

U. Wisconsin CS/ECE 752 Advanced Computer Architecture I

Prof. David A. Wood

Unit 8: Storage Hierarchy I: Caches

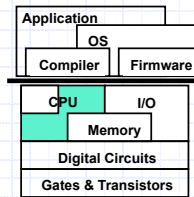
Slides developed by Amir Roth of University of Pennsylvania with sources that included University of Wisconsin slides by Mark Hill, Guri Sohi, Jim Smith, and David Wood.

Slides enhanced by Milo Martin, Mark Hill, and David Wood with sources that included Profs. Asanovic, Falsafi, Hoe, Lipasti, Shen, Smith, Sohi, Vijaykumar, and Wood

CS/ECE 752 (Wood): Caches

1

This Unit: Caches



- Memory hierarchy concepts
- Cache organization
- High-performance techniques
- Low power techniques
- Some example calculations

CS/ECE 752 (Wood): Caches

2

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 50-100ns
 - Naive implementation: loads/stores can be 300x 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

CS/ECE 752 (Wood): Caches

3

Types of Memory

- **Static RAM (SRAM)**
 - 6-10 transistors per bit
 - 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 (>20ns internal access, >40ns pin-to-pin)
 - Different fabrication steps (does not mix well with logic)
- Nonvolatile storage: Magnetic disk, Flash, STT MRAM, etc.

CS/ECE 752 (Wood): Caches

4

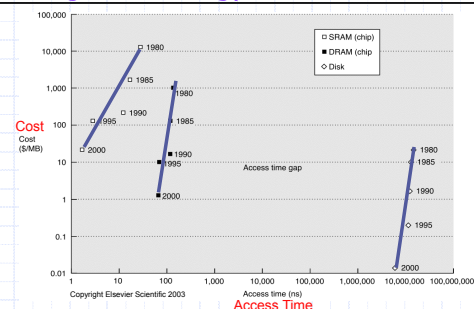
Storage Technology

- **Cost** - what can \$300 buy today?
 - SRAM - 500MB (volume chips)
 - DRAM - 32GB - 64x cheaper than SRAM
 - FLASH - 512 GB - 16X cheaper than DRAM
 - Disk - 8500 GB - 16X cheaper than FLASH
- **Latency**
 - SRAM - <1 to 5ns (on chip)
 - DRAM - ~40ns --- 100x or more slower
 - Disk - 10,000,000ns or 10ms --- 100,000x slower (mechanical)
- **Bandwidth**
 - SRAM - 10-100GB/sec
 - DRAM - ~6GB/sec per channel
 - Disk - 100MB/sec (0.1 GB/sec) - sequential access only

CS/ECE 752 (Wood): Caches

5

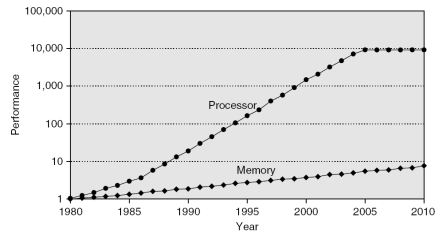
Storage Technology Trends



CS/ECE 752 (Wood): Caches

6

The “Memory Wall”



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

CS/ECE 752 (Wood): Caches

7

Locality to the Rescue

- Locality of memory references**
 - Property of real programs, few exceptions
 - Books and library analogy
- 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

CS/ECE 752 (Wood): Caches

8

Known From the Beginning

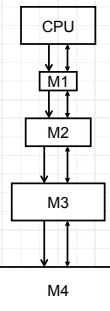
“Ideally, one would desire an infinitely large memory capacity such that any particular word would be immediately available ... We are forced to recognize the possibility of constructing a hierarchy of memories, each of which has a greater capacity than the preceding but which is less quickly accessible.”

