

# CIS 501

## Computer Architecture

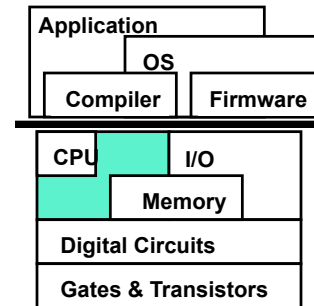
### Unit 4: Caches

Slides originally developed by Amir Roth with contributions by Milo Martin at University of Pennsylvania with sources that included University of Wisconsin slides by Mark Hill, Guri Sohi, Jim Smith, and David Wood.

## Readings

- H+P
  - Appendix C.1–C.3
  - Chapter 5.1–5.2
- Paper:
  - Jouppi, "Improving Direct-Mapped Cache Performance by the Addition of a Small Fully-Associative Cache and Prefetch Buffers", ISCA 1990
    - ISCA's "most influential paper award" awarded 15 years later

## This Unit: Caches



- Basic memory hierarchy concepts
  - Speed vs capacity
- Caches
- Later
  - Organizing an entire memory hierarchy
  - Main memory
  - Virtual memory
- Note: much of this should be review, some new stuff (hopefully)

## Motivation

- Processor can compute only as fast as memory
  - A 3Ghz processor can execute an "add" operation in 0.33ns
  - Today's "Main memory" latency is more than 33ns
  - Naïve implementation: loads/stores can be 100x slower than other operations
- Unobtainable goal:
  - Memory that operates at processor speeds
  - Memory as large as needed for all running programs
  - Memory that is cost effective
- Can't achieve all of these goals at once

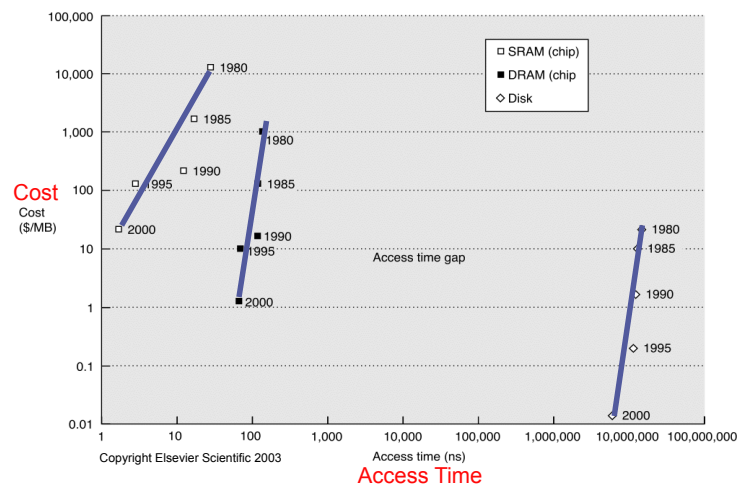
## Review: Types of Memory

- **Static RAM (SRAM)**
  - 6 transistors per bit (two inverters, two other transistors for off/on)
  - Optimized for speed (first) and density (second)
  - Fast (sub-nanosecond latencies for small SRAM)
    - Speed proportional to its area
  - Mixes well with standard processor logic
- **Dynamic RAM (DRAM)**
  - 1 transistor + 1 capacitor per bit
  - Optimized for density (in terms of cost per bit)
  - Slow (>40ns internal access, ~100ns pin-to-pin)
  - Different fabrication steps (does not mix well with logic)
- Nonvolatile storage: Magnetic disk, Flash RAM

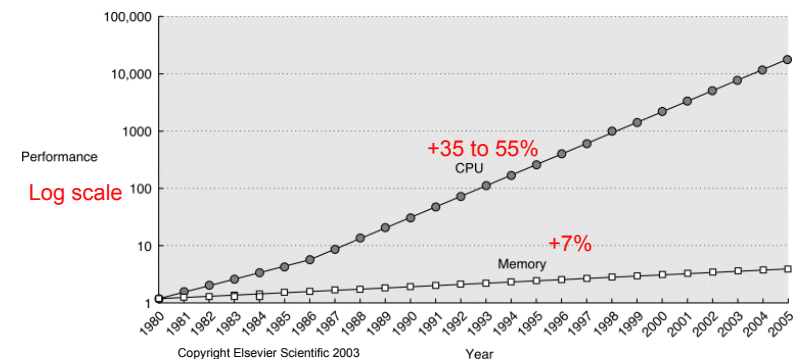
## Memory & Storage Technologies

- **Cost** - what can \$200 buy today?
  - SRAM: 16MB
  - DRAM: 4,000MB (4GB) – 250x cheaper than SRAM
  - Flash: 64,000MB (64GB) – 16x cheaper than DRAM
  - Disk: 2,000,000MB (2TB) – 32x vs. Flash (512x vs. DRAM)
- **Latency**
  - SRAM: <1 to 2ns (on chip)
  - DRAM: ~50ns – 100x or more slower than SRAM
  - Flash: 75,000ns (75 microseconds) – 1500x vs. DRAM
  - Disk: 10,000,000ns (10ms) – 133x vs Flash (200,000x vs DRAM)
- **Bandwidth**
  - SRAM: 300GB/sec (e.g., 12-port 8-byte register file @ 3Ghz)
  - DRAM: ~25GB/s
  - Flash: 0.25GB/s (250MB/s), 100x less than DRAM
  - Disk: 0.1 GB/s (100MB/s), 250x vs DRAM, **sequential** access only

## Memory Technology Trends



## The "Memory Wall"



- Processors are getting faster more quickly than memory (note log scale)
  - Processor speed improvement: 35% to 55%
  - Memory latency improvement: 7%

## Known From the Beginning

“Ideally, one would desire an infinitely large memory capacity such that any particular word would be immediately available ... We are forced to recognize the possibility of constructing a hierarchy of memories, each of which has a greater capacity than the preceding but which is less quickly accessible.”

Burks, Goldstine, VonNeumann

“Preliminary discussion of the logical design of an electronic computing instrument”

IAS memo 1946

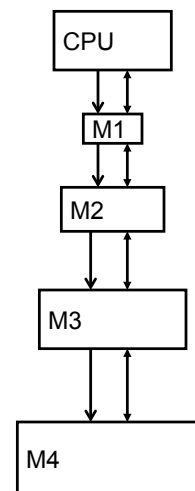
## Library Analogy

- Consider books in a library
- Library has lots of books, but it is slow to access
  - Far away (time to walk to the library)
  - Big (time to walk within the library)
- How can you avoid these latencies?
  - Check out books, take them home with you
    - Put them on desk, on bookshelf, etc.
  - But desks & bookshelves have limited capacity
    - Keep recently used books around (**temporal locality**)
    - Grab books on related topic at the same time (**spatial locality**)
    - Guess what books you'll need in the future (prefetching)

## Locality to the Rescue

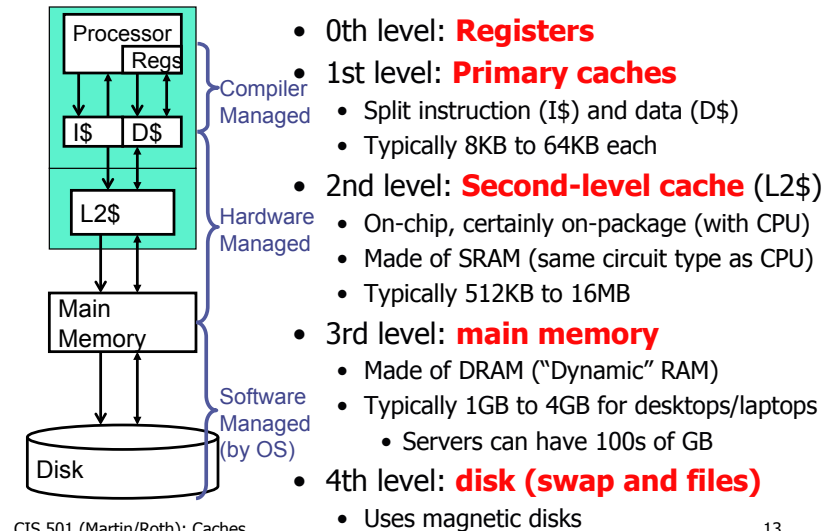
- **Locality of memory references**
  - Property of real programs, few exceptions
  - Books and library analogy (next slide)
- **Temporal locality**
  - Recently referenced data is likely to be referenced again soon
  - **Reactive**: cache recently used data in small, fast memory
- **Spatial locality**
  - More likely to reference data near recently referenced data
  - **Proactive**: fetch data in large chunks to include nearby data
- Holds for data and instructions

## Exploiting Locality: Memory Hierarchy



- Hierarchy of memory components
  - Upper components
    - Fast ↔ Small ↔ Expensive
  - Lower components
    - Slow ↔ Big ↔ Cheap
- Connected by “buses”
  - Which also have latency and bandwidth issues
- Most frequently accessed data in M1
  - M1 + next most frequently accessed in M2, etc.
  - Move data up-down hierarchy
- Optimize average access time
  - $latency_{avg} = latency_{hit} + \%_{miss} * latency_{miss}$
  - Attack each component

## Concrete Memory Hierarchy



CIS 501 (Martin/Roth): Caches

13

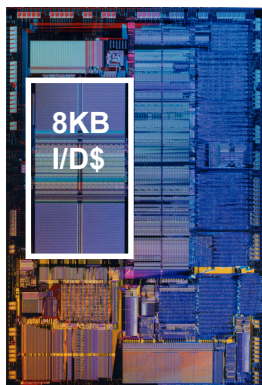
## Library Analogy Revisited

- Registers  $\leftrightarrow$  books on your desk
  - Actively being used, small capacity
- Caches  $\leftrightarrow$  bookshelves
  - Moderate capacity, pretty fast to access
- Main memory  $\leftrightarrow$  library
  - Big, holds almost all data, but slow
- Disk (swap)  $\leftrightarrow$  inter-library loan
  - Very slow, but hopefully really uncommon

