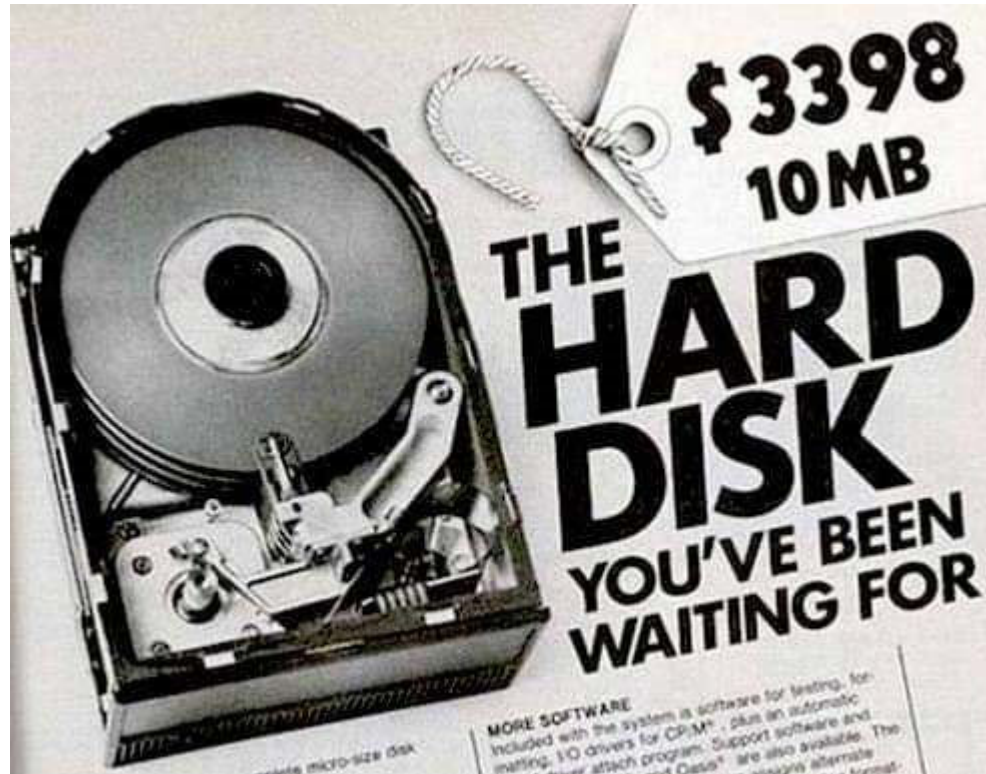


A look at the past...



A look at today...

WD - easystore® 8TB External USB 3.0 Hard Drive - Black

Model: WDBCKA0080HBK-NESN SKU: 5792401

★★★★★ 4.7 (3,225 Customers) | 88 Answered Questions

Only @ Best Buy



\$199.99 or **\$33.34/mo.***
Save \$80 suggested payments with
Was \$279.99 **6-Month Financing**
[Show me how >](#)

Free item with purchase

Geek Protect your product
Learn about Replacement Plans

2 Years
\$19.99

No plan
selected

Add to Cart

Save for Later

☐ Compare

Hard Drive Capacity:

4000GB

8000GB

10000GB

FREE Shipping: Get it by Mon, Dec 3

Want it faster? Pick it up at Salem or see more shipping options in checkout to 97330.

Cardmember Offers

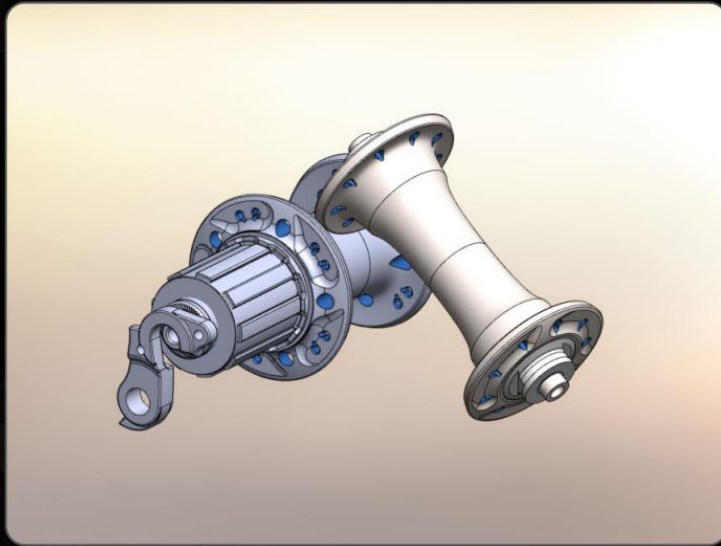
6 Month Financing

Get 5% Back in Rewards

Moving into GPUs...