Burks, Goldstine, VonNeumann
 “Preliminary discussion of the logical design of an electronic computing instrument”
 IAS memo 1946

CS/ECE 752 (Wood): Caches

9

Exploiting Locality: Memory Hierarchy

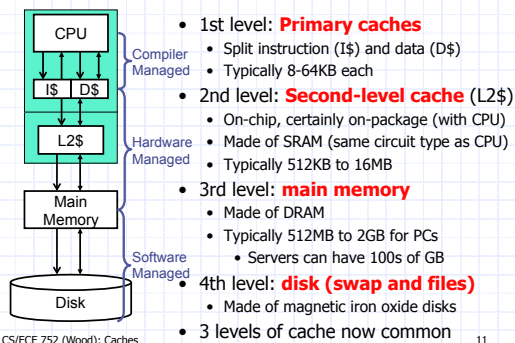


- Hierarchy of memory components
 - Upper components
 - Fast ↔ Small ↔ Expensive
 - Lower components
 - Slow ↔ Big ↔ Cheap
- Connected by buses
 - Which also have latency and bandwidth issues
- Most frequently accessed data in M1
 - M1 + next most frequently accessed in M2, etc.
 - Move data up-down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

CS/ECE 752 (Wood): Caches

10

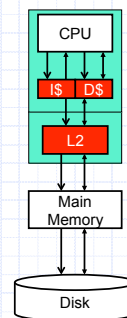
Concrete Memory Hierarchy



CS/ECE 752 (Wood): Caches

11

This Unit: Caches

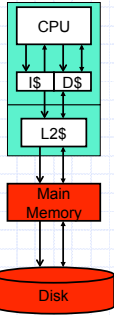


- Cache organization
 - ABC
 - Miss classification
- High-performance techniques
 - Reducing misses
 - Improving miss penalty
 - Improving hit latency
- Low-power techniques
- Some example performance calculations

CS/ECE 752 (Wood): Caches

12

Looking forward: Memory and Disk



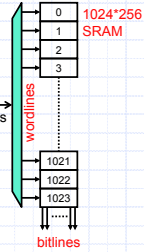
- Main memory
 - Virtual memory
 - DRAM-based memory systems
- Disks and Storage
 - Properties of disks
 - Disk arrays (for performance and reliability)

CS/ECE 752 (Wood): Caches

13

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

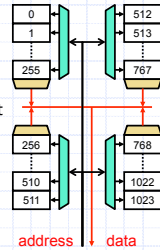
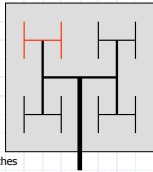


CS/ECE 752 (Wood): Caches

14

Physical Cache 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

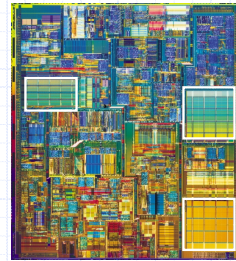


CS/ECE 752 (Wood): Caches

15

Physical Cache Layout

- Arrays and h-trees make caches easy to spot in μ graphs

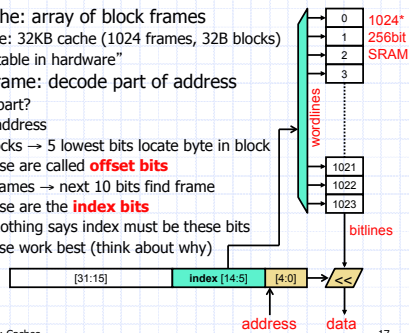


CS/ECE 752 (Wood): Caches

16

Basic Cache Structure

- 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)

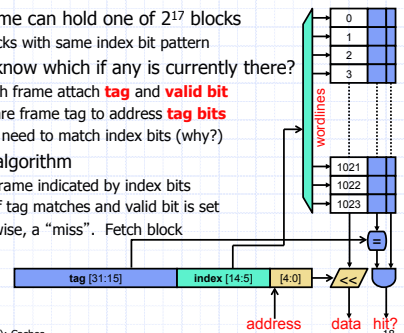


CS/ECE 752 (Wood): Caches

17

Basic Cache Structure

- 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". Fetch block



CS/ECE 752 (Wood): Caches

18

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 → 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
- What about 64-bit addresses?
 - Tag increases to 49bits, ~20% overhead

