## Slide 1

# Lecture 4
# Basic Caches

**Fall 2016**

**Pejman Lotfi-Kamran**

Adapted from slides originally developed by Profs. Falsafi, Hill, Hoe, Lipasti, Shen, Smith, Sohi, and Vijaykumar of Carnegie Mellon University, EPFL, Purdue University, and University of Wisconsin.

## Slide 2

# Where Are We?

| Fr | Sa | Su | Mo | Tu |
|----|----|----|----|----|
|  | 27-Shahrivar |  | 29-Shahrivar |  |
|  | 3-Mehr |  | 5-Mehr |  |
|  | 10-Mehr |  | 12-Mehr |  |
|  | 17-Mehr |  | 19-Mehr |  |
|  | 24-Mehr |  | 26-Mehr |  |
|  | 1-Aban |  | 3-Aban |  |
|  | 8-Aban |  | 10-Aban |  |
|  | 15-Aban |  | 17-Aban |  |
|  | 22-Aban |  | 24-Aban |  |
|  | 29-Aban |  | 1-Azar |  |
|  | 6-Azar |  | 8-Azar |  |
|  | 13-Azar |  | 15-Azar |  |
|  | 20-Azar |  | 22-Azar |  |
|  | 27-Azar |  | 29-Azar |  |
|  | 4-Dey |  | 6-Dey |  |

This Lecture
  ▫ Basic caches

Next Lecture:
  ▫ Low miss-ratio caches

## Slide 3

# Memory Systems

Basic caches
  ▫ introduction
  ▫ fundamental questions
  ▫ cache size, block size, associativity

today

Advanced caches

Main memory

Virtual memory

## Slide 4

# Motivation

CPU can only go as fast as memory!
  ▫ memory reference/inst x bytes-per-reference x IPC/cycle time
  ▫ In 1990: (1+0.2) x 4 x 1 / 2ns = 2.4 GB/s
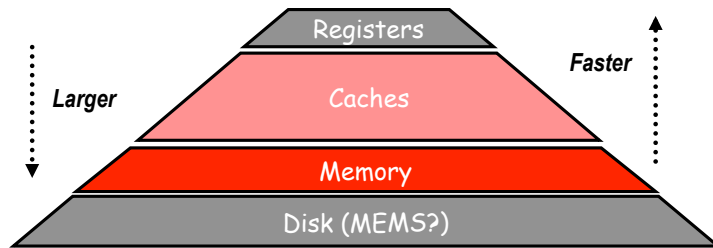  ▫ In 2000: (1x4+0.2x8) x 3 / 0.3ns = 56 GB/s
  ▫

Want storage memory:
  ▫ as fast as CPU
  ▫ as large as required by all of the running applications

## Memory Hierarchy

Make common case fast:

- common: temporal & spatial locality
- fast: smaller more expensive memory



Registers

*Faster*

*Larger*

Caches

Memory

Disk (MEMS?)

---

## Storage Hierarchies

Storages are layered by hierarchies away from the CPU in the order of

- increasing latency $(t_i)$     $t_i < t_{i+1}$
- increasing size $(s_i)$
  - $\Rightarrow$ decrease unit cost $(c_i)$   $s_i < s_{i+1}$, $c_i < c_{i+1}$
- decreasing bandwidth $(b_i)$    $b_i > b_{i+1}$
- increasing xfer unit $(x_i)$     $x_i < x_{i+1}$

Level 0  Registers _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ ISA feature _ _ _ _ _ _ _

                                                           Memory Abstractions

Level 1  (n levels of) Caches
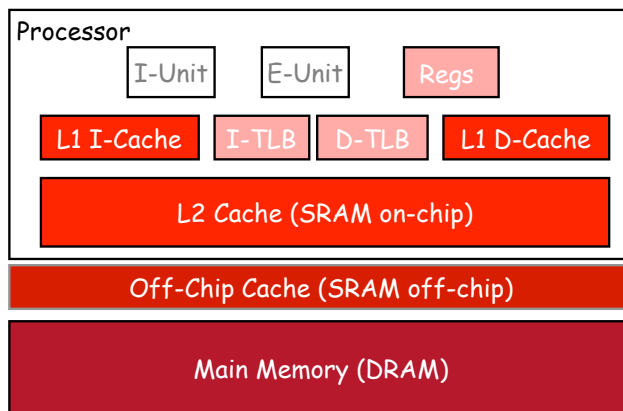
Level 2  Main Memory (Primary Storage)