NIGHT AND DAY DIFFERENCE

Without GPU



With GPU



33 NVIDIA

Graphical Processing Units

- Given the hardware invested to do graphics well, how can be supplement it to improve performance of a wider range of applications?



- Basic idea:
 - Heterogeneous execution model
 - CPU is the *host*, GPU is the *device*
 - Develop a C-like programming language for GPU
 - Unify all forms of GPU parallelism as *CUDA thread*
 - Programming model is “Single Instruction Multiple Thread”

NVIDIA GPU Architecture

- Similarities to vector machines:
 - Works well with data-level parallel problems
 - Scatter-gather transfers
 - Mask registers
 - Large register files
- Differences:
 - No scalar processor
 - Uses multithreading to hide memory latency
 - Has many functional units, as opposed to a few deeply pipelined units like a vector processor

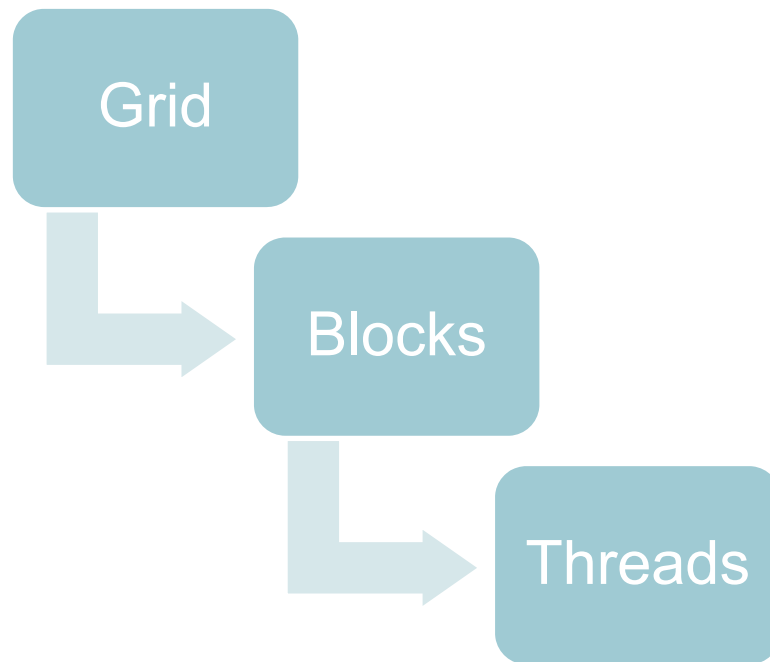
Terminology

- *Threads of SIMD instructions*
 - Each has its own PC
 - Thread scheduler uses scoreboard to dispatch
 - No data dependencies between threads!
 - Keeps track of up to 48 threads of SIMD instructions
 - Hides memory latency
- Thread block scheduler schedules blocks to SIMD processors
- Within each SIMD processor:
 - 32 SIMD lanes
 - Wide and shallow compared to vector processors

Threads and Blocks

- A thread is associated with each data element
 - Threads are organized into blocks
 - Blocks are organized into a grid
-
- GPU hardware handles thread management, not applications or OS

Graphical Explanation



Thread Block 0	SIMD Thread0	A[0] = B [0] * C[0]	
		A[1] = B [1] * C[1]	
		
	SIMD Thread1	A[31] = B [31] * C[31]	
		A[32] = B [32] * C[32]	
		A[33] = B [33] * C[33]	
		
		A[63] = B [63] * C[63]	
		A[64] = B [64] * C[64]	
		
		A[479] = B [479] * C[479]	
		SIMD Thread15	A[480] = B [480] * C[480]
	A[481] = B [481] * C[481]		
		
