# ECE 571 – Advanced Microprocessor-Based Design Lecture 10

Vince Weaver

http://web.eece.maine.edu/~vweaver

vincent.weaver@maine.edu

23 February 2017

# Announcements

- HW#5 due

- HW#6 will be posted

# Oh No, More Caches!

# Cache Associativity

- direct-mapped – an address maps to only one cache line

- fully-associative (content-addressable memory, CAM) – an address can map to any cache line

- set-associative – an address can map to multiple "ways"

- scratchpad – software managed (seen in DSPs and some CPUs)

# Cache Terms

- Line – which row of a cache being accessed

- Blocks – size of data chunk stored by a cache

- Tags – used to indicate high bits of address; used to detect cache hits

- Sets (or ways) – parts of an associative cache

# Replacement Policy

- FIFO

- LRU

- Round-robin

- Random

- Pseudo-LRU

- Spatial

# Load Policy

- Critical Word First – when loading a multiple-byte line, bring in the bytes of interest first

# Consistency

Need to make sure Memory eventually matches what we have in cache.

- write-back – keeps track of dirty blocks, only writes back at eviction time. poor interaction on multi-processor machines
- write-through – easiest for consistency, potentially more bandwidth needed, values written that are discarded
- write-allocate – Usually in conjunction with write-back Load cacheline from memory before writing.

# Inclusiveness

- Inclusive – every item in L1 also in L2
  simple, but wastes cache space (multiple copies)

- Exclusive – item cannot be in multiple levels at a time

# Other Cache Types

- Victim Cache – store last few evicted blocks in case brought back in again, mitigate smaller associativity

- Assist Cache – prefetch into small cache, avoid problem where prefetch kicks out good values

- Trace Cache – store predecoded program traces instead of (or in addition to) instruction cache

# Virtual vs Physical Addressing

Programs operate on Virtual addresses.

- PIPT, PIVT (Physical Index, Physical/Virt Tagged) – easiest but requires TLB lookup to translate in critical path

- VIPT, VIVT (Virtual Index, Physical/Virt Tagged) – No need for TLB lookup, but can have aliasing between processes. Can use page coloring, OS support, or ASID (address space id) to keep things separate

# Cache Miss Types

- Compulsory (Cold) — miss because first time seen

- Capacity — wouldn't have been a miss with larger cache

- Conflict — miss caused by conflict with another address (would not have been miss with fully assoc cache)

- Coherence — miss caused by other processor

# Fixing Compulsory Misses

Prefetching

- Hardware Prefetchers – very good on modern machines. Automatically bring in nearby cachelines.

- Software – loading values before needed
  also special instructions available

- Large-blocksize of caches. A load brings in all nearby values in the rest of the block.

# Fixing Capacity Misses

- Build Bigger Caches

# Fixing Conflict Misses

- More Ways in Cache

- Victim Cache

- Code/Variable Alignment, Cache Conscious Data Placement

# Fixing Coherence Misses

- False Sharing – independent values in a cache line being accessed by multiple cores

# Cache Parameters Example 1

32kB cache ($2^{15}$), direct mapped ($2^0$)
32 Byte linesize ($2^5$), 32-bit address size ($2^{32}$)

offset $= log_2(linesize) = 5$ bits
lines $= log_2((cachesize/\#ways)/linesize) = 1024$ lines
(10 bits)
tag $=$ addresssize - (offset bits + line bits) $= 17$ bits

| tag | line | offset | |
|---|---|---|---|
| 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 | 14 13 12 11 10 9  8  7  6  5 | 4  3  2  1  0 | |

# Cache Parameters Example 2

32kB cache ($2^{15}$), 4-way ($2^2$)
32 Byte linesize ($2^5$), 32-bit address size ($2^{32}$)

offset $= log_2(linesize) = 5$ bits
lines $= log_2((cachesize/\#ways)/linesize) = 256$ lines
(8 bits)
tag $=$ addresssize - (offset bits + line bits) $= 19$ bits

| tag | | line | offset | |
|---|---|---|---|---|
| 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 | 12 11 10 9 | 8 7 6 5 | 4 3 2 1 0 | |