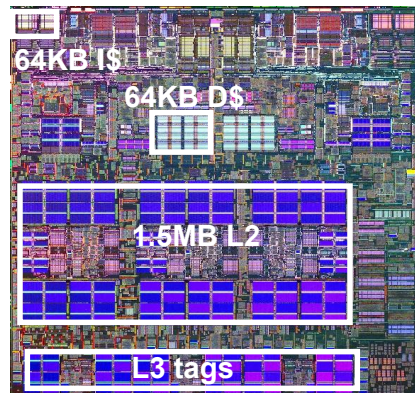
CIS 501 (Martin/Roth): Caches

14

## Evolution of Cache Hierarchies



Intel 486



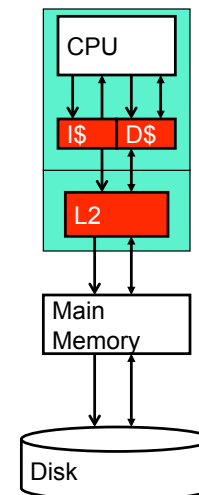
IBM Power5 (dual core)

- Chips today are 30–70% cache by area

CIS 501 (Martin/Roth): Caches

15

## This Unit: Caches

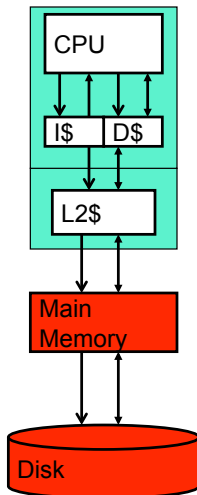


- "Cache": hardware managed
  - Hardware automatically retrieves missing data
  - Built from fast SRAM, usually on-chip today
  - In contrast to off-chip, DRAM "main memory"
- Cache organization
  - ABC
  - Miss classification
- High-performance techniques
  - Reducing misses
  - Improving miss penalty
  - Improving hit latency
- Some example performance calculations

CIS 501 (Martin/Roth): Caches

16

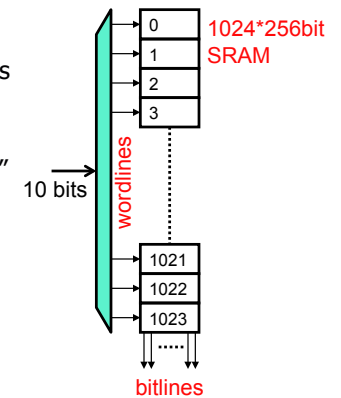
## Looking forward: Memory and Disk



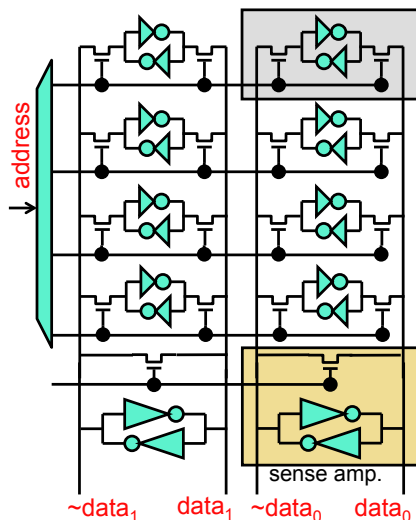
- Main memory
  - DRAM-based memory systems
  - Virtual memory
- Disks and Storage
  - Properties of disks
  - Disk arrays (for performance and reliability)

## Basic Memory Array Structure

- Number of entries
  - $2^n$ , where  $n$  is number of address bits
  - Example: 1024 entries, 10 bit address
  - Decoder changes  $n$ -bit address to  $2^n$  bit "one-hot" signal
  - One-bit address travels on "wordlines"
- Size of entries
  - Width of data accessed
  - Data travels on "bitlines"
  - 256 bits (32 bytes) in example



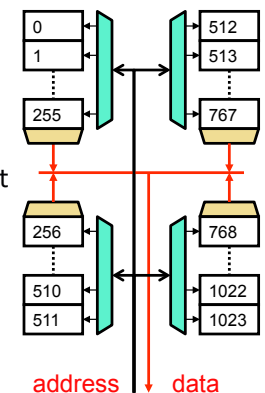
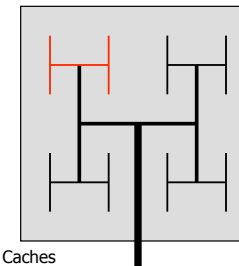
## SRAM Circuit Implementation



- SRAM:
  - Six transistors (6T) cells
  - 4 for the cross-coupled inverters
  - 2 access transistors
- "Static"
  - Cross-coupled inverters hold state
- To read
  - Equalize, swing, amplify
- To write
  - Overwhelm

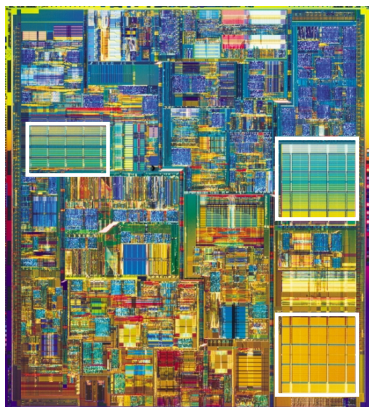
## FYI: Physical Memory Layout

- Logical layout
  - Arrays are vertically contiguous
- Physical layout - roughly square
  - Vertical partitioning to minimize wire lengths
  - **H-tree**: horizontal/vertical partitioning layout
    - Applied recursively
    - Each node looks like an H



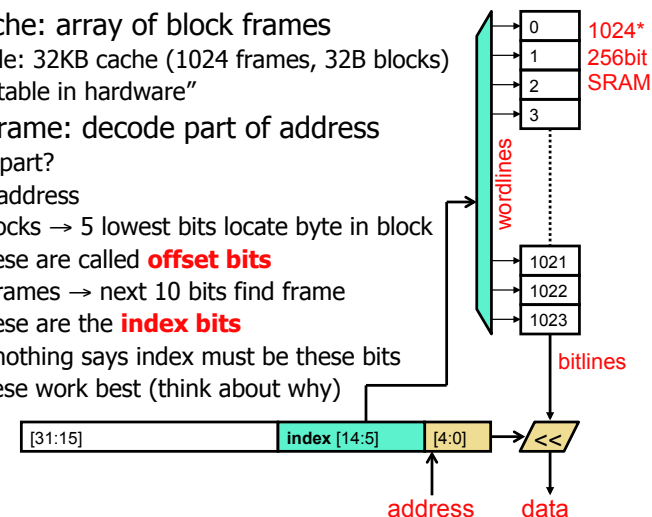
## Physical Cache Layout

- Arrays and h-trees make caches easy to spot in  $\mu$ graphs



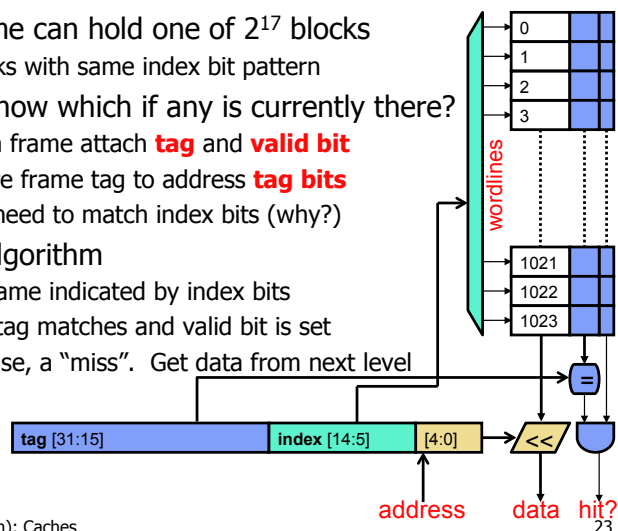
## Caches: Finding Data via Indexing

- Basic cache: array of block frames
  - Example: 32KB cache (1024 frames, 32B blocks)
  - "Hash table in hardware"
- To find frame: decode part of address
  - Which part?
  - 32-bit address
  - 32B blocks  $\rightarrow$  5 lowest bits locate byte in block
    - These are called **offset bits**
  - 1024 frames  $\rightarrow$  next 10 bits find frame
    - These are the **index bits**
  - Note: nothing says index must be these bits
  - But these work best (think about why)



## Knowing that You Found It: Tags

- Each frame can hold one of  $2^{17}$  blocks
  - All blocks with same index bit pattern
- How to know which if any is currently there?
  - To each frame attach **tag** and **valid bit**
  - Compare frame tag to address **tag bits**
    - No need to match index bits (why?)
- Lookup algorithm
  - Read frame indicated by index bits
  - "Hit" if tag matches and valid bit is set
  - Otherwise, a "miss". Get data from next level



