

Next Generation GPU Architecture Code-named 'Fermi'

The Soul of a Supercomputer in the Body of a GPU

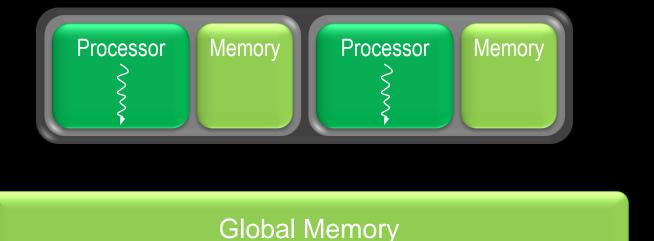
Why is NVIDIA at Super Computing?



- Graphics is a throughput problem
 - paint every pixel within frame time = 120 Mpix / sec
 - every pixel is independent tremendous inherent parallelism
- Create, run, & retire lots of threads very rapidly
 - measured 14.8 Gthread/s on trivial (x[i]+=1) kernel
- Use multithreading to hide latency
 - 1 stalled thread is OK if 100 are ready to run
- Much of scientific computing is also a throughput problem!

Multicore CPU: Run ~10 Threads Fast

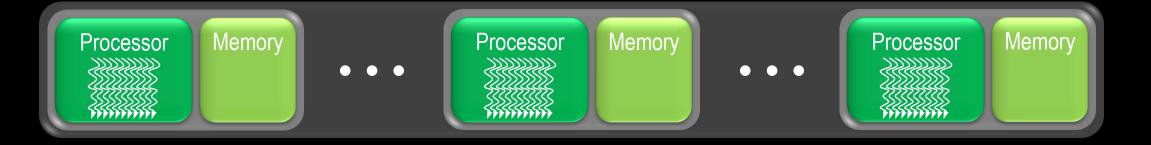




- Few processors, each supporting 1–2 hardware threads
- On-chip memory/cache near processors
- Shared global memory space (external DRAM)

Manycore GPU: Run ~10,000 Threads Fast





Global Memory

- Hundreds of processors, each supporting hundreds of hardware threads
- On-chip memory/cache near processors
- Shared global memory space (external DRAM)

Why do CPUs and GPUs differ?

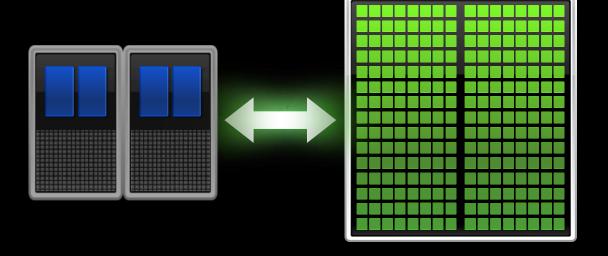


- Different goals produce different designs
 - GPU assumes work load is highly parallel
 - CPU must be good at everything, parallel or not
- CPU: minimize latency experienced by 1 thread
 - lots of big on-chip caches
 - extremely sophisticated control
- GPU: maximize throughput of all threads
 - lots of big ALUs
 - multithreading can hide latency ... so skip the big caches
 - simpler control, cost amortized over ALUs via SIMD