Cache Performance Simulation

- Parameters: 8-bit addresses, 32B cache, 4B blocks
 - Nibble notation (base 4) tag (3 bits) index (3 bits) 2 bits
 - Initial contents: 0000, 0010, 0020, 0030, 0100, 0110, 0120, 0130

Cache contents (prior to access)	Address	Outcome
0000, 0010, 0020, 0030, 0100, 0110, 0120, 0130	3020	Miss
0000, 0010, 3020 , 0030, 0100, 0110, 0120, 0130	3030	Miss
0000, 0010, 3020, 3030 , 0100, 0110, 0120, 0130	2100	Miss
0000, 0010, 3020, 3030, 2100 , 0110, 0120, 0130	0012	Hit
0000, 0010 , 3020, 3030, 2100, 0110, 0120, 0130	0020	Miss
0000, 0010, 0020 , 3030, 2100, 0110, 0120, 0130	0030	Miss
0000, 0010, 0020, 0030 , 2100, 0110, 0120, 0130	0110	Hit
0000, 0010, 0020, 0030, 2100, 0110 , 0120, 0130	0100	Miss
0000, 1010, 0020, 0030, 0100 , 0110, 0120, 0130	2100	Miss
1000, 1010, 0020, 0030, 2100 , 0110, 0120, 0130	3020	Miss

Hill's 3C Miss Rate Classification

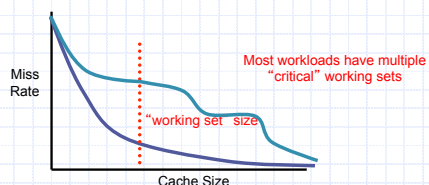
- Compulsory
 - Miss caused by initial access
- Capacity
 - Miss caused by finite capacity
 - I.e., would not miss in infinite cache
- Conflict
 - Miss caused by finite associativity
 - I.e., would not miss in a fully-associative cache
- Coherence (4th C, added by Jouppi)
 - Miss caused by invalidation to enforce coherence

Miss Rate: ABC

- **Associativity**
 - + Decreases conflict misses
 - Increases latency_{hit}
- **Block size**
 - Increases conflict/capacity misses (fewer frames)
 - + Decreases compulsory/capacity misses (spatial prefetching)
 - No effect on latency_{hit}
 - May increase latency_{miss}
- **Capacity**
 - + Decreases capacity misses
 - Increases latency_{hit}

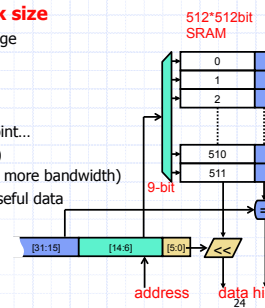
Increase Cache Size

- Biggest caches always have better miss rates
 - However latency_{hit} increases
- Diminishing returns



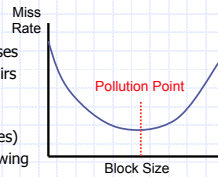
Block Size

- Given capacity, manipulate %_{miss} by changing organization
- One option: increase **block size**
 - Notice index/offset bits change
 - Tag remain the same
- Ramifications
 - + Exploit **spatial locality**
 - Caveat: past a certain point...
 - + Reduce tag overhead (why?)
 - Useless data transfer (needs more bandwidth)
 - Premature replacement of useful data
 - Fragmentation



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
- Both effects always present
 - Spatial prefetching dominates initially
 - Depends on size of the cache
 - Good block size is 16–256B
 - Program dependent