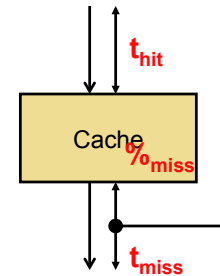
## Calculating Tag Overhead

- "32KB cache" means cache holds 32KB of data
  - Called **capacity**
  - Tag storage is considered overhead
- Tag overhead of 32KB cache with 1024 32B frames
  - 32B frames  $\rightarrow$  5-bit offset
  - 1024 frames  $\rightarrow$  10-bit index
  - 32-bit address  $-$  5-bit offset  $-$  10-bit index = 17-bit tag
  - $(17\text{-bit tag} + 1\text{-bit valid}) \times 1024 \text{ frames} = 18\text{Kb tags} = 2.2\text{KB tags}$
  - $\sim 6\%$  overhead
- What about 64-bit addresses?
  - Tag increases to 49 bits,  $\sim 20\%$  overhead (worst case)

## Handling a Cache Miss

- What if requested data isn't in the cache?
  - How does it get in there?
- Cache controller**: finite state machine
  - Remembers miss address
  - Accesses next level of memory
  - Waits for response
  - Writes data/tag into proper locations
- All of this happens on the **fill path**
- Sometimes called **backside**

## Cache Performance Equation



- For a cache
  - Access**: read or write to cache
  - Hit**: desired data found in cache
  - Miss**: desired data not found in cache
    - Must get from another component
    - No notion of "miss" in register file
  - Fill**: action of placing data into cache
- $\%_{\text{miss}}$  (miss-rate): #misses / #accesses
- $t_{\text{hit}}$ : time to read data from (write data to) cache
- $t_{\text{miss}}$ : time to read data into cache

- Performance metric: average access time

$$t_{\text{avg}} = t_{\text{hit}} + \%_{\text{miss}} * t_{\text{miss}}$$

## CPI Calculation with Cache Misses

- Parameters
  - Simple pipeline with base CPI of 1
  - Instruction mix: 30% loads/stores
  - I\$:  $\%_{\text{miss}} = 2\%$ ,  $t_{\text{miss}} = 10$  cycles
  - D\$:  $\%_{\text{miss}} = 10\%$ ,  $t_{\text{miss}} = 10$  cycles
- What is new CPI?
  - $\text{CPI}_{\text{I\$}} = \%_{\text{missI\$}} * t_{\text{miss}} = 0.02 * 10 \text{ cycles} = 0.2 \text{ cycle}$
  - $\text{CPI}_{\text{D\$}} = \%_{\text{load/store}} * \%_{\text{missD\$}} * t_{\text{missD\$}} = 0.3 * 0.1 * 10 \text{ cycles} = 0.3 \text{ cycle}$
  - $\text{CPI}_{\text{new}} = \text{CPI} + \text{CPI}_{\text{I\$}} + \text{CPI}_{\text{D\$}} = 1 + 0.2 + 0.3 = 1.5$

## Measuring Cache Performance

- Ultimate metric is  $t_{\text{avg}}$ 
  - Cache capacity and circuits roughly determines  $t_{\text{hit}}$
  - Lower-level memory structures determine  $t_{\text{miss}}$
  - Measure  $\%_{\text{miss}}$ 
    - Hardware performance counters
    - Simulation (homework)
    - Paper simulation (next)



## Cache Miss Paper Simulation

- 4-bit addresses → 16B memory
  - Simpler cache diagrams than 32-bits
- 8B cache, 2B blocks
  - Figure out number of sets: 4 (capacity / block-size)
  - Figure out how address splits into offset/index/tag bits
    - Offset: least-significant  $\log_2(\text{block-size}) = \log_2(2) = 1 \rightarrow 0000$
    - Index: next  $\log_2(\text{number-of-sets}) = \log_2(4) = 2 \rightarrow 0000$
    - Tag: rest =  $4 - 1 - 2 = 1 \rightarrow 0000$
- Cache diagram
  - 0000|0001 are addresses of bytes in this block, values don't matter

Cache contents				Address	Outcome
Set00	Set01	Set10	Set11		
0000 0001	0010 0011	0100 0101	0110 0111		

CIS 501 (Martin/Roth): Caches

29

## Cache Miss Paper Simulation

- 8B cache, 2B blocks

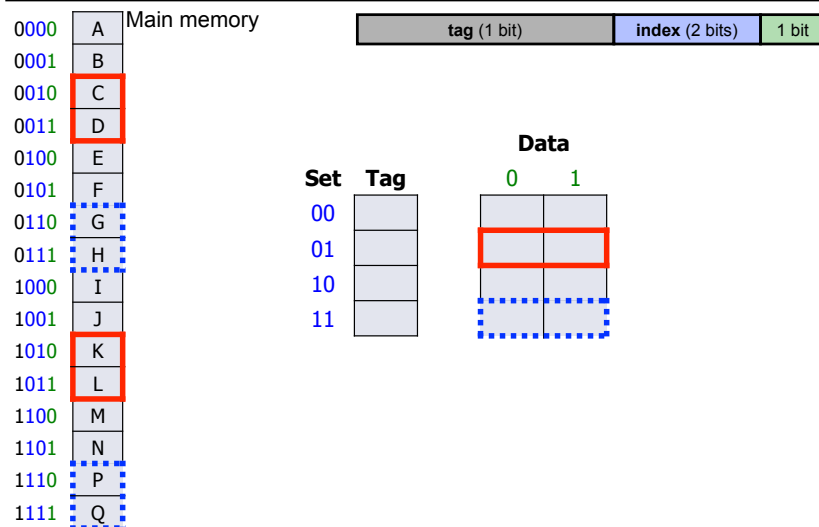
Cache contents (prior to access)				Address	Outcome
Set00	Set01	Set10	Set11		
0000 0001	0010 0011	0100 0101	0110 0111	1100	Miss
0000 0001	0010 0011	1100 1101	0110 0111	1110	Miss
0000 0001	0010 0011	1100 1101	1110 1111	1000	Miss
1000 1001	0010 0011	1100 1101	1110 1111	0011	Hit
1000 1001	0010 0011	1100 1101	1110 1111	1000	Hit
1000 1001	0010 0011	1100 1101	1110 1111	0000	Miss
0000 0001	0010 0011	1100 1101	1110 1111	1000	Miss

- How to reduce %<sub>miss</sub>? And hopefully t<sub>avg</sub>?

CIS 501 (Martin/Roth): Caches

30

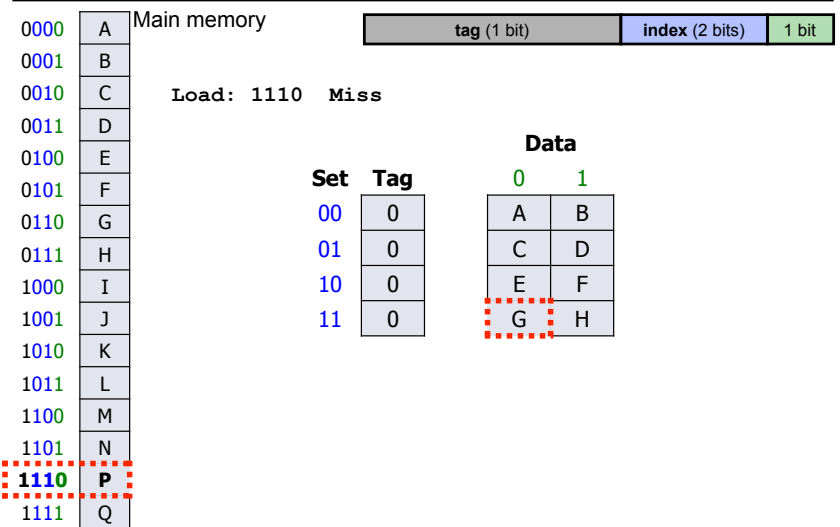
## 4-bit Address, 8B Cache, 2B Blocks



CIS 501 (Martin/Roth): Caches

31

## 4-bit Address, 8B Cache, 2B Blocks

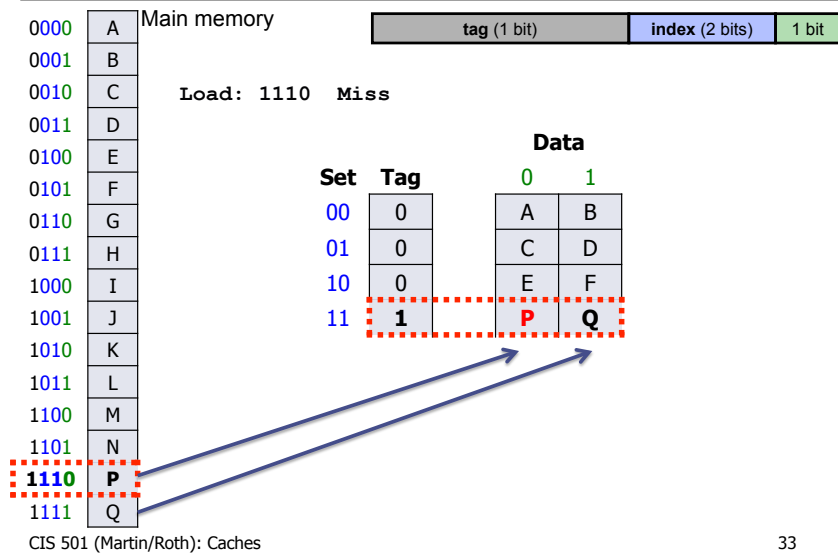


CIS 501 (Martin/Roth): Caches

32



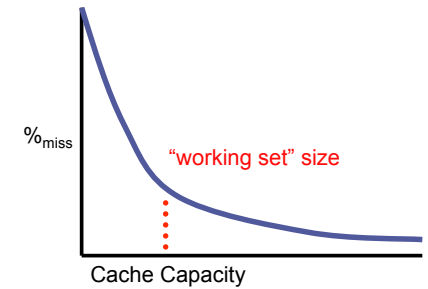
## 4-bit Address, 8B Cache, 2B Blocks



33

## Capacity and Performance

- Simplest way to reduce  $\%_{\text{miss}}$ : increase capacity
  - Miss rate decreases monotonically
    - "Working set"**: insns/data program is actively using
    - Diminishing returns
- However  $t_{\text{hit}}$  increases
  - Latency proportional to  $\sqrt{\text{capacity}}$
- $t_{\text{avg}}$ ?