# Cache Example

512 Byte cache, 2-Way Set Associative, with 16 byte lines, LRU replacement.

24-bit tag, 16 lines (4 bits), 4-bit offset.

| tag | line | offset | |
|---|---|---|---|
| 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9  8 | 7  6  5  4 | 3  2  1  0 | |

# Cache Example 1

# Cache Example – Instruction 1

`ldb r1, 0x00000000`

| | Way 0 | | | |
|---|---|---|---|---|
| line | V | D | LRU | Tag |
| 0 | 1 | 0 | 0 | 0000 00 |
| 1 | 0 | | | |
| 2 | 0 | | | |
| 3 | 0 | | | |
| 4 | 0 | | | |
| 5 | 0 | | | |
| . . . | | | | |
| b | 0 | | | |
| c | 0 | | | |
| d | 0 | | | |
| e | 0 | | | |
| f | 0 | | | |

| | Way 1 | | | |
|---|---|---|---|---|
| | V | D | LRU | Tag |
| | 0 | | | |
| | 0 | | | |
| | 0 | | | |
| | 0 | | | |
| | 0 | | | |
| | 0 | | | |
| | | | | |
| | 0 | | | |
| | 0 | | | |
| | 0 | | | |
| | 0 | | | |
| | 0 | | | |

## Miss, Cold

# Cache Example – Instruction 2

`ldb r1, 0x00000001`

| line | | Way 0 | | | | Way 1 | | |
|------|---|---|-----|---------|---|---|-----|-----|
| | V | D | LRU | Tag | V | D | LRU | Tag |
| 0 | 1 | 0 | 0 | 0000 00 | 0 | | | |
| 1 | 0 | | | | 0 | | | |
| 2 | 0 | | | | 0 | | | |
| 3 | 0 | | | | 0 | | | |
| 4 | 0 | | | | 0 | | | |
| 5 | 0 | | | | 0 | | | |
| . . . | | | | | | | | |
| b | 0 | | | | 0 | | | |
| c | 0 | | | | 0 | | | |
| d | 0 | | | | 0 | | | |
| e | 0 | | | | 0 | | | |
| f | 0 | | | | 0 | | | |

## Hit

# Cache Example – Instruction 3

`ldb r1, 0x00000010`

|  | Way 0 | | | | | Way 1 | | | |
|------|---|---|-----|---------|---|---|---|-----|-----|
| line | V | D | LRU | Tag | | V | D | LRU | Tag |
| 0 | 1 | 0 | 0 | 0000 00 | | 0 | | | |
| 1 | 1 | 0 | 0 | 0000 00 | | 0 | | | |
| 2 | 0 | | | | | 0 | | | |
| 3 | 0 | | | | | 0 | | | |
| 4 | 0 | | | | | 0 | | | |
| 5 | 0 | | | | | 0 | | | |
| . . . | | | | | | | | | |
| b | 0 | | | | | 0 | | | |
| c | 0 | | | | | 0 | | | |
| d | 0 | | | | | 0 | | | |
| e | 0 | | | | | 0 | | | |
| f | 0 | | | | | 0 | | | |

## Miss, Cold

# Cache Example – Instruction 4

`ldb r1, 0x80000010`

| line | Way 0 | | | | Way 1 | | | |
|------|---|---|-----|---------|---|---|-----|---------|
|      | V | D | LRU | Tag     | V | D | LRU | Tag     |
| 0    | 1 | 0 | 0   | 0000 00 | 0 |   |     |         |
| 1    | 1 | 0 | 1   | 0000 00 | 1 | 0 | 0   | 8000 00 |
| 2    | 0 |   |     |         | 0 |   |     |         |
| 3    | 0 |   |     |         | 0 |   |     |         |
| 4    | 0 |   |     |         | 0 |   |     |         |
| 5    | 0 |   |     |         | 0 |   |     |         |
| ...  |   |   |     |         |   |   |     |         |
| b    | 0 |   |     |         | 0 |   |     |         |
| c    | 0 |   |     |         | 0 |   |     |         |
| d    | 0 |   |     |         | 0 |   |     |         |
| e    | 0 |   |     |         | 0 |   |     |         |
| f    | 0 |   |     |         | 0 |   |     |         |