		A[511] = B [511] * C[511]	
A[512] = B [512] * C[512]			
Grid		A[7679] = B [7679] * C[7679]	
	SIMD Thread0	A[7680] = B [7680] * C[7680]	
		A[7681] = B [7681] * C[7681]	
... ..			
SIMD Thread1	A[7711] = B [7711] * C[7711]		
	A[7712] = B [7712] * C[7712]		
	A[7713] = B [7713] * C[7713]		
		
	A[7743] = B [7743] * C[7743]		
	A[7744] = B [7744] * C[7744]		
		
	A[8159] = B [8159] * C[8159]		
	SIMD Thread15	A[8160] = B [8160] * C[8160]	
A[8161] = B [8161] * C[8161]			
... ..			
		A[8191] = B [8191] * C[8191]	

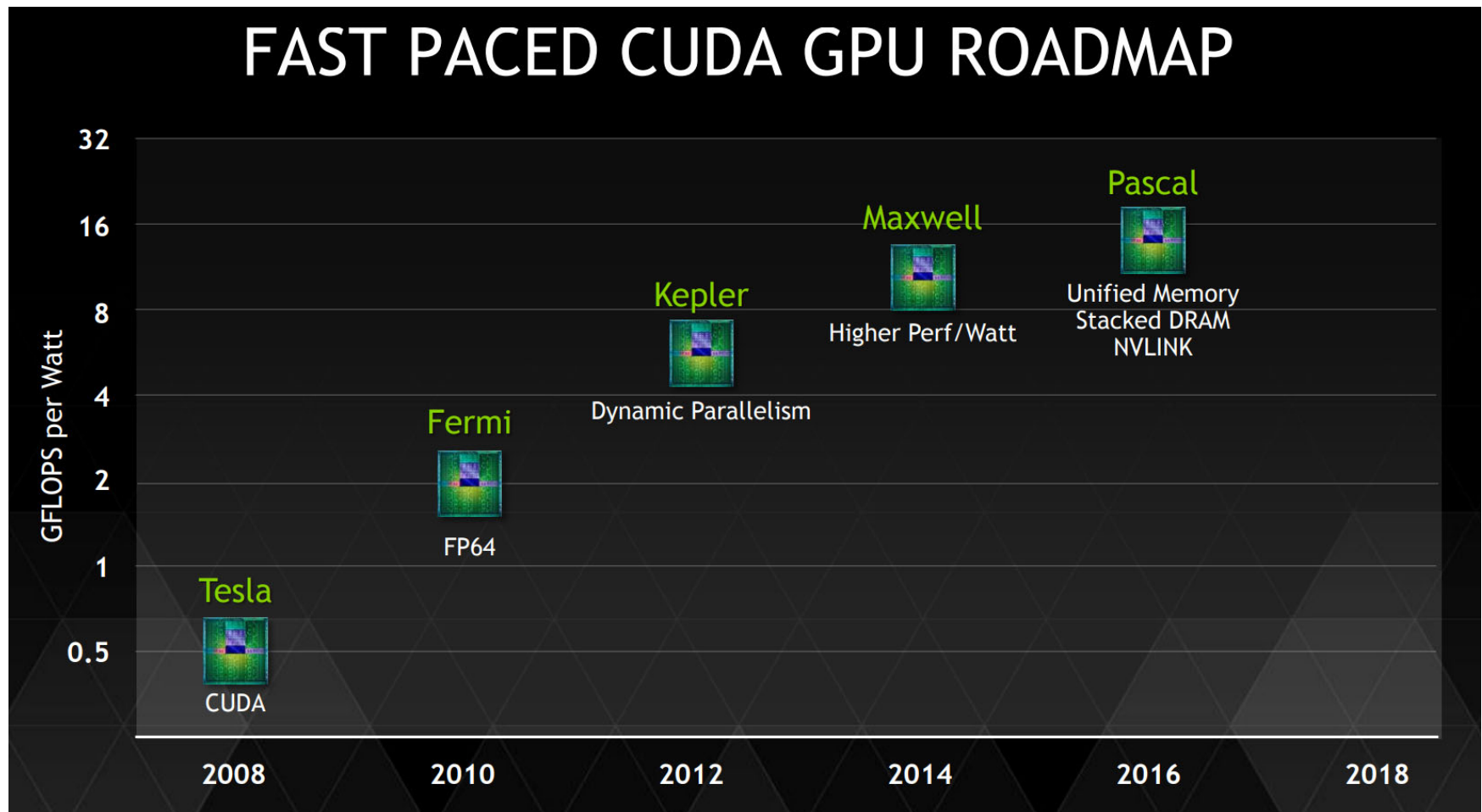
Example

- Multiply two vectors of length 8192
 - Code that works over all elements is the grid
 - Thread blocks break this down into manageable sizes
 - Up to 512 elements per block
 - SIMD instruction executes 32 elements at a time
 - Thus grid size = 16 blocks
 - Block is analogous to a strip-mined vector loop with vector length of 32
 - Block is assigned to a *multithreaded SIMD processor* by the *thread block scheduler*
 - Fermi GPUs have 7-15 multithreaded SIMD processors

Example

- NVIDIA GPU has 32,768 registers
 - Divided into lanes
 - Each SIMD thread is limited to 64 registers
 - SIMD thread has up to:
 - 64 vector registers of 32 32-bit elements
 - 32 vector registers of 32 64-bit elements
- Fermi has 16 physical SIMD lanes, each containing 2048 registers

GPUs Continue to Advance



NVIDIA Instruction Set Arch.

- ISA is an abstraction of the hardware instruction set
 - “Parallel Thread Execution (PTX)”
 - Uses virtual registers