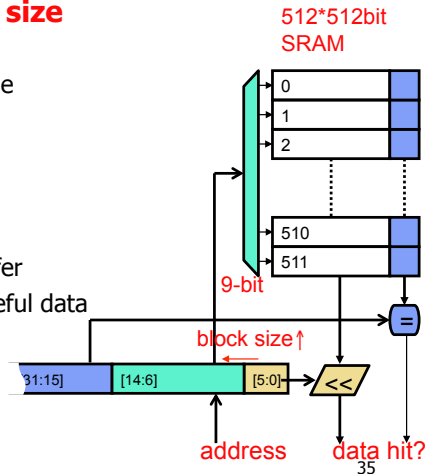
- Given capacity, manipulate  $\%_{\text{miss}}$  by changing **organization**

CIS 501 (Martin/Roth): Caches

34

## Block Size

- Given capacity, manipulate  $\%_{\text{miss}}$  by changing organization
- One option: increase **block size**
  - Exploit **spatial locality**
  - Notice index/offset bits change
  - Tag remain the same
- Ramifications
  - Reduce  $\%_{\text{miss}}$  (up to a point)
  - Reduce tag overhead (why?)
  - Potentially useless data transfer
  - Premature replacement of useful data
  - Fragmentation



CIS 501 (Martin/Roth): Caches

35

## Block Size and Tag Overhead

- Tag overhead of 32KB cache with 1024 32B frames
  - 32B frames → 5-bit offset
  - 1024 frames → 10-bit index
  - 32-bit address – 5-bit offset – 10-bit index = 17-bit tag
  - (17-bit tag + 1-bit valid) \* 1024 frames = 18Kb tags = 2.2KB tags
  - ~6% overhead
- Tag overhead of 32KB cache with 512 64B frames
  - 64B frames → 6-bit offset
  - 512 frames → 9-bit index
  - 32-bit address – 6-bit offset – 9-bit index = 17-bit tag
  - (17-bit tag + 1-bit valid) \* 512 frames = 9Kb tags = 1.1KB tags
  - + ~3% overhead

CIS 501 (Martin/Roth): Caches

36

## Block Size Cache Miss Paper Simulation

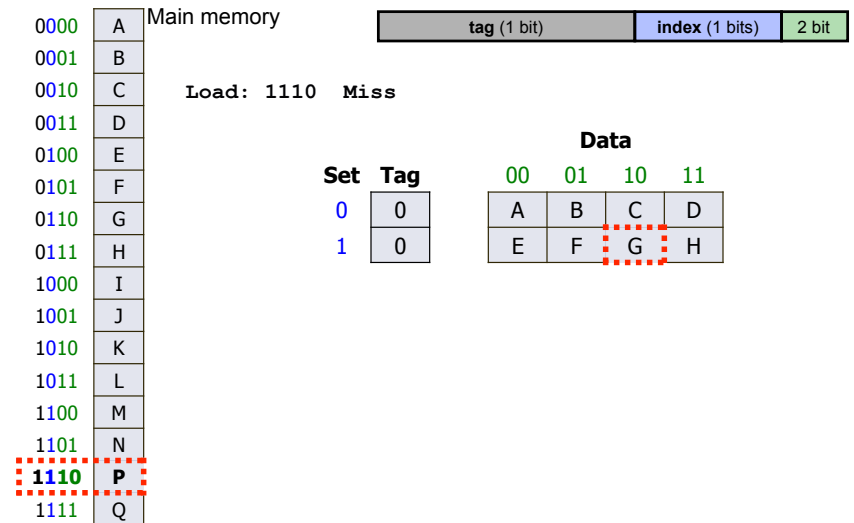
- 8B cache, **4B blocks**

tag (1 bit)	index (1 bit)	2 bits
-------------	---------------	--------

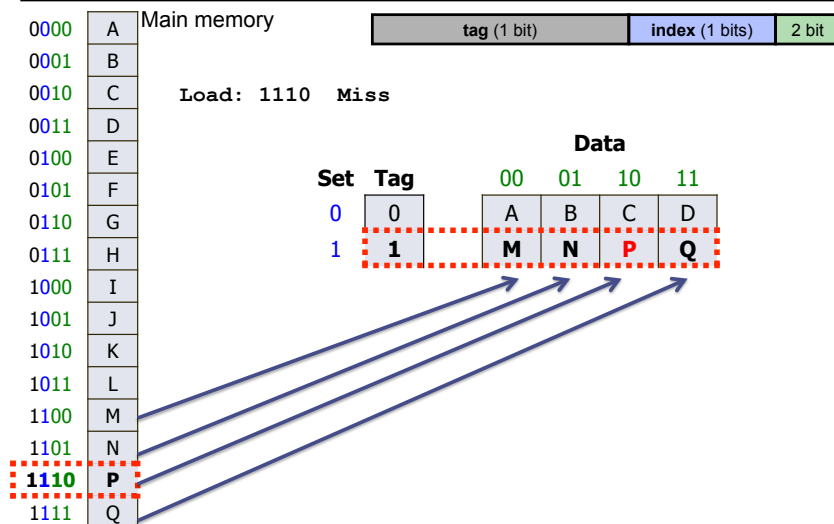
Cache contents (prior to access)		Address	Outcome
Set0	Set1		
0000 0001 0010 0011	0100 0101 0110 0111	1100	Miss
0000 0001 0010 0011	<b>1100 1101 1110 1111</b>	1110	<b>Hit (spatial locality)</b>
0000 0001 0010 0011	1100 1101 1110 1111	1000	Miss
<b>1000 1001 1010 1011</b>	1100 1101 1110 1111	0011	<b>Miss</b>
0000 0001 0010  <b>0011</b>	1100 1101 1110 1111	1000	<b>Miss</b>
<b>1000 1001 1010 1011</b>	1100 1101 1110 1111	0000	Miss
<b>0000 0001 0010 0011</b>	1100 1101 1110 1111	1000	Miss

- + **Spatial "prefetching"**: miss on 1100 brought in 1110
- **Conflicts**: miss on 1000 kicked out 0011

## 4-bit Address, 8B Cache, 4B Blocks

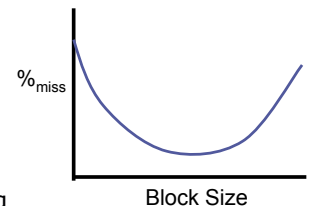


## 4-bit Address, 8B Cache, 4B Blocks



## Effect of Block Size on Miss Rate

- Two effects on miss rate
  - + **Spatial prefetching (good)**
    - For blocks with adjacent addresses
    - Turns miss/miss into miss/hit pairs
  - **Interference (bad)**
    - For blocks with non-adjacent addresses (but in adjacent frames)
    - Turns hits into misses by disallowing simultaneous residence
    - Consider entire cache as one big block
- Both effects always present
  - Spatial prefetching dominates initially
    - Depends on size of the cache
  - Good block size is 16–128B
    - Program dependent



## Block Size and Miss Penalty

- Does increasing block size increase  $t_{\text{miss}}$ ?
  - Don't larger blocks take longer to read, transfer, and fill?
  - They do, but...
- $t_{\text{miss}}$  of an isolated miss is not affected
  - Critical Word First / Early Restart (CRF/ER)**
  - Requested word fetched first, pipeline restarts immediately
  - Remaining words in block transferred/filled in the background
- $t_{\text{miss}}$ 'es of a cluster of misses will suffer
  - Reads/transfers/fills of two misses can't happen at the same time
  - Latencies can start to pile up
  - This is a bandwidth problem (more later)

## Conflicts

- 8B cache, 2B blocks

tag (1 bit)				index (2 bits)		1 bit
Cache contents (prior to access)				Address	Outcome	
Set00	Set01	Set10	Set11			
0000 0001	0010 0011	0100 0101	0110 0111	1100	Miss	
0000 0001	0010 0011	1100 1101	0110 0111	1110	Miss	
0000 0001	0010 0011	1100 1101	1110 1111	1000	Miss	
1000 1001	0010 0011	1100 1101	1110 1111	0011	Hit	
1000 1001	0010 0011	1100 1101	1110 1111	1000	Hit	
1000 1001	0010 0011	1100 1101	1110 1111	0000	Miss	
0000 0001	0010 0011	1100 1101	1110 1111	1000	Miss	

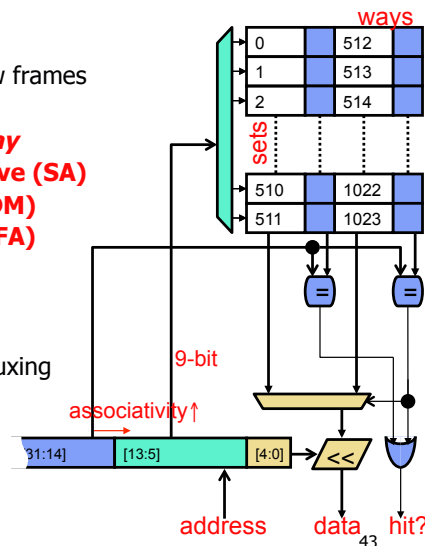
- Pairs like 0000/1000 **conflict**
  - Regardless of block-size (assuming capacity < 16)
  - Q: can we allow pairs like these to simultaneously reside?
  - A: yes, reorganize cache to do so