## Miss, Cold

# Cache Example – Instruction 5

`ldb r1, 0xC0000010`

|  | Way 0 | | | | | Way 1 | | | |
| line | V | D | LRU | Tag | | V | D | LRU | Tag |
| 0 | 1 | 0 | 0 | 0000 00 | | 0 | | | |
| 1 | 1 | 0 | 0 | C000 00 | | 1 | 0 | 1 | 8000 00 |
| 2 | 0 | | | | | 0 | | | |
| 3 | 0 | | | | | 0 | | | |
| 4 | 0 | | | | | 0 | | | |
| 5 | 0 | | | | | 0 | | | |
| . . . | | | | | | | | | |
| b | 0 | | | | | 0 | | | |
| c | 0 | | | | | 0 | | | |
| d | 0 | | | | | 0 | | | |
| e | 0 | | | | | 0 | | | |
| f | 0 | | | | | 0 | | | |

Miss, Cold (never in cache previously)

# Cache Example – Instruction 6

`ldb r1, 0xC0000002`

| line | Way 0 | | | | | Way 1 | | | |
|---|---|---|---|---|---|---|---|---|---|
| | V | D | LRU | Tag | | V | D | LRU | Tag |
| 0 | 1 | 0 | 1 | 0000 00 | | 1 | 0 | 0 | c000 00 |
| 1 | 1 | 0 | 0 | C000 00 | | 1 | 0 | 1 | 8000 00 |
| 2 | 0 | | | | | 0 | | | |
| 3 | 0 | | | | | 0 | | | |
| 4 | 0 | | | | | 0 | | | |
| 5 | 0 | | | | | 0 | | | |
| ... | | | | | | | | | |
| b | 0 | | | | | 0 | | | |
| c | 0 | | | | | 0 | | | |
| d | 0 | | | | | 0 | | | |
| e | 0 | | | | | 0 | | | |
| f | 0 | | | | | 0 | | | |

## Miss, Cold

# Cache Example – Instruction 7

`ldb r1, 0x00000010`

|  | Way 0 | | | | Way 1 | | | |
|------|---|---|-----|---------|---|---|-----|---------|
| line | V | D | LRU | Tag | V | D | LRU | Tag |
| 0 | 1 | 0 | 1 | 0000 00 | 1 | 0 | 0 | c000 00 |
| 1 | 1 | 0 | 1 | C000 00 | 1 | 0 | 0 | 0000 00 |
| 2 | 0 | | | | 0 | | | |
| 3 | 0 | | | | 0 | | | |
| 4 | 0 | | | | 0 | | | |
| 5 | 0 | | | | 0 | | | |
| . . . | | | | | | | | |
| b | 0 | | | | 0 | | | |
| c | 0 | | | | 0 | | | |
| d | 0 | | | | 0 | | | |
| e | 0 | | | | 0 | | | |
| f | 0 | | | | 0 | | | |

## Miss, Conflict

# Cache Example – Instruction 8

## stb r1, 0x00000005

<table>
<tr><td rowspan="2">line</td><td colspan="4" align="center">Way 0</td><td colspan="4" align="center">Way 1</td></tr>
<tr><td>V</td><td>D</td><td>LRU</td><td>Tag</td><td>V</td><td>D</td><td>LRU</td><td>Tag</td></tr>
<tr><td>0</td><td>1</td><td>1</td><td>0</td><td>0000 00</td><td>1</td><td>0</td><td>1</td><td>c000 00</td></tr>
<tr><td>1</td><td>1</td><td>0</td><td>1</td><td>C000 00</td><td>1</td><td>0</td><td>0</td><td>0000 00</td></tr>
<tr><td>2</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
<tr><td>3</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
<tr><td>4</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
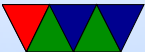<tr><td>5</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
<tr><td>...</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr>
<tr><td>b</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
<tr><td>c</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
<tr><td>d</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
<tr><td>e</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
<tr><td>f</td><td>0</td><td></td><td></td><td></td><td>0</td><td></td><td></td><td></td></tr>
</table>

## Hit