GPU Computing: The Best of Both Worlds



Multicore CPU Manycore GPU

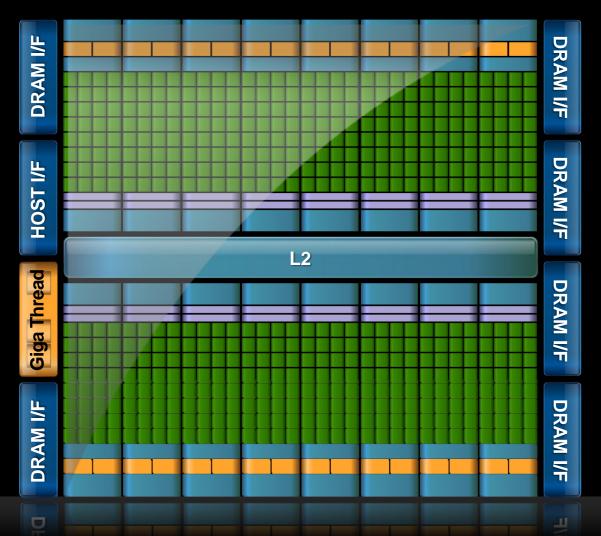


Heterogeneous Computing

Goals of the 'Fermi' Architecture



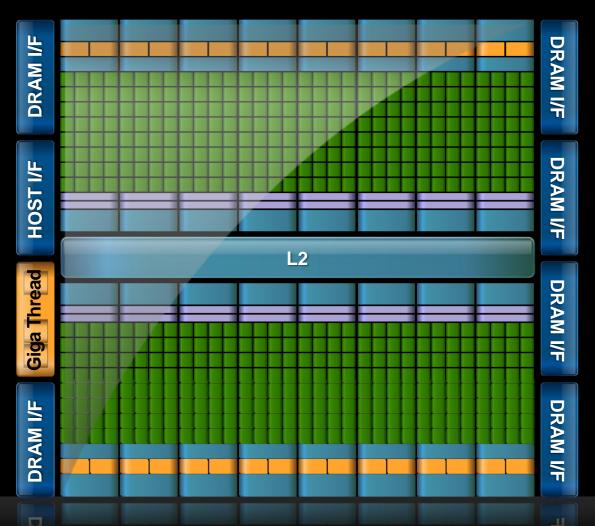
- Improved peak performance
- Improved throughput efficiency
- Broader applicability
- Full integration within modern software development environment



Fermi Focus Areas



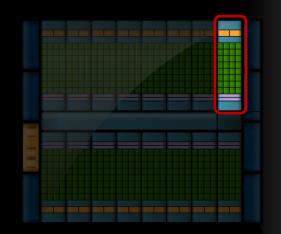
- Expand performance sweet spot of the GPU
 - Caching
 - Concurrent kernels
 - FP64
 - 512 cores
 - GDDR5 memory
- Bring more users, more applications to the GPU
 - C++
 - Visual Studio Integration
 - ECC



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)

	FP32	FP64	INT	SFU	LD/ST
Ops / clk	32	16	32	4	16



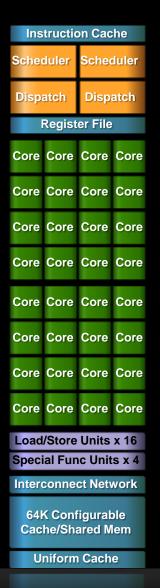


Instruction Cache						
Sche	duler	Scheduler				
Disp	atch	Dispatch				
Register File						
Core	Core	Core	Core			
Core	Core	Core	Core			
Core	Core	Core	Core			
Core	Core	Core	Core			
Core	Core	Core	Core			
Core	Core	Core	Core			
Core	Core	Core	Core			
Core	Core	Core	Core			
Load/Store Units x 16						
Special Func Units x 4						
Interconnect Network						
64K Configurable Cache/Shared Mem						
Uniform Cache						

Fermi Instruction Set Architecture

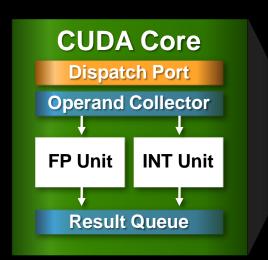
NVIDIA.