## Set-Associativity

- Set-associativity**
  - Block can reside in one of few frames
  - Frame groups called **sets**
  - Each frame in set called a **way**
  - This is **2-way set-associative (SA)**
  - 1-way → **direct-mapped (DM)**
  - 1-set → **fully-associative (FA)**

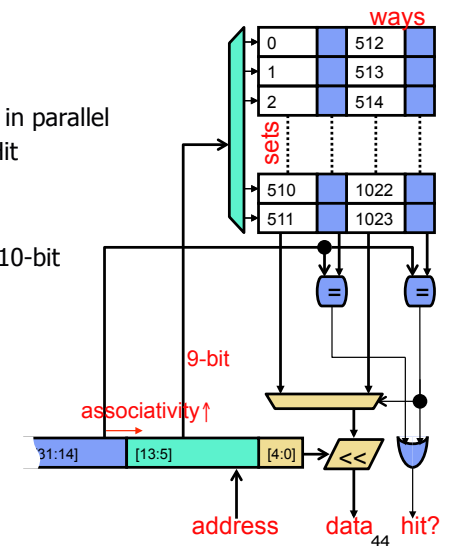
- + Reduces conflicts
- Increases latency<sub>hit</sub>:
  - additional tag match & muxing

- Note: valid bit not shown



## Set-Associativity

- Lookup algorithm
  - Use index bits to find set
  - Read data/tags in all frames in parallel
  - Any** (match and valid bit), Hit
- Notice tag/index/offset bits
  - Only 9-bit index (versus 10-bit for direct mapped)
- Notice block numbering



## Associativity and Miss Paper Simulation

- 8B cache, 2B blocks, **2-way set-associative**

Cache contents (prior to access)				Address	Outcome
Set0.Way0	Set0.Way1	Set1.Way0	Set1.Way1		
0000 0001	0100 0101	0010 0011	0110 0111	1100	Miss
<b>1100 1101</b>	0100 0101	0010 0011	0110 0111	1110	Miss
1100 1101	0100 0101	<b>1110 1111</b>	0110 0111	1000	Miss
1100 1101	<b>1000 1001</b>	1110 1111	0110 0111	0011	<b>Miss (new conflict)</b>
1100 1101	1000 1001	1110 1111	<b>0010 0011</b>	1000	Hit
1100 1101	1000 1001	1110 1111	0010 0011	0000	Miss
<b>0000 0001</b>	1000 1001	1110 1111	0010 0011	1000	<b>Hit (avoid conflict)</b>

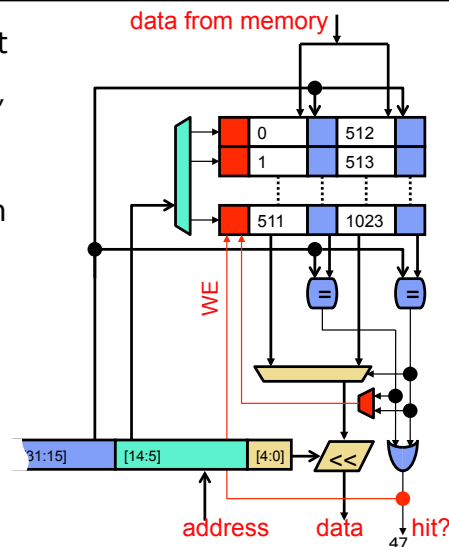
- + **Avoid conflicts**: 0000 and 1000 can both be in set 0
- **Introduce some new conflicts**: notice address re-arrangement
  - Happens, but conflict avoidance usually dominates

## Replacement Policies

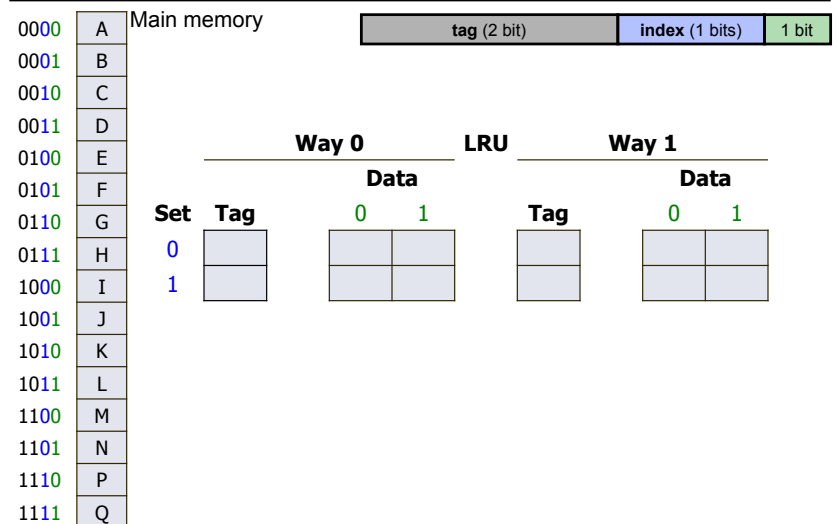
- Set-associative caches present a new design choice
  - On cache miss, which block in set to replace (kick out)?
- Some options
  - Random**
  - FIFO (first-in first-out)**
  - LRU (least recently used)**
    - Fits with temporal locality, LRU = least likely to be used in future
  - NMRU (not most recently used)**
    - An easier to implement approximation of LRU
    - Is LRU for 2-way set-associative caches
  - Belady's**: replace block that will be used furthest in future
    - Unachievable optimum
- Which policy is simulated in previous example?

## NMRU and Miss Handling

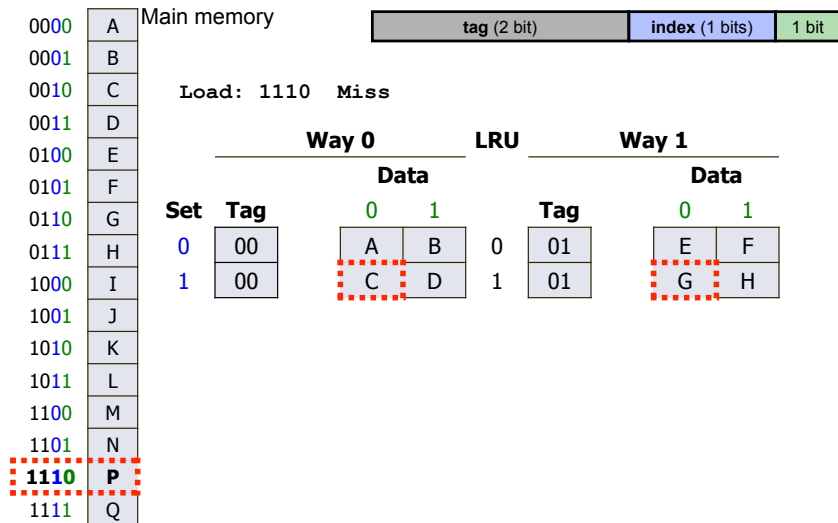
- Add **LRU** field to each set
  - "Least recently used"
  - LRU data is encoded "way"
  - Hit? update MRU
- LRU bits updated on each access



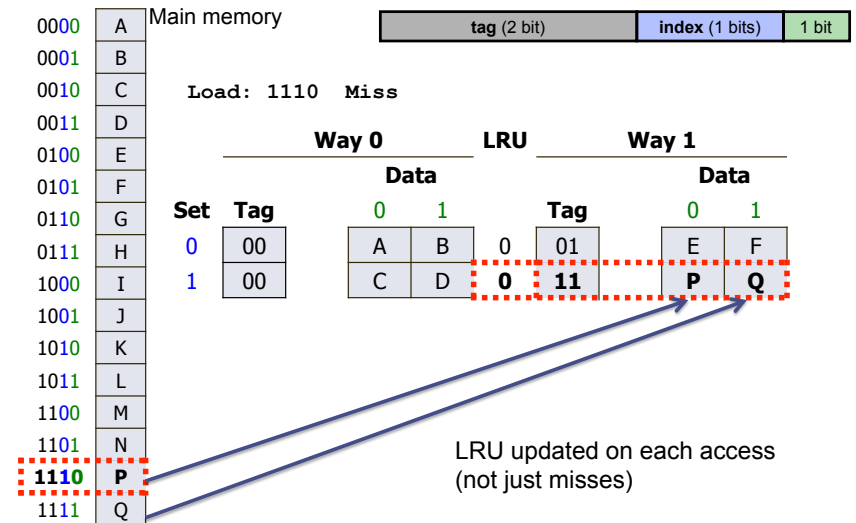
## 4-bit Address, 8B Cache, 2B Blocks, 2-way