Level 3  Disks (Secondary Storage)

Level 4  Tape Backup (Tertiary Storage)

---

## Processor/Memory Boundaries

Processor

| I-Unit | E-Unit | Regs |

| L1 I-Cache | I-TLB | D-TLB | L1 D-Cache |

L2 Cache (SRAM on-chip)

Off-Chip Cache (SRAM off-chip)

Main Memory (DRAM)

---

## Caches

An automatically managed hierarchy

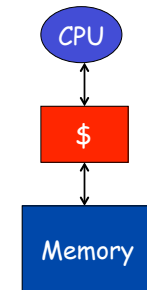> *"A hiding place, esp. of goods, treasure, etc." -- OED*

Keep recently accessed block

- temporal locality

Break memory into blocks (several bytes)

and transfer data to/from cache in blocks

- spatial locality

*A lot of architectures opt for software*

*managed scratch-pad memory instead*

*  e.g. Cray-1, embedded processors, Why??*

CPU

$

Memory

2

## Cache Performance

Assume
- ❑ Cache access time is equal to 1 cycle
- ❑ Cache miss ratio is 0.01
- ❑ Cache miss penalty is 20 cycles

Mean access time

= Cache access time + miss ratio * miss penalty

= 1 + 0.01 * 20 = 1.2

Typically
- ❑ level-1 is 16K-64K, level-2 is 512K-4M,memory is 8G-2TB
- ❑ level-1 as fast as the processor *(increasingly 3-cycles)*
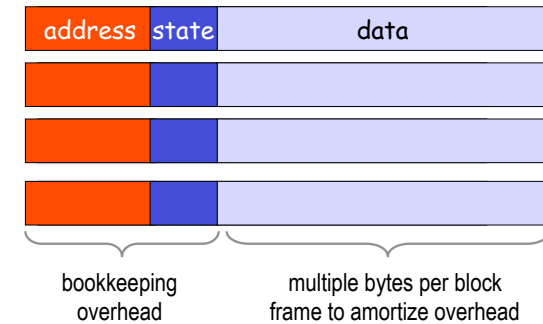- ❑ level-1 is 1/10000 capacity but contains 98% of references

*Memoization & amortization*

## Cache (Abstractly)

Keep recently accessed block in "block frame"
- ❑ state (e.g., valid)
- ❑ address tag
- ❑ data



bookkeeping overhead      multiple bytes per block frame to amortize overhead

## Cache (Abstractly)

On memory read
if incoming address corresponds to one of the stored address tag then
- ○ HIT
- ○ return data

else
- ○ MISS
- ○ choose & displace a current block in use
- ○ fetch new (referenced) block from memory into frame
- ○ return data

- *Where and how to look for a block? (Block placement)*
- *Which block is replaced on a miss? (Block replacement)*
- *What happens on a write? Write strategy (Later)*