- Enables C++ features
 - virtual functions, new/delete, try/catch
- Unified load/store addressing
- 64-bit addressing for large problems
- Optimized for CUDA C, OpenCL & Direct Compute
 - native (x,y)-based LD/ST with format conversion
- Enables system call functionality



CUDA Core Architecture

- New IEEE 754-2008 floating-point standard, surpassing even the most advanced CPUs
- Fused multiply-add (FMA) instruction for both single and double precision
- Newly designed integer ALU optimized for 64-bit and extended precision operations



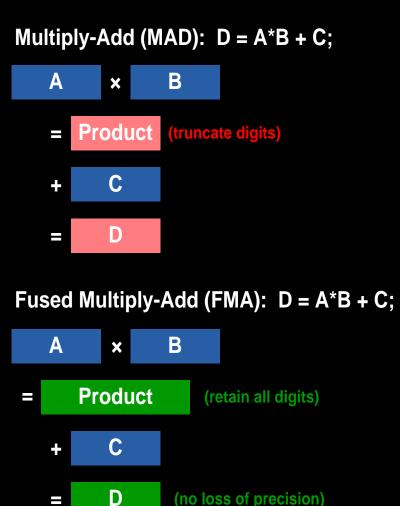




IEEE 754-2008 Floating Point



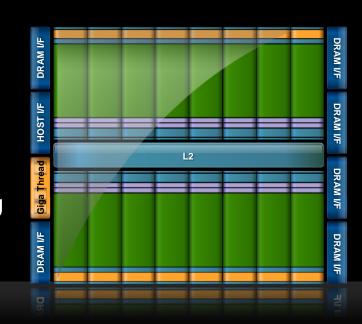
- IEEE 754-2008 results
 - 64-bit double precision
 - 32-bit single precision
 - full-speed subnormal operands & results
 - NaNs, +/- Infinity
- IEEE 754-2008 rounding
 - nearest even, zero, +inf, -inf
- IEEE 754-2008 Fused Multiply-Add (FMA)
 - D = A*B + C;
 - No loss of precision
 - IEEE divide & sqrt use FMA



Cached Memory Hierarchy



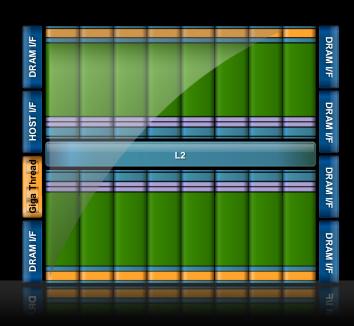
- True cache hierarchy + on-chip shared RAM
 - On-chip shared memory: good fit for regular memory access
 - dense linear algebra, image processing, ...
 - Caches: good fit for irregular or unpredictable memory access
 - ray tracing, sparse matrix multiply, physics ...
- Separate L1 Cache for each SM (16/48 KB)
 - Improves bandwidth and reduces latency
- Unified L2 Cache for all SMs (768 KB)
 - Fast, coherent data sharing across all cores in the GPU



Larger, Faster Memory Interface



- GDDR5 memory interface
 - 2x speed of GDDR3
- Up to 1 Terabyte of memory attached to GPU
 - Operate on large data sets



ECC Memory Protection



- All major internal memories protected by ECC
 - Register file
 - L1 cache
 - L2 cache
- External DRAM protected by ECC
 - GDDR5 as well as SDDR3 memory configurations
- ECC is a must have for many computing applications
 - Clear customer feedback

Faster Atomic Operations (Read / Modify / Write)



- Significantly more efficient chip-wide inter-block synchronization
- Much faster parallel aggregation
 - ray tracing, histogram computation, clustering and pattern recognition, face recognition, speech recognition, BLAS, etc
- Accelerated by cached memory hierarchy
- Fermi increases atomic performance by 5x to 20x