## 4-bit Address, 8B Cache, 2B Blocks, 2-way

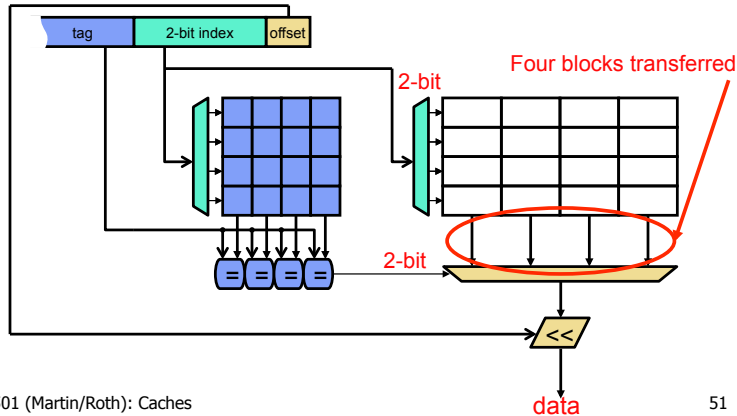


## 4-bit Address, 8B Cache, 2B Blocks, 2-way



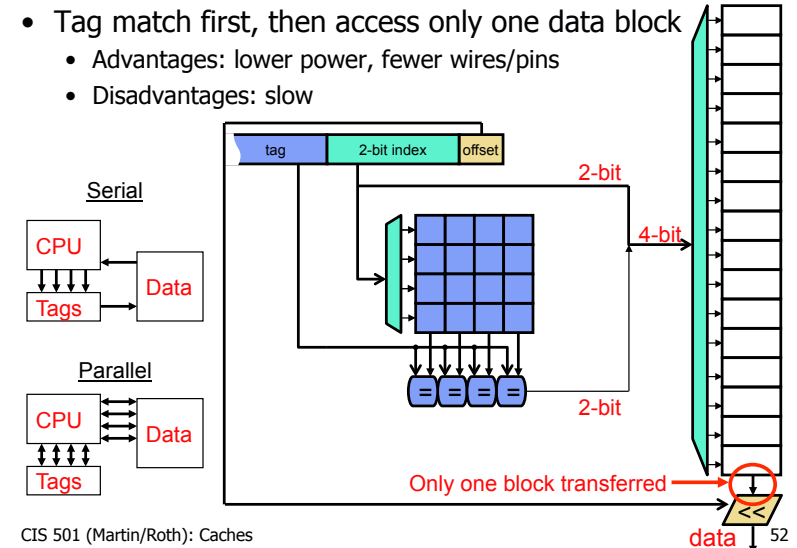
## Parallel or Serial Tag Access?

- Note: data and tags actually physically separate
  - Split into two different arrays
- Parallel access example:



## Serial Tag Access

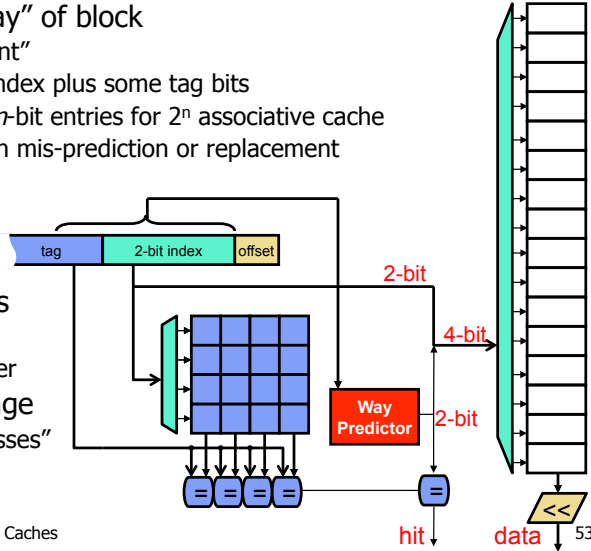
- Tag match first, then access only one data block
  - Advantages: lower power, fewer wires/pins
  - Disadvantages: slow



## Best of Both? Way Prediction

- Predict "way" of block
  - Just a "hint"
  - Use the index plus some tag bits
  - Table of  $n$ -bit entries for  $2^n$  associative cache
  - Update on mis-prediction or replacement

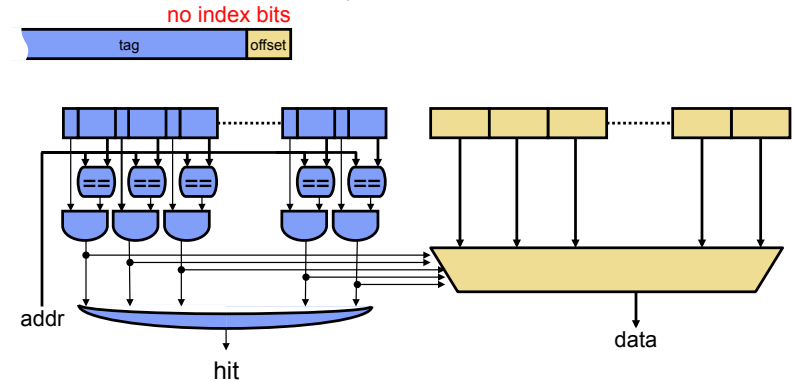
- Advantages
  - Fast
  - Low-power
- Disadvantage
  - More "misses"



CIS 501 (Martin/Roth): Caches

## Full Associativity

- How to implement full (or at least high) associativity?
  - This way is terribly inefficient
  - 1K matches are unavoidable, but 1K data reads + 1K-to-1 mux?

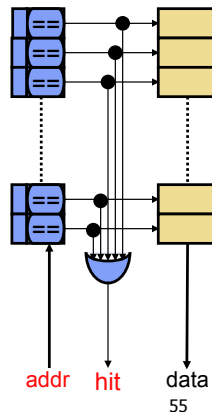


CIS 501 (Martin/Roth): Caches

54

## Full Associativity with CAMs

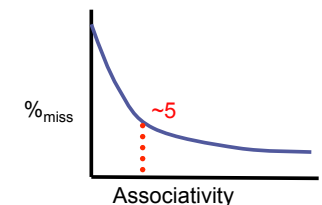
- CAM**: content addressable memory
  - Array of words with built-in comparators
  - Input is data (tag)
  - Output is 1-hot encoding of matching slot
- Fully associative cache
  - Tags as CAM, data as RAM
  - Effective but somewhat expensive
    - But cheaper than any other way
  - Upshot: used for 16-/32-way associativity
    - No good way to build 1024-way associativity
    - + No real need for it, either
- CAMs are used elsewhere, too



CIS 501 (Martin/Roth): Caches

## Associativity and Performance

- Higher associative caches
  - + Have better (lower)  $\%_{\text{miss}}$ 
    - Diminishing returns
  - However  $t_{\text{hit}}$  increases
    - The more associative, the slower
  - What about  $t_{\text{avg}}$ ?



- Block-size and number of sets should be powers of two
  - Makes indexing easier (just rip bits out of the address)
- 3-way set-associativity? No problem

CIS 501 (Martin/Roth): Caches

56

## Classifying Misses: 3C Model (Hill)

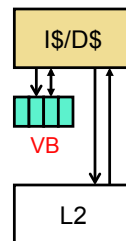
- Divide cache misses into three categories
  - **Compulsory (cold)**: never seen this address before
    - **Would miss even in infinite cache**
  - **Capacity**: miss caused because cache is too small
    - **Would miss even in fully associative cache**
    - Identify? Consecutive accesses to block separated by access to at least N other distinct blocks (N is number of frames in cache)
  - **Conflict**: miss caused because cache associativity is too low
    - Identify? **All other misses**
  - **(Coherence)**: miss due to external invalidations
    - Only in shared memory multiprocessors (later)
- Calculated by multiple simulations
  - Simulate infinite cache, fully-associative cache, normal cache
  - Subtract to find each count

## Miss Rate: ABC

- Why do we care about 3C miss model?
  - So that we know what to do to eliminate misses
  - If you don't have conflict misses, increasing associativity won't help
- **Associativity**
  - + Decreases conflict misses
  - Increases latency<sub>hit</sub>
- **Block size**
  - Increases conflict/capacity misses (fewer frames)
  - + Decreases compulsory/capacity misses (spatial locality)
  - No significant effect on latency<sub>hit</sub>
- **Capacity**
  - + Decreases capacity misses
  - Increases latency<sub>hit</sub>

## Reducing Conflict Misses: Victim Buffer

- Conflict misses: not enough associativity
  - High-associativity is expensive, but also rarely needed
    - 3 blocks mapping to same 2-way set and accessed (XYZ)+
- **Victim buffer (VB)**: small fully-associative cache
  - Sits on I\$/D\$ miss path
  - Small so very fast (e.g., 8 entries)
  - Blocks kicked out of I\$/D\$ placed in VB
  - On miss, check VB: hit? Place block back in I\$/D\$
  - 8 extra ways, shared among all sets
    - + Only a few sets will need it at any given time
  - + Very effective in practice
  - Does VB reduce %<sub>miss</sub> or latency<sub>miss</sub>?



