# 6.823 SPring15 Quiz 4 Review (L20-L24)

TA: Po-An Tsai, Hsin-Jung Yang

# **Transactional memory**

• Atomicity (all or nothing)

At commit, all memory writes take effect at onceOn abort, none of the writes appear to take effect

Isolation

-No other code can observe writes before commit

• Serializability

Transactions seem to commit in a single serial orderThe exact order is not guaranteed

# **Data Management Policy**

 Eager versioning (undo-log based) Update memory location directly Maintain undo info in a log Fast commits Slow aborts

2. Lazy versioning (write-buffer based) Buffer data until commit in a write buffer Update actual memory locations at commit Fast aborts Slow commits

### **Conflict Detection**

• 1. Pessimistic detection

Check for conflicts during loads or stores

• 2. Optimistic detection

Detect conflicts when a transaction attempts to commit

#### Reliability

- What happens if a bit flip
  - DUE : Detected Unrecoverable Error
  - SDC : Silent Data Corruption

- What structures are important
  - ACE : Architecturally Correct Execution
  - AVF : Architectural Vulnerability Factor

#### **DUE and SDC**



#### **ACE and AVF**

• Do the questions to make sure you understand how to calculate or argue for them

#### **VLIW: Very Long Instruction Word**



- Software (compiler) packs independent instructions into a large instruction
- Deterministic latency for all operations
  - Latency measured in 'instructions'
- Hardware does not perform dependency checks
  - Software adds explict NOP instructions

#### **Loop Execution**



How many FP ops/cycle?

1 fadd / 8 cycles = 0.125

- How to get denser packing?
  - Loop unrolling
  - Software pipelining

5/12/2015

# **Loop Unrolling**

• Unroll the inner loop to perform M iterations at once

To get more independent instructions

 Need to handle the case where the total iteration is not a multiple of M



# **Software Pipelining**

• Execute different iterations in parallel



Unroll 4 iterations and then do software pipelining: How many FP ops/cycle? 4 fadds / 4 cycles = 1

### **Predicated Execution**

- Eliminate hard to predict branches with predicated execution
- Example: (p1) ADD r1, r2, r3
  - Execute (commit) add if the predicate register p1 is true



# **Vector Machine**



#### The basic structure of a vector architecture

- Single Instruction Multiple Data (SIMD)
  - Data level parallelism (DLP)
  - Single instruction for multiple loop iterations
  - Vector Length Register (VLR)
  - Conditional execution using the vector mask (VM) register

#### **Vector Arithmetic Execution**

- Use deep pipeline (=> fast clock) to execute element operations
- Simplifies control of deep pipeline because elements in vector are independent (=> no hazards!)



V3 <- v1 \* v2

#### **Multiple Lanes**



# **Vector Instruction Parallelism**

• Can overlap the execution of multiple vector instructions



# **Vector Chaining**

 Allow a vector operation to start as soon as the individual elements of its vector source operand become available



#### **Graphical Processing Units (GPU)**

- GPU consists of many multithreaded SIMD processors
- Terminology
  - Thread: a CUDA thread, which is associated with each data element
  - Warp: a traditional tread that contains SIMD instructions executed on a SIMD processor. Each SIMD instruction operates with multiple (CUDA) threads (ex: 32 threads/warp).
  - Programming abstractions:
    - **Thread block:** A group of threads (ex: 512 threads) mapped to a SIMD processor to execute a vectorized loop.
    - **Grid:** A grid has multiple thread blocks that can be mapped to different SIMD processors and execute independently.

#### **Multithreaded SIMD Processor**



#### **Multithreaded SIMD Processor**



Example:

- 16 physical lanes
- 10s of warps with 32 threads per warp
- Warp scheduler issues an SIMD instruction (for a specific warp) when its elements (threads in the warp) are all ready

# **GPU vs Vector Machine**

#### • Similarities:

- Works well with data-level parallel problems
- Mask registers
- Multiple lanes
- Large register files

#### • Differences:

- GPU uses multithreading to hide memory latency
- GPU has many functional units, as opposed to a few deeply pipelined units like a vector processor

#### The end

#### Good luck!! 🙂