## Terminology

block (cache line) — minimum unit that may be present

hit — block is found in the cache

miss — block is not found in the cache

miss ratio — fraction of references that miss

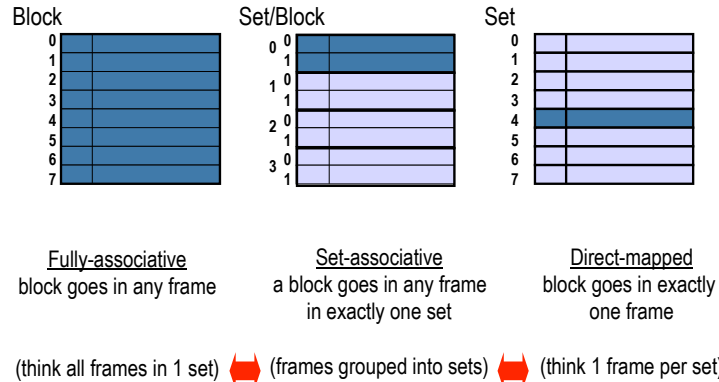hit time — time to access the cache

miss penalty
- ❑ time to replace block in the cache + deliver to upper level
- ❑ access time — time to get first word
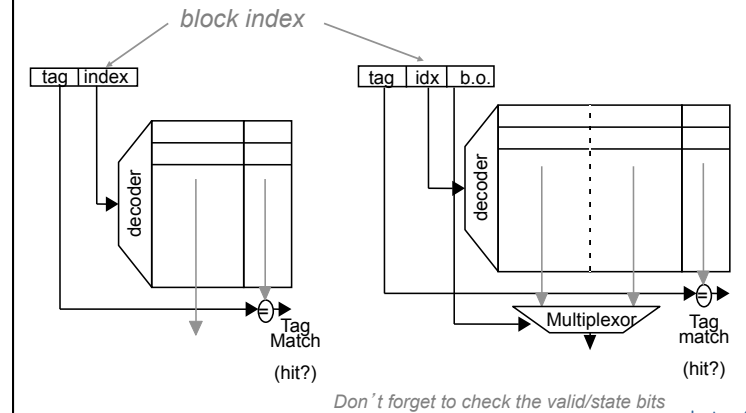- ❑ transfer time — time for remaining words

# Block Placement

Where does block 12 (b' 1100) go?



| Block | Set/Block | Set |

Fully-associative
block goes in any frame

Set-associative
a block goes in any frame
in exactly one set

Direct-mapped
block goes in exactly
one frame

(think all frames in 1 set) ⬌ (frames grouped into sets) ⬌ (think 1 frame per set)

---

# Direct Mapped Caches

*block index*



tag | index

tag | idx | b.o.

decoder

decoder

Tag Match

Multiplexor

Tag match
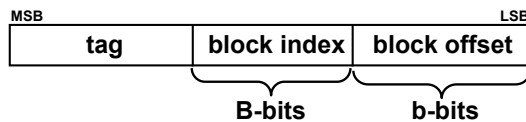
(hit?)

(hit?)

*Don't forget to check the valid/state bits*

---

# Cache Block Size

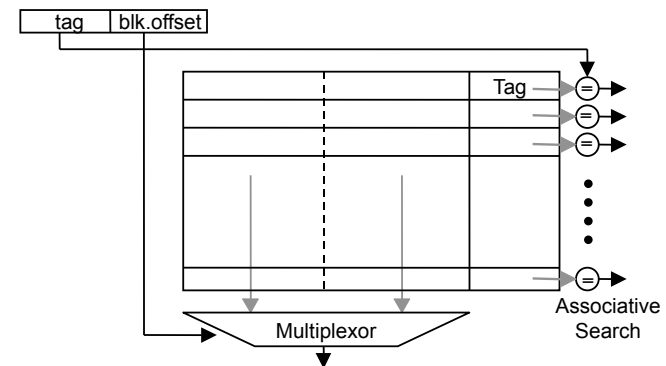Each cache block frame or (cache line) has only one tag but can hold multiple "chunks" of data

  ▫ reduce tag storage overhead

   In 32-bit addressing, an 1-MB direct-mapped cache has 12 bits of tags

   4-byte cache block $\Rightarrow$ 256K blocks $\Rightarrow$ ~384KB of tag
   128-byte cache block $\Rightarrow$ 8K blocks $\Rightarrow$ ~12KB of tag

  ▫ the entire cache block is transferred to and from memory all at once

   *good for spatial locality because if you access address i, you will probably want i+1 as well (prefetching effect)*

