

A Few Words About Final Project	<u>Example Projec</u>
 Purpose: A chance to dig in deeper into a parallel programming model and explore concepts. Present results to work on communication of technical ideas Write a non-trivial parallel program that combines two parallel programming languages/models. In some cases, just do two separate implementations. OpenMP + SSE-3 OpenMP + CUDA (but need to do this in separate parts of the code) TBB + SSE-3 MPI + OpenMP MPI + SSE-3 MPI + CUDA 	 Look in the textbox Recall Red/Blue fr Implement in MP Implement main Algorithms from C SOR from Ch. 7 CUDA implement FFT from Ch. 10 Jacobi from Ch. 10 Graph algorithms Image and signal p Other domains
Present results in a poster session on the last day of class CS4961	CS







- David Kirk and Wen-mei Hwu manuscript (in progress)
 - http://www.toodoc.com/CUDA-textbook-by-David-Kirkfrom-NVIDIA-and-Prof-Wen-mei-Hwu-pdf.html
- · CUDA 2.x Manual, particularly Chapters 2 and 4 (download from nvidia.com/cudazone)
- Nice series from Dr. Dobbs Journal by Rob Farber - http://www.ddj.com/cpp/207200659

```
11/05/09
```

CUDA Programming Model: A Highly Multithreaded Coprocessor The GPU is viewed as a compute device that: - Is a coprocessor to the CPU or host Has its own DRAM (device memory) Runs many threads in parallel Data-parallel portions of an application are executed on the device as kernels which run in parallel on many threads Differences between GPU and CPU threads

- - GPU threads are extremely lightweight
 - Very little creation overhead - GPU needs 1000s of threads for full efficiency
 - Multi-core CPU needs only a few

11/05/09

UNIVERS

















4







A Very Simple Execution Model

No branch prediction

- Just evaluate branch targets and wait for resolution
- But wait is only a small number of cycles

• No speculation - Only execute useful instructions



UNIVERSITY OF UTAH









7

Cod transformations	Application	
	Conventional Architectures	GPU
Tiling	•Manage reuse in limited storage	•Manage reuse in limited storage •Partition parallel execution at 2 levels
Data-copy	•Eliminate conflict misses in cache	•Copy data to specialized memory structures
Permutation	•Reorder loop structure to enable other optimizations	•Reorder loop structure to enable other optimizations
Unrolling	•Exposes fine-grain parallelism by replicating the loop body	•Exposes fine-grain parallelism by replicating the loop body
•••		







Bandwidth to Shared Memory: Par<u>allel Memory Accesses</u>

- $\boldsymbol{\cdot}$ Consider each thread accessing a different location in shared memory
- Bandwidth maximized if each one is able to proceed *in parallel*
- Hardware to support this
 Banked memory: each bank can support an access on every memory cycle









Shared memory bank conflicts

- Shared memory is as fast as registers if there are no bank conflicts
- The fast case:
 - If all threads of a half-warp access different banks, there is no bank conflict
 - If all threads of a half-warp access the identical address, there is no bank conflict (broadcast)
- The slow case:
 - Bank Conflict: multiple threads in the same half-warp access the same bank

UNIVERSIT

- Must serialize the accesses
- Cost = max # of simultaneous accesses to a single bank

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009 ECE 498AL, University of Illinois, Urbana-Champaign Global Memory Accesses
each thread issues memory accesses to data types of arying sizes, perhaps as small as 1 byte entities
Given an address to load or store, memory returns/ uptes "segments" of either 32 bytes, 64 bytes or 128 bytes
Maximizing bandwidth:
Operate on an *entire* 128 byte segment for each memory transfer