 - Translation to machine code is performed in software
 - Example:

```

shl.s32    R8, blockldx, 9    ; Thread Block ID * Block size (512 or 29)
add.s32    R8, R8, threadldx ; R8 = i = my CUDA thread ID
ld.global.f64 RD0, [X+R8]    ; RD0 = X[i]
ld.global.f64 RD2, [Y+R8]    ; RD2 = Y[i]
mul.f64    RD0, RD0, RD4      ; Product in RD0 = RD0 * RD4 (scalar a)
add.f64    RD0, RD0, RD2      ; Sum in RD0 = RD0 + RD2 (Y[i])
st.global.f64 [Y+R8], RD0    ; Y[i] = sum (X[i]*a + Y[i])
  
```

Conditional Branching

- Like vector architectures, GPU branch hardware uses internal masks
- Also uses
 - Branch synchronization stack
 - Entries consist of masks for each SIMD lane
 - I.e. which threads commit their results (all threads execute)
 - Instruction markers to manage when a branch diverges into multiple execution paths
 - Push on divergent branch
 - ...and when paths converge
 - Act as barriers
 - Pops stack
- Per-thread-lane 1-bit predicate register, specified by programmer

Example

if (X[i] != 0)

 X[i] = X[i] - Y[i];

else X[i] = Z[i];

ld.global.f64 RD0, [X+R8]
setp.neq.s32 P1, RD0, #0
@!P1, bra ELSE1, *Push

; RD0 = X[i]
; P1 is predicate register 1
; Push old mask, set new mask bits
; if P1 false, go to ELSE1

1 ld.global.f64 RD2, [Y+R8]
2 sub.f64 RD0, RD0, RD2
3 st.global.f64 [X+R8], RD0
→ @P1, bra ENDIF1, *Comp

; RD2 = Y[i]
; Difference in RD0
; X[i] = RD0
; complement mask bits

X[i] - Y[i]

ELSE1: ld.global.f64 RD0, [Z+R8]
 st.global.f64 [X+R8], RD0

; if P1 true, go to ENDIF1
; RD0 = Z[i]
; X[i] = RD0

= Z[i]