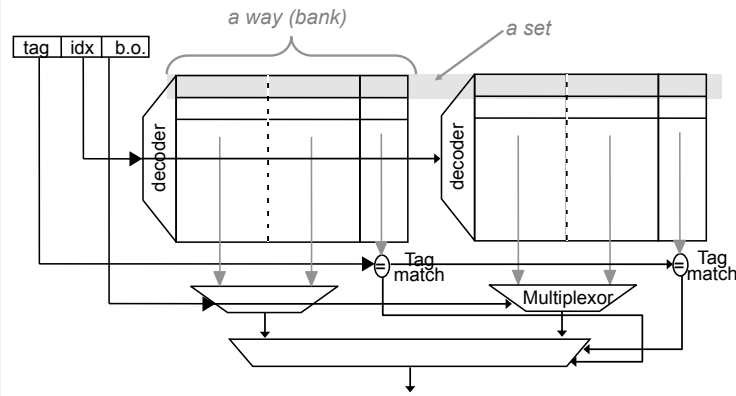Block size = $2^b$; Direct Mapped Cache Size = $2^{B+b}$

| MSB | | | LSB |
|---|---|---|---|
| **tag** | **block index** | **block offset** | |

**B-bits**     **b-bits**

---

# Fully Associative Cache

tag | blk.offset



Tag

Multiplexor

Associative
Search

4

## N-Way Set Associative Cache



*a way (bank)*     *a set*

tag | idx | b.o.

decoder    decoder

Tag match    Tag match

Multiplexor

Cache Size = N x $2^{B+b}$

---

## Associative Block Replacement

Which block in a set to replace on a miss?

Ideally — Belady's algorithm, replace the block that "will" be accessed the furthest in the future

- How do you implement it?

Approximations:

Least recently used — LRU
- optimized (assume) for temporal locality
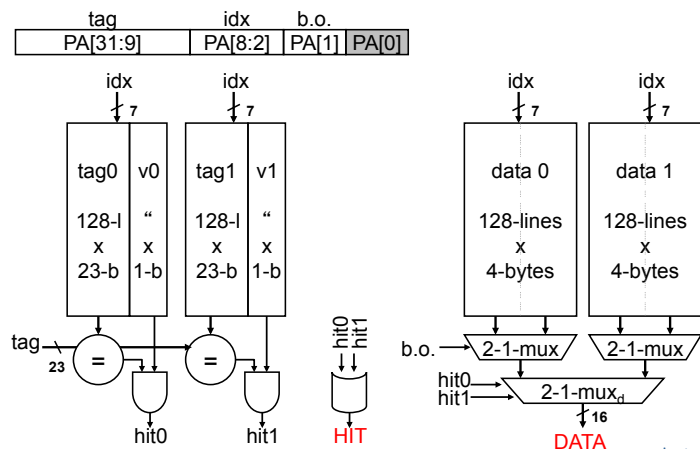  *(expensive for more than 2-way)*

Not most recently used — NMRU
- track MRU, random select from others, good compromise

Random
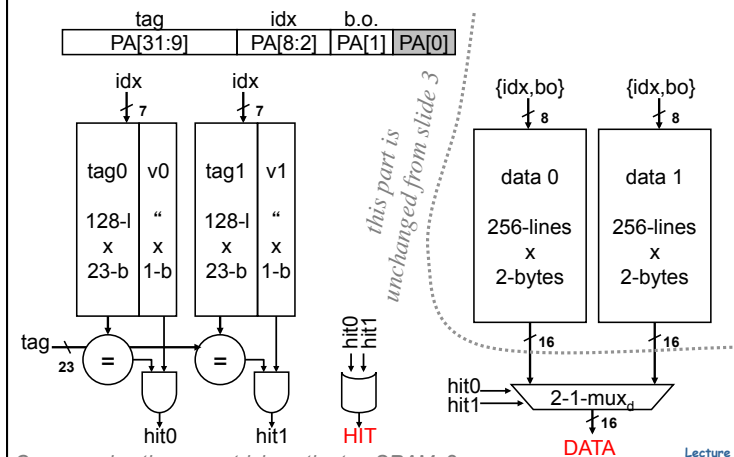- nearly as good as LRU, simpler *(usually pseudo-random)*

---

## Example: a=2, C=1kb, b=4b, word-size=2b
## Basic Solution



tag     idx    b.o.