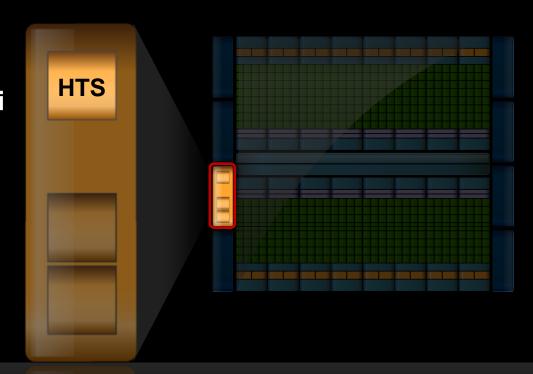
GigaThread[™] Hardware Thread Scheduler



Hierarchically manages tens of thousands of simultaneously active threads

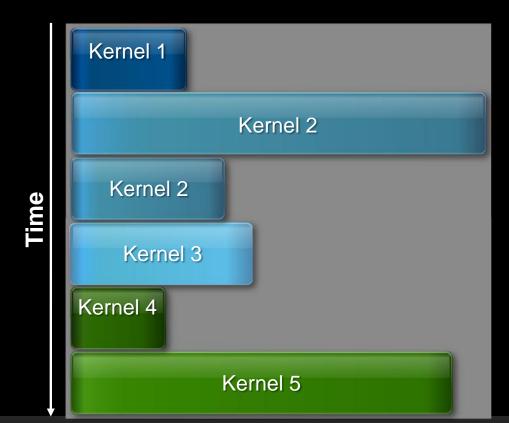
10x faster context switching on Fermi

Concurrent kernel execution



Concurrent Kernel Execution







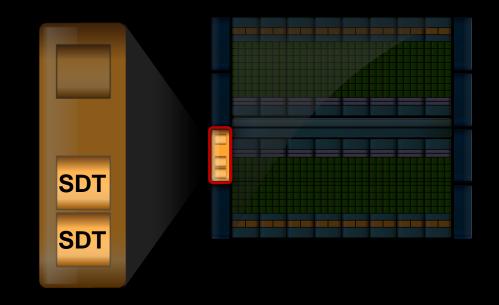
Serial Kernel Execution

Parallel Kernel Execution

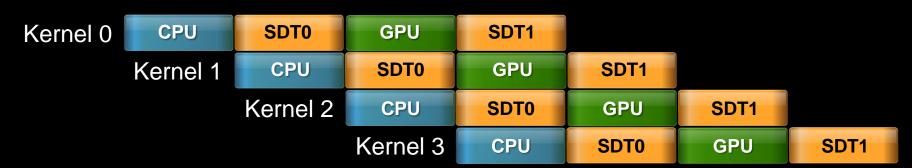
GigaThread Streaming Data Transfer Engine



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



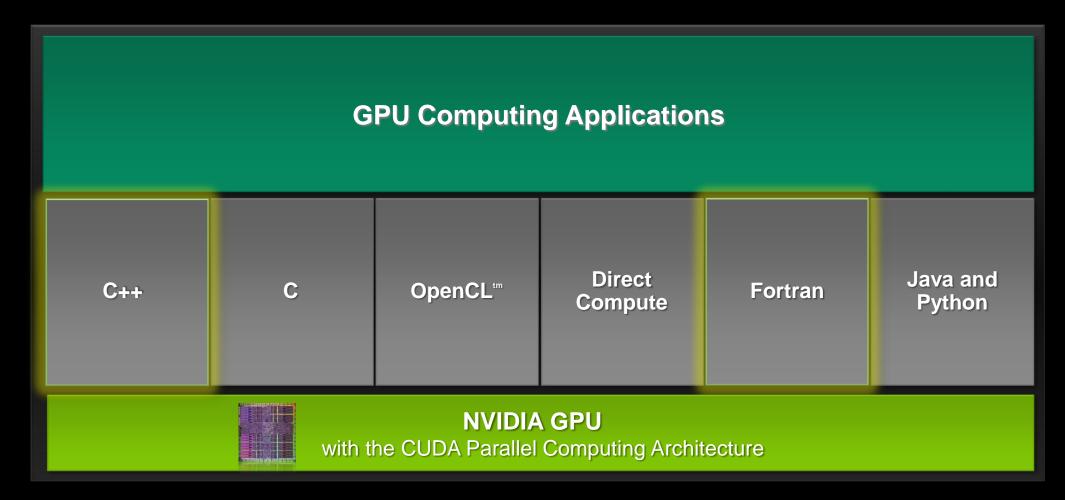
Fully overlapped with CPU/GPU processing