ENDIF1: <next instruction>, *Pop ; pop to restore old mask

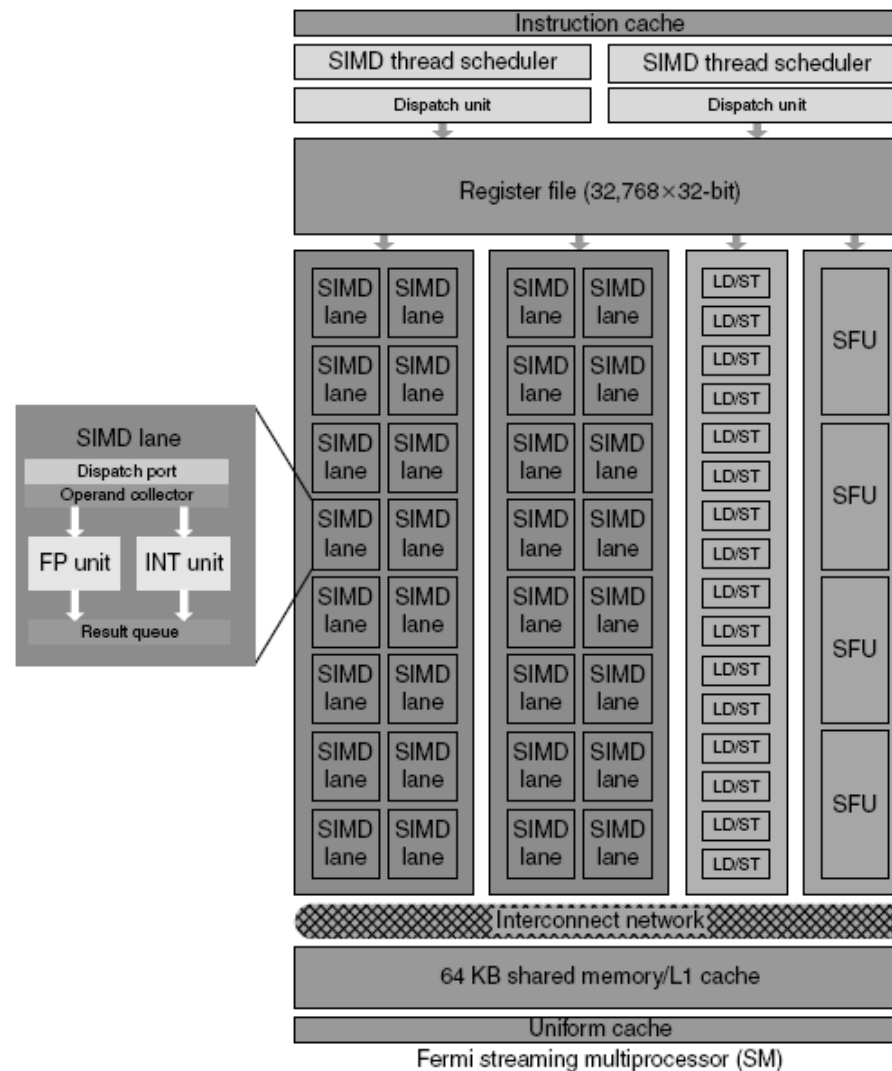
NVIDIA GPU Memory Structures

- Each SIMD Lane has private section of off-chip DRAM
 - “Private memory”
 - Contains stack frame, spilling registers, and private variables
- Each multithreaded SIMD processor also has local memory
 - Shared by SIMD lanes / threads within a block
- Memory shared by SIMD processors is GPU Memory
 - Host can read and write GPU memory

Fermi Architecture Innovations

- Each SIMD processor has
 - Two SIMD thread schedulers, two instruction dispatch units
 - 16 SIMD lanes (SIMD width=32, chime=2 cycles), 16 load-store units, 4 special function units
 - Thus, two threads of SIMD instructions are scheduled every two clock cycles
- Fast double precision
- Caches for GPU memory
- 64-bit addressing and unified address space
- Error correcting codes
- Faster context switching
- Faster atomic instructions

Fermi Multithreaded SIMD Proc.



Loop-Level Parallelism

- Focuses on determining whether data accesses in later iterations are dependent on data values produced in earlier iterations
 - Loop-carried dependence
- Example 1:
for (i=999; i>=0; i=i-1)
 x[i] = x[i] + s;
- No loop-carried dependence

Loop-Level Parallelism

- Example 2:

```
for (i=0; i<100; i=i+1) {  
    A[i+1] = A[i] + C[i]; /* S1 */  
    B[i+1] = B[i] + A[i+1]; /* S2 */  
}
```

- S1 and S2 use values computed by S1 in previous iteration
- S2 uses value computed by S1 in same iteration

Loop-Level Parallelism

- Example 3:

```
for (i=0; i<100; i=i+1) {
    A[i] = A[i] + B[i]; /* S1 */
    B[i+1] = C[i] + D[i]; /* S2 */
}
```

- S1 uses value computed by S2 in previous iteration but dependence is not circular so loop is parallel

- Transform to:

```
A[0] = A[0] + B[0];
for (i=0; i<99; i=i+1) {
    B[i+1] = C[i] + D[i];
    A[i+1] = A[i+1] + B[i+1];
}
B[100] = C[99] + D[99];
```