| PA[31:9] | PA[8:2] | PA[1] | PA[0] |

idx   idx      idx      idx

7   7     7      7

tag0 | v0   tag1 | v1     data 0   data 1

128-l   "   128-l   "    128-lines   128-lines
x   x   x   x    x     x
23-b | 1-b   23-b | 1-b    4-bytes   4-bytes

tag   =   =   hit0 hit1    b.o. → 2-1-mux   2-1-mux

23

hit0   hit1   HIT    hit0 hit1 → 2-1-mux$_d$

16

DATA

---

## The same cache parameters
## but tune for "narrower" data SRAMs



tag     idx    b.o.

| PA[31:9] | PA[8:2] | PA[1] | PA[0] |

idx   idx      {idx,bo}   {idx,bo}

7   7     8      8

*this part is unchanged from slide 3*

tag0 | v0   tag1 | v1     data 0   data 1

128-l   "   128-l   "    256-lines   256-lines
x   x   x   x    x     x
23-b | 1-b   23-b | 1-b    2-bytes   2-bytes

tag   =   =   hit0 hit1    16   16

23

hit0   hit1   HIT    hit0 hit1 → 2-1-mux$_d$

16

DATA

*Can you play the same trick on the tag SRAMs?*

## The same cache parameters but tune for "fatter" data SRAMs

tag — idx — b.o.

| PA[31:9] | PA[8:2] | PA[1] | PA[0] |

PA[8:3]   PA[8:3]

idx   idx   6   6

idx 7   idx 7

tag0 | v0 | tag1 | v1

data 0   data 1

128-l x 23-b | " x 1-b | 128-l x 23-b | " x 1-b

64-lines x 8-bytes   64-lines x 8-bytes

this part is unchanged from slide 3

tag

23   =   =

{PA[2],b.o.} → 4-1-mux   4-1-mux

hit0 hit1

hit0 → 2-1-mux_d

hit1 →   16

hit0   hit1   HIT   DATA

*Can you play the same trick on the tag SRAMs?*

---

## The same cache parameters but each block frame is interleaved over the 2 SRAM banks

tag — idx — b.o.

| PA[31:9] | PA[8:2] | PA[1] | PA[0] |

idx 7   idx 7

idx 7   idx 7

tag0 | v0 | tag1 | v1

data 0   data 1

128-l x 23-b | " x 1-b | 128-l x 23-b | " x 1-b

128-lines x 4-bytes   128-lines x 4-bytes

*this part is unchanged from slide 3*

h0·b̄o   h1·bo   h1·b̄o   h0·bo

tag

23   =   =

h0 h1

b.o. → 2-1-mux   2-1-mux

h0·b̄o+h1·bo
h1·b̄o+h0·bo → 2-1-mux_d   16

h0   h1   HIT   DATA

---

## Miss Classification (3+1 C's)

compulsory
- "cold miss" on first access to a block
  — defined as: miss in infinite cache

capacity
- misses occur because cache not large enough
  — defined as: miss in fully-associative cache

conflict
- misses occur because of restrictive mapping strategy
- only in set-associative or direct-mapped cache
  — defined as: not attributable to compulsory or capacity

coherence
- misses occur because of sharing among multiprocessors

---

## Fundamental Cache Parameters that affects miss rate

Cache size            *(C)*

Block size            *(b)*

Cache associativity   *(a)*

6

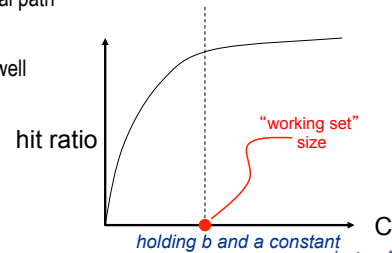## Cache Size

Cache size in the total data (not including tag) capacity
- bigger can exploit temporal locality better
- not ALWAYS better

Too large a cache
- smaller is faster => bigger is slower
- access time may degrade critical path

Too small a cache
- don't exploit temporal locality well
- useful data constantly replaced

hit ratio

"working set" size

C

*holding b and a constant*

## Block Size

Block size is the data that is both
- associated with an address tag
- not necessarily the unit of transfer between hierarchies *(remember sub-blocking)*