## Overlapping Misses: Lockup Free Cache

- **Lockup free**: allows other accesses while miss is pending
  - Consider: Load [r1] -> r2; Load [r3] -> r4; Add r2, r4 -> r5
  - **Handle misses in parallel**
    - "memory-level parallelism"
  - Makes sense for...
    - Processors that can go ahead despite D\$ miss (out-of-order)
  - Implementation: **miss status holding register (MSHR)**
    - Remember: miss address, chosen frame, requesting instruction
    - When miss returns know where to put block, who to inform
  - Common scenario: "hit under miss"
    - Handle hits while miss is pending
    - Easy
  - Less common, but common enough: "miss under miss"
    - A little trickier, but common anyway
    - Requires multiple MSHRs: search to avoid frame conflicts



## Software Restructuring: Data

- Capacity misses: poor spatial or temporal locality
  - Several code restructuring techniques to improve both
    - Compiler must know that restructuring preserves semantics
- Loop interchange:** spatial locality
  - Example: row-major matrix:  $x[i][j]$  followed by  $x[i][j+1]$
  - Poor code:  $x[i][j]$  followed by  $x[i+1][j]$ 

```
for (j = 0; j < NCOLS; j++)
    for (i = 0; i < NROWS; i++)
        sum += X[i][j];    // say
```
  - Better code

```
for (i = 0; i < NROWS; i++)
    for (j = 0; j < NCOLS; j++)
        sum += X[i][j];
```

## Software Restructuring: Data

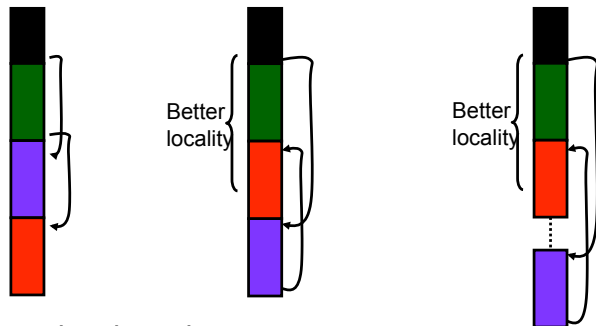
- Loop blocking:** temporal locality
  - Poor code

```
for (k=0; k < NITERATIONS; k++)
    for (i=0; i < NELEMS; i++)
        sum += X[i];    // say
```
  - Better code
    - Cut array into `CACHE_SIZE` chunks
    - Run all phases on one chunk, proceed to next chunk

```
for (i=0; i < NELEMS; i+=CACHE_SIZE)
    for (k=0; k < NITERATIONS; k++)
        for (ii=0; ii < i+CACHE_SIZE-1; ii++)
            sum += X[ii];
```
- Assumes you know `CACHE_SIZE`, do you?
- Loop fusion: similar, but for multiple consecutive loops

## Software Restructuring: Code

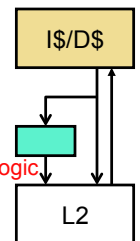
- Compiler can layout code for temporal and spatial locality
  - If (a) { **code1;** } else { **code2;** } **code3;**
  - But, code2 case never happens (say, error condition)



- Fewer taken branches, too
- Intra-procedure, inter-procedure

## Prefetching

- Prefetching:** put blocks in cache proactively/speculatively
  - Key: anticipate upcoming miss addresses accurately
    - Can do in software or hardware
  - Simple example: **next block prefetching**
    - Miss on address  $X \rightarrow$  anticipate miss on  $X + \text{block-size}$
    - + Works for insns: sequential execution
    - + Works for data: arrays
  - Timeliness:** initiate prefetches sufficiently in advance
  - Coverage:** prefetch for as many misses as possible
  - Accuracy:** don't pollute with unnecessary data
    - It evicts useful data



## Software Prefetching

---

- Use a special “prefetch” instruction
  - Tells the hardware to bring in data, doesn’t actually read it
  - Just a hint
- Inserted by programmer or compiler
- Example

```
for (i = 0; i < NROWS; i++)
    for (j = 0; j < NCOLS; j += BLOCK_SIZE) {
        __builtin_prefetch(&X[i][j] + BLOCK_SIZE);
        for (jj = j; jj < j + BLOCK_SIZE - 1; jj++)
            sum += x[i][jj];
    }
```
- Multiple prefetches bring multiple blocks in parallel
  - Using lockup-free caches
  - “Memory-level” parallelism

## Hardware Prefetching

---

- What to prefetch?
  - **Stride-based sequential prefetching**
    - Can also do N blocks ahead to hide more latency
    - + Simple, works for sequential things: insns, array data
    - + Works better than doubling the block size
  - **Address-prediction**
    - Needed for non-sequential data: lists, trees, etc.
    - Use a hardware table to detect strides, common patterns
- When to prefetch?
  - On every reference?
  - On every miss?

## More Advanced Address Prediction

---

- “Next-block” prefetching is easy, what about other options?
- **Correlating predictor**
  - Large table stores (miss-addr → next-miss-addr) pairs
  - On miss, access table to find out what will miss next
    - It’s OK for this table to be large and slow
- Content-directed or dependence-based prefetching
  - Greedily chases pointers from fetched blocks
- Jump pointers
  - Augment data structure with prefetch pointers
- Make it easier to prefetch: cache-conscious layout/malloc
  - Lays lists out serially in memory, makes them look like array
- Active area of research

## Write Issues

---

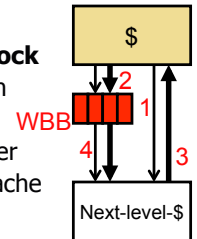
- So far we have looked at reading from cache
  - Instruction fetches, loads
- What about writing into cache
  - Stores, not an issue for instruction caches (why they are simpler)
- Several new issues
  - Tag/data access
  - Write-through vs. write-back
  - Write-allocate vs. write-not-allocate
  - Hiding write miss latency

## Tag/Data Access

- Reads: read tag and data in parallel
  - Tag mis-match → data is garbage (OK, stall until good data arrives)
- Writes: read tag, write data in parallel?
  - Tag mis-match → clobbered data (oops)
  - For associative caches, which way was written into?
- Writes are a pipelined two step (multi-cycle) process
  - Step 1: match tag
  - Step 2: write to matching way
  - Bypass (with address check) to avoid load stalls
  - May introduce structural hazards

## Write Propagation

- When to propagate new value to (lower level) memory?
- **Option #1: Write-through**: immediately
  - On hit, update cache
  - Immediately send the write to the next level
- **Option #2: Write-back**: when block is replaced
  - Requires additional “dirty” bit per block
    - Replace **clean** block: **no extra traffic**
    - Replace **dirty** block: **extra “writeback” of block**
- + **Writeback-buffer (WBB)**: keep it off critical path
  - Step#1: Send “fill” request to next-level
  - Step#2: While waiting, write dirty block to buffer
  - Step#3: When new blocks arrives, put it into cache
  - Step#4: Write buffer contents to next-level



## Write Propagation Comparison

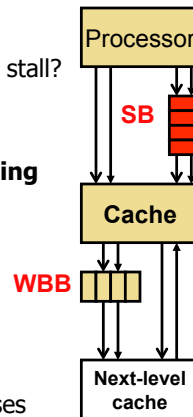
- **Write-through**
  - Requires additional bus bandwidth
    - Consider repeated write hits
  - Next level must handle small writes (1, 2, 4, 8-bytes)
  - + No need for dirty bits in cache
  - + No need to handle “writeback” operations
    - Simplifies miss handling (no write-back buffer)
  - Sometimes used for L1 caches (for example, by IBM)
- **Write-back**
  - + Key advantage: uses less bandwidth
  - Reverse of other pros/cons above
  - Used by Intel and AMD
  - Second-level and beyond are generally write-back caches

## Write Miss Handling

- How is a write miss actually handled?
- **Write-allocate**: fill block from next level, then write it
  - + Decreases read misses (next read to block will hit)
  - Requires additional bandwidth
  - Commonly used (especially with write-back caches)
- **Write-non-allocate**: just write to next level, no allocate
  - Potentially more read misses
  - + Uses less bandwidth
  - Use with write-through

## Write Misses and Store Buffers

- Read miss?
  - Load can't go on without the data, it must stall
- Write miss?
  - Technically, no instruction is waiting for data, why stall?
- Store buffer**: a small buffer
  - Stores put address/value to store buffer, **keep going**
  - Store buffer writes stores to D\$ in the background
  - Loads must search store buffer (in addition to D\$)
  - + Eliminates stalls on write misses (mostly)
  - Creates some problems (later)
- Store buffer vs. writeback-buffer
  - Store buffer: "in front" of D\$, for hiding store misses
  - Writeback buffer: "behind" D\$, for hiding writebacks



73

CIS 501 (Martin/Roth): Caches

## Store Buffer Examples (1 of 2)

- Example #1:
  - Store "1" into address A
    - Miss in cache, put in store buffer (initiate miss)
  - Load from address B
    - Hit in cache, read value from cache
  - Wait for miss to fill, write a "1" to A when done
- Example #2:
  - Store "1" into address A
    - Miss, put in store buffer (initiate miss)
  - Load from address A
    - Miss in cache, but do we stall? Don't need to stall
    - Just bypass load value from the store buffer

CIS 501 (Martin/Roth): Caches

74

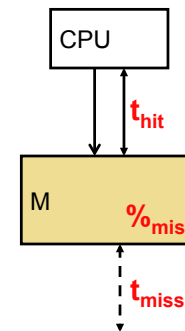
## Store Buffer Examples (2 of 2)

- Example #3:
  - Store byte value "1" into address A
    - Miss in cache, put in store buffer (initiate miss)
  - Load word from address A, A+1, A+2, A+3
    - Hit in cache, read value from where?
    - Read both cache and store buffer (byte-by-byte merge)
- Store buffer holds address, data, and **per-byte** valid bits
- Example #4:
  - Store byte value "1" into address A (initiate miss)
  - Store byte value "2" into address B (initiate miss)
  - Store byte value "3" into Address A (???)
    - Can the first and last store share the same entry?
  - What if "B" fills first?
    - Can the second store leave before the first store?

