



## Parts of Exam

## I. Definitions

- A list of 10 terms you will be asked to define
- II. Constraints and Architecture
  - Understand constraints on numbers of threads, blocks, warps, size of storage Understand basic GPU architecture: processors and memory
  - hierarchy
- III. Problem Solving

  - Analyze data dependences and data reuse in code and use this to guide CUDA parallelization and memory hierarchy mapping Given some CUDA code, indicate whether global memory accesses will be coalesced and whether there will be bank conflicts in shared memory

  - shared memory Given some CUDA code, add synchronization to derive a correct implementation Given some CUDA code, provide an optimized version that will have fewer divergent branches Given some CUDA code, derive a partitioning into threads and blocks that does not exceed various hardware limits

L15: Review for Midterm

- IV. (Brief) Essay Question
  - Pick one from a set of 4

UNIVERSITY

## How Much? How Many?

- How many threads per block? Max 512
- How many blocks per grid? Max 65535
- How many threads per warp? 32 •
- How many warps per multiprocessor? 24 •
- How much shared memory per streaming ٠ multiprocessor? 16Kbytes
- How many registers per streaming multiprocessor? 8192 (G80)/16384 (GTX/ Tesla')
- Size of constant cache: 8Kbytes

UNIVERSITY

## Syllabus Syllabus L6 & L7: Memory Hierarchy III: Memory Bondwidth Optimization Leftover from L5 -- Tiling (for registers) Bandwidth - maximize utility of each memory cycle Memory accesses in scheduling (half-ware) Understanding shared memory coalescing (for compute capability < 1.2 and > 1.2) Understanding shared memory bank conflicts L8: Control Flow Divergent branches Execution model Warp vote functions L9: Flacting Point Single precision versus double precision L1: Introduction and CUDA Overview Not much there... L2: Hardware Execution Model Difference between a parallel programming model and a hardware execution model SIMD, MIMD, SIMT, SPMD How are warps selected for execution of a warp? How are warps selected for execution (scorebaarding)? L3: Whiting Correct Programs Race condition, dependence What is a reduction computation and why is it a good match for a GPU? What does \_\_syncthreads () do? (barrier synchronization) Atomic operations Memory Fence Instructions Device emulation mode L4 & L5: Memory Hierarchy: Locality and Data Placement Memory latency and memory bandwidth optimizations Rease callity What are the different memory spaces on the device, who can read/write them? How do you tell the compiler that something belongs in a particular memory space? Tiling transformation (to fit data into constrained storage): Safety and profitability L1: Introduction and CUDA Overview Warp vote functions LiP Floating Point Single precision version Single precision version which operations are compliant? Instrinsics vs. arithmetic operations, what is more precise? What operations can be performed in 4 cycles, and what operations take longer? LiO & LII: Dense Linear Algebra on 6PUs What operatines and be performed in 4 cycles, and what operations take longer? LiO & LII: Dense Linear Algebra on 6PUs What operatines and be performed in 4 cycles, and what operations take longer? LiO & LII: Dense Linear Algebra on 6PUs Transpose in shared memory plus padding trick LI2: Sparse Linear Algebra on 6PUS Different sparse matrix representations Stencil computations using sparse matrix representations Host tiling for constant cache (plus data structure reorganization) Replacing trig function intrinsic calls with hardware implementations Global synchronization for MPM/GLMP

UNIVERSITY

L15: Review for Midterm

UNIVERSITY

2