Too small blocks
- don't exploit spatial locality well
- have inordinate tag overhead

Too large blocks
- useless data transferred
- useful data permanently replaced
  — too few total # blocks

hit ratio

b

*holding C and a constant*

## Associativity

Partition cache frames into
- equivalence classes of frames called sets

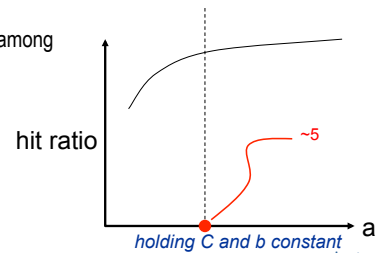Typical values for associativity
- 1, 2-, 4-, 8-way associative
  *5-way associative 20KByte on SuperSPARC, Why?*

Larger associativity
- lower miss ratio, less variation among programs
- only important for small "C/b"

Smaller associativity
- lower cost, faster hit time

hit ratio

~5

a

*holding C and b constant*

## Write Policies

Writes are more interesting
- on reads, data can be accessed in parallel with tag compare
- on writes, needs two steps
- is turn-around time important on for writes?
  *cache optimization often defer writes for reads*

Choices of Write Policies
- On write hits, update memory?
  - Yes: write-through
    *+no coherence issue, +immediate observability, -more bandwidth*
  - No: write-back
- On write misses, allocate a cache block frame?
  - Yes: write-allocate
  - No: no-write-allocate

## Write Policies (Cont.)
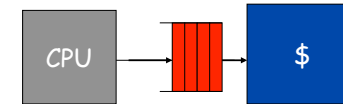
Write-through
- ❑ update memory on each write
- ❑ keeps memory up-to-date
- ❑ traffic/reference = $f_{writes}$, e.g. 0.20
  - independent of cache performance (miss ratio)

Write-back
- ❑ update memory only on block replacement
- ❑ many cache lines are only read and never written to
- ❑ add "dirty" bit to status word
  - ○ originally cleared after replacement
  - ○ set when a block frame is written to
  - ○ only write back a dirty block, and "drop" clean blocks w/o memory update
- ❑ traffic/reference = $f_{dirty}$ x miss x B
  - ○ e.g., traffic/reference = 1/2 x 0.05 x 4 = 0.1

## Write Buffers



Buffer CPU writes
- ❑ allows reads to proceed
- ❑ stall only when full
- ❑ data dependence?
  - ○ What happens on dependent loads/stores?

## Write Buffers (Cont.)

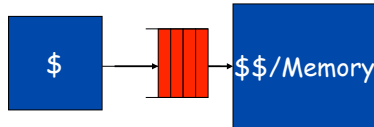| Write Policy | Write Alloc | Hit/Miss | Write Buffer writes to |
|---|---|---|---|
| Back | Yes | Both | Cache |
| Back | No | Hit | Cache |
| Back | No | Miss | Memory |
| Through | Yes | Both | Both |
| Through | No | Hit | Both |
| Through | No | Miss | Memory |

## Write Buffers (Cont.)

Design issues:
- ❑ Design for bursts
- ❑ Coalesce adjacent writes?
- ❑ Sixteen entries is typical

# Writeback Buffers



Between write-back cache and next level
1. Move replaced, dirty blocks to buffer
2. Read new line
3. Move replaced data to memory

Usually need 8 write-back buffer entries

---

# Mark Hill's DM vs. SA: "Bigger & Dumber is Better"

$t_{avg} = t_{hit}$ + miss ratio x $t_{miss}$
- comparable DM and SA caches with same $t_{miss}$
- but, associativity that minimizes $t_{avg}$ is often smaller than associativity that minimizes miss ratio

remember:

$$diff(t_{cache}) = t_{cache}(SA) - t_{cache}(DM) \geq 0$$

$$diff(miss) = miss(SA) - miss(DM) \leq 0$$

e.g.,
assuming $diff(t_{cache})$ = 0 => SA better

assuming $diff$(miss) = -1%, $t_{miss}$ = 20
$\Rightarrow$ if $diff(t_{cache})$ > 0.2 cycle then SA loses

