modules 👼 Watch 1			Memory 1	🖧 Call Stack 🔽	Breakpoints 🖃 Output 📑	Pending	Checkins	

Fermi Summary

Future: ExaScale Computing

- ExaScale (10^18) is the next milestone following
 PetaScale (10^15) computing
- DARPA, HP, France and others are planning for ExaScale computing facilities.
- DARPA report (28th Sept., 2008) described four major challenges (Power consumption, memory, Concurrency/Locality, Resiliency)
- DARPA reports that the number one issue is Power. Extrapolation of Power indicates over 100MW for Exaflop.

ExaScale Computing Study: Technology Challenges in Achieving Exascale Systems

Peter Kogge, Editor & Study Lead Keren Bergman Shekhar Borkar Dan Campbell William Carlson William Dally Monty Denneau **Paul Franzon** William Harrod Kerry Hill Jon Hiller Sherman Kart Stephen Keckler Dean Klein Robert Luca: Mark Richard Al Scarpelli Steven Scott Allan Snavely Thomas Sterling R. Stanley Williams Katherine Yelick

September 28, 2008

This work was sponsored by DARPA IPTO in the ExaStable Computing Study with Dr. William Harod as Program Manager, AFEL contract mumber FA8626-07-C-7224. This report is published in the interest of scientific and inclucing information exchange and its publication form not constitute the Government's approval or elonger end of this idea or findings.

NOTICE

Using Government denvique, specifications, or other data included in this document for any parsyos other than Greenment proceedenses of the one of any any style high the U.S. Government. The fact that the Government formation or supplied the drawings, specifications, or other data does not locate the holder or any other person or corporation, or convey any ights or permission to manufacture, use, or call any parameted arrandom that may relate to them.

APPROVED FOR PUBLIC RELEASE, DISTRIBUTION UNLIMITED.

Available at www.darpa.mil/personnel/docs/ExaScale_Study_Initial.pdf

A Possible NVIDIA ExaScale Machine in 2017

by Bill Dally: Chief Scientist & Sr. VP of Research, NVIDIA A projection based on Moore's Law and does not represent a committed roadmap

- GPU Node (GPU + CPU + memory + supply) ~300W
 - 2,400 throughput-cores (7,200FPUs), 16 CPUs single chip
 - 40TFLOPS (SP), 13TFLOPS (DP)
- Node Memory
 - 128GB DRAM, 2TB/s bandwidth
 - 512GB Phase-change Flash for checkpoint and scratch
- Cabinet ~100kW
 - 384-Nodes 15.7PFLOPS (SP), 50TB DRAM
 - Dragonfly network 1TB/sec node bandwidth
- System ~ 10MW
 - 128 cabinets 2 ExaFLOPS (SP), 6.8PB DRAM
 - Dragonfly network with active optical links

- Performance of Single-threaded processor is saturating
- Performance = Parallelism
- NVIDIA are committed to offer massive parallel GPUs for now and beyond
- CUDA GPU computing abstracts parallel programming
- 'Fermi' is the key milestone in GPU Computing

Thank you for your attention