	G80	GT200	Fermi
Transistors	681 million	1.4 billion	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
© 2009 NVIDIA Corporation			

CUDA Parallel Computing Architecture





C with CUDA 3.0



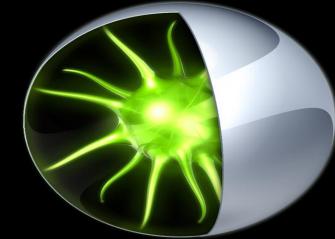
- Unified addressing for C and C++ pointers
 - Global, shared, local addresses
 - Enables 3rd party GPU callable libraries, dynamic linking
 - One 40-bit address space for load/store instructions
- Compiling for native 64-bit addressing
- IEEE 754-2008 single & double precision
 - C99 math.h support
- Concurrent Kernels

NVIDIA Integrated Development Environment Code-named 'Nexus'



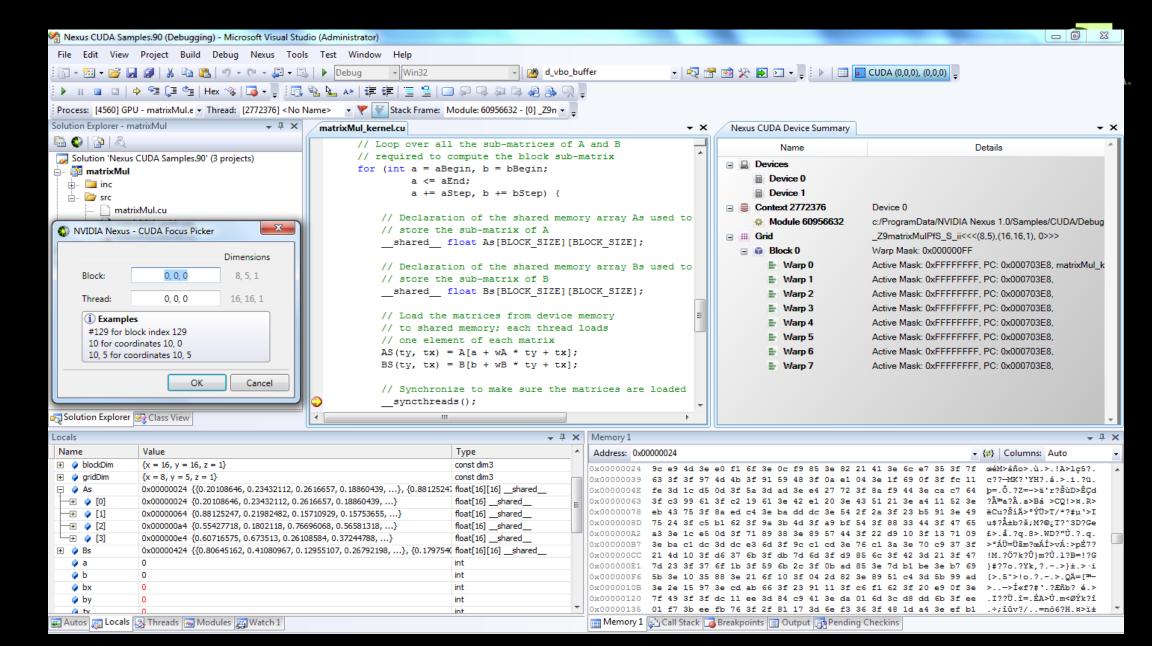
Industry's 1st IDE for massively parallel applications

Accelerate development of CPU + GPU co-processing applications



- Complete Visual Studio-integrated development environment
 - C, C++, OpenCL, & DirectCompute platforms
 - DirectX and OpenGL graphics APIs





Fermi Summary



- Third generation CUDA architecture
- Improved peak performance
- Improved throughput efficiency
- Broader applicability
- Full integration within modern development environment