---

# "Harvard" vs. "Princeton"

Unified *(sometimes known as Princeton)*
- less costly, dynamic response, handles writes to instructions

Split I and D *(sometimes known as Harvard)*
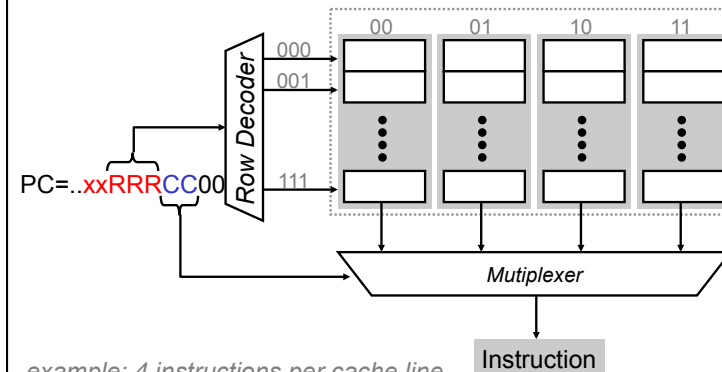- most of the time code and data don't mix
- 2x bandwidth, place close to I/D ports
- can customize size (I-footprint generally smaller than d-footprint), no interference between I/D
- self-modifying code can cause "coherence" problems

Caches should be split for frequent simultaneous I & D access
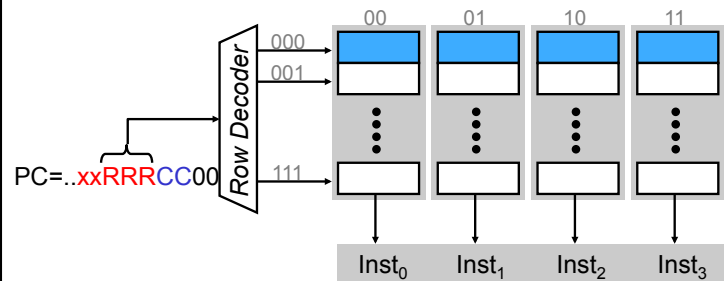- no longer a question in "high-performance" on-chip L-1 caches

---

# L1 Instruction Cache Issues



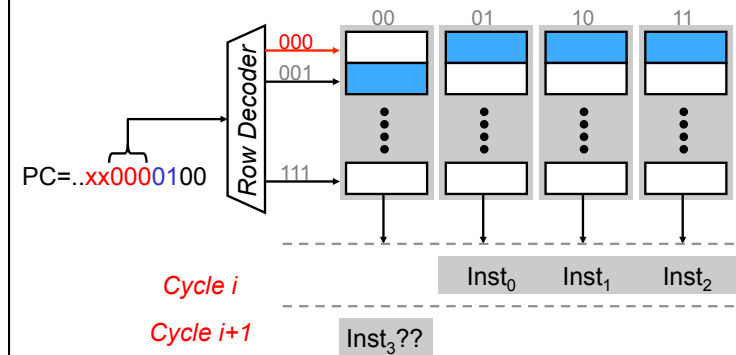PC=..xxRRRCC00

*example: 4 instructions per cache line*

Instruction

9

## Spatial Locality and Fetch Bandwidth



PC=..xxRRRCC00

00    01    10    11

000
001

111

$Inst_0$   $Inst_1$   $Inst_2$   $Inst_3$

Lecture 4
Slide 37

## Fetch Group Misalignment



PC=..xx0000100

00    01    10    11

000
001

111

Cycle i

Cycle i+1

$Inst_0$   $Inst_1$   $Inst_2$

$Inst_3$??

Lecture 4
Slide 38

## Auto Alignment



..xx0000100

00    01    10    11

rotation

$Inst_0$   $Inst_1$   $Inst_2$   $Inst_3$

- A block frame physically spans multiple SRAM lines (**What is a block frame?*)
- Do you always get the maximum number of instructions?

Lecture 4
Slide 39