75

CIS 501 (Martin/Roth): Caches

## Memory Performance Equation



- For memory component M
  - Access**: read or write to M
  - Hit**: desired data found in M
  - Miss**: desired data not found in M
    - Must get from another (slower) component
  - Fill**: action of placing data in M
- $\%_{\text{miss}}$  (miss-rate): #misses / #accesses
- $t_{\text{hit}}$ : time to read data from (write data to) M
- $t_{\text{miss}}$ : time to read data into M

- Performance metric

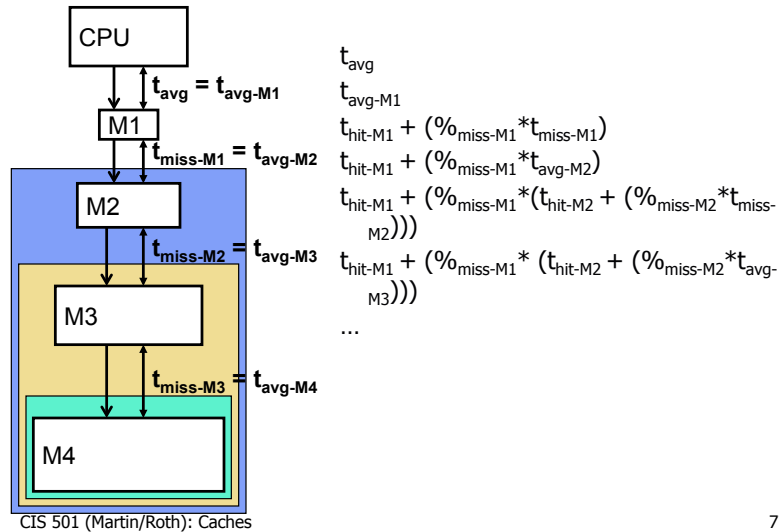
- $t_{\text{avg}}$ : average access time

$$t_{\text{avg}} = t_{\text{hit}} + \%_{\text{miss}} * t_{\text{miss}}$$

CIS 501 (Martin/Roth): Caches

76

## Hierarchy Performance



77

## Local vs Global Miss Rates

- Local hit/miss rate:
  - Percent of references to cache hit (e.g., 90%)
  - Local miss rate is (100% - local hit rate), (e.g., 10%)
- Global hit/miss rate:
  - Misses per instruction (1 miss per 30 instructions)
  - Instructions per miss (3% of instructions miss)
  - Above assumes loads/stores are 1 in 3 instructions
- Consider second-level cache hit rate
  - L1: 2 misses per 100 instructions
  - L2: 1 miss per 100 instructions
  - L2 "local miss rate" -> 50%

CIS 501 (Martin/Roth): Caches

78

## Performance Calculation I

- Parameters
  - Base pipeline CPI = 1
    - In this case, already incorporates  $t_{hit}$
  - Instruction mix: 30% loads/stores
  - I\$:  $\%_{miss} = 2\%$ ,  $t_{miss} = 10$  cycles
  - D\$:  $\%_{miss} = 10\%$ ,  $t_{miss} = 10$  cycles
- What is new CPI?
  - $CPI_{I\$} = \%_{missI\$} * t_{miss} = 0.02 * 10 \text{ cycles} = 0.2 \text{ cycle}$
  - $CPI_{D\$} = \%_{memory} * \%_{missD\$} * t_{missD\$} = 0.30 * 0.10 * 10 \text{ cycles} = 0.3 \text{ cycle}$
  - $CPI_{new} = CPI + CPI_{I\$} + CPI_{D\$} = 1 + 0.2 + 0.3 = 1.5$

CIS 501 (Martin/Roth): Caches

79

## Performance Calculation II

- Parameters
  - Reference stream: all loads
  - D\$:  $t_{hit} = 1\text{ns}$ ,  $\%_{miss} = 5\%$
  - L2:  $t_{hit} = 10\text{ns}$ ,  $\%_{miss} = 20\%$  (local miss rate)
  - Main memory:  $t_{hit} = 50\text{ns}$
- What is  $t_{avgD\$}$  without an L2?
  - $t_{missD\$} = t_{hitM}$
  - $t_{avgD\$} = t_{hitD\$} + \%_{missD\$} * t_{hitM} = 1\text{ns} + (0.05 * 50\text{ns}) = 3.5\text{ns}$
- What is  $t_{avgD\$}$  with an L2?
  - $t_{missD\$} = t_{avgL2}$
  - $t_{avgL2} = t_{hitL2} + \%_{missL2} * t_{hitM} = 10\text{ns} + (0.2 * 50\text{ns}) = 20\text{ns}$
  - $t_{avgD\$} = t_{hitD\$} + \%_{missD\$} * t_{avgL2} = 1\text{ns} + (0.05 * 20\text{ns}) = 2\text{ns}$

CIS 501 (Martin/Roth): Caches

80

## Performance Calculation III

- Memory system parameters
  - D\$:  $t_{\text{hit}} = 1\text{ns}$ ,  $\%_{\text{miss}} = 10\%$ , 50% dirty, writeback-buffer, write-buffer
  - Main memory:  $t_{\text{hit}} = 50\text{ns}$
  - 32-byte block size
- Reference stream: 20% stores, 80% loads
- What is  $t_{\text{avgD\$}}$ ?
  - Write-buffer  $\rightarrow$  hides store misses latency
  - Writeback-buffer  $\rightarrow$  hides dirty writeback latency
  - $t_{\text{missD\$}} = t_{\text{hitM}}$
  - $t_{\text{avgD\$}} = t_{\text{hitD\$}} + \%_{\text{loads}} * \%_{\text{missD\$}} * t_{\text{hitM}} = 1\text{ns} + (0.8 * 0.10 * 50\text{ns}) = 5\text{ns}$

## Bandwidth Calculation

- Memory system parameters
  - D\$:  $t_{\text{hit}} = 1\text{ns}$ ,  $\%_{\text{miss}} = 10\%$ , 50% dirty, writeback-buffer, write-buffer
  - Main memory:  $t_{\text{hit}} = 50\text{ns}$
  - 32-byte block size
- Reference stream: 20% stores, 80% loads
- What is the average bytes transferred per miss?
  - All misses: 32-byte blocks
  - Dirty evictions: 50% of the time \* 32-byte block
  - 48B per miss
- Average bytes transfer per memory operation?
  - 10% of memory operations miss, so 4.8B per memory operation

## Designing a Cache Hierarchy

- For any memory component:  $t_{\text{hit}}$  vs.  $\%_{\text{miss}}$  tradeoff
- Upper components (I\$, D\$) emphasize low  $t_{\text{hit}}$ 
  - Frequent access  $\rightarrow t_{\text{hit}}$  important
  - $t_{\text{miss}}$  is not bad  $\rightarrow \%_{\text{miss}}$  less important
  - Low capacity/associativity (to reduce  $t_{\text{hit}}$ )
  - Small-medium block-size (to reduce conflicts)
- Moving down (L2, L3) emphasis turns to  $\%_{\text{miss}}$ 
  - Infrequent access  $\rightarrow t_{\text{hit}}$  less important
  - $t_{\text{miss}}$  is bad  $\rightarrow \%_{\text{miss}}$  important
  - High capacity/associativity/block size (to reduce  $\%_{\text{miss}}$ )

## Memory Hierarchy Parameters

Parameter	I\$/D\$	L2	L3	Main Memory
$t_{\text{hit}}$	2ns	10ns	30ns	100ns
$t_{\text{miss}}$	10ns	30ns	100ns	10ms (10M ns)
Capacity	8KB–64KB	256KB–8MB	2–16MB	1–4GBs
Block size	16B–64B	32B–128B	32B–256B	NA
Associativity	1–4	4–16	4–16	NA

- Some other design parameters
  - Split vs. unified insns/data
  - Inclusion vs. exclusion vs. nothing
  - On-chip, off-chip, or partially on-chip?

## Split vs. Unified Caches

- **Split I\$/D\$**: insns and data in different caches
  - To minimize structural hazards and  $t_{hit}$
  - Larger unified I\$/D\$ would be slow, 2nd port even slower
  - Optimize I\$ for wide output (superscalar), no writes
  - Why is 486 I/D\$ unified?
- **Unified L2, L3**: insns and data together
  - To minimize  $\%_{miss}$
  - + Fewer capacity misses: unused insn capacity can be used for data
  - More conflict misses: insn/data conflicts
    - A much smaller effect in large caches
  - Insn/data structural hazards are rare: simultaneous I\$/D\$ miss
  - Go even further: unify L2, L3 of multiple cores in a multi-core

## Hierarchy: Inclusion versus Exclusion

- **Inclusion**
  - A block in the L1 is always in the L2
  - Good for write-through L1s (why?)
- **Exclusion**
  - Block is either in L1 or L2 (never both)
  - Good if L2 is small relative to L1
    - Example: AMD's Duron 64KB L1s, 64KB L2
- **Non-inclusion**
  - No guarantees

## Current Cache Research

- "Drowsy Caches"
  - Data/tags allowed to leak away (power)
- "Frequent Value Cache"/"Compressed Cache"
  - Frequent values like 0, 1 compressed (performance, power)
- "Direct Address Cache" + "Cool Cache"
  - Support tag-unchecked loads in compiler and hardware (power)
- "Distance Associative Cache" + "NUCA"
  - Moves frequently used data to closer banks/subarrays
  - Like an associative cache in which not all ways are equal
- "Cache Compression"
  - Use simple compression to increase capacity of secondary caches
  - Saves bandwidth, too

## Summary

- **Average access time** of a memory component
  - $latency_{avg} = latency_{hit} + \%_{miss} * latency_{miss}$
  - Hard to get low  $latency_{hit}$  and  $\%_{miss}$  in one structure → hierarchy
- **Memory hierarchy**
  - Cache (SRAM) → memory (DRAM) → swap (Disk)
  - Smaller, faster, more expensive → bigger, slower, cheaper
- Cache ABCs (**capacity, associativity, block size**)
  - 3C miss model: compulsory, capacity, conflict
- **Performance optimizations**
  - $\%_{miss}$ : victim buffer, prefetching
  - $latency_{miss}$ : critical-word-first/early-restart, lockup-free design
- **Write issues**
  - Write-back vs. write-through/write-allocate vs. write-no-allocate