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

Block Size and Performance

- Parameters: 8-bit addresses, 32B cache, **8B blocks**
 - Initial contents : 0000(0010), 0020(0030), 0100(0110), 0120(0130)

	tag (3 bits)	index (2 bits)	3 bits
Cache contents (prior to access)			
0000(0010), 0020(0030), 0100(0110), 0120(0130)	3020		Miss
0000(0010), 3020(3030) , 0100(0110), 0120(0130)	3030		Hit (spatial locality)
0000(0010), 3020(3030), 0100(0110), 0120(0130)	2100		Miss
0000(0010), 3020(3030), 2100(2110) , 0120(0130)	0012		Hit
0000(0010), 3020(3030), 2100(2110), 0120(0130)	0020		Miss
0000(0010), 0020(0030) , 2100(2110), 0120(0130)	0030		Hit (spatial locality)
0000(0010), 0020(0030), 2100(2110), 0120(0130)	0110		Miss (conflict)
0000(0010), 0020(0030), 0100(0110) , 0120(0130)	0100		Hit (spatial locality)
0000(0010), 0020(0030), 0100(0110), 0120(0130)	2100		Miss
0000(0010), 0020(0030), 2100(2110) , 0120(0130)	3020		Miss

Large Blocks and Superblocks

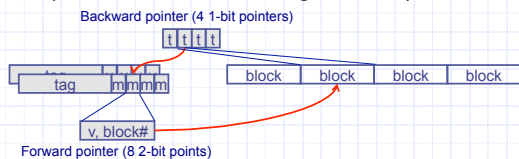
- Large cache blocks can take a long time to refill
 - refill cache line *critical word first*
 - restart cache access before complete refill
- Large cache blocks can waste bus bandwidth if block size is larger than spatial locality
 - group blocks into SuperBlocks (aka Sectors)
 - associate separate valid (coherence) bits for each block
- Sparse access patterns can use 1/S of the cache
 - S is blocks per superblock
 - Internal fragmentation!

SuperBlock: One tag, multiple blocks, one-to-one mapping



Decoupled SuperBlocks

- Level of indirection between tags and blocks
 - Forward pointers
 - Tag points to multiple disjoint blocks
 - Backward pointers
 - Block points to matching tag (and block #)
- Separates notions of data and tag associativity....



Conflicts

- What about pairs like 3030/0030, 0100/2100?
 - These will **conflict** in any sized cache (regardless of block size)
 - Will keep generating misses
- Can we allow pairs like these to simultaneously reside?
 - Yes, reorganize cache to do so

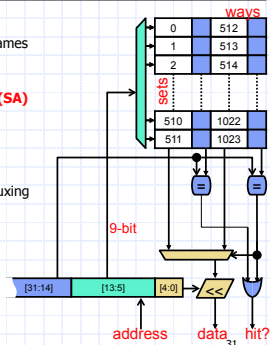
	tag (3 bits)	index (3 bits)	2 bits
Cache contents (prior to access)			
0000, 0010, 0020, 0030, 0100, 0110, 0120, 0130	3020		Miss
0000, 0010, 3020, 0030, 0100, 0110, 0120, 0130	3030		Miss
0000, 0010, 3020, 3030 , 0100, 0110, 0120, 0130	2100		Miss
0000, 0010, 3020, 3030, 2100, 0110, 0120, 0130	0012		Hit
0000, 0010, 3020, 3030, 2100, 0110, 0120, 0130	0020		Miss
0000, 0010, 0020, 3030, 2100, 0110, 0120, 0130	0030		Miss
0000, 0010, 0020, 0030 , 2100, 0110, 0120, 0130	0110		Hit

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 $\text{latency}_{\text{hit}}$ additional muxing

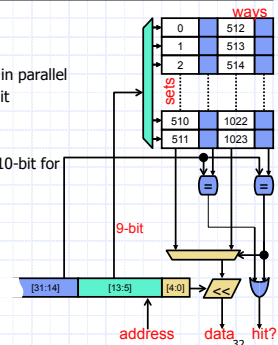
- Note: valid bit not shown



CS/ECE 752 (Wood): Caches

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



CS/ECE 752 (Wood): Caches

Associativity and Performance

- Parameters: 32B cache, 4B blocks, **2-way set-associative**
 - Initial contents : 0000, 0010, 0020, 0030, 0100, 0110, 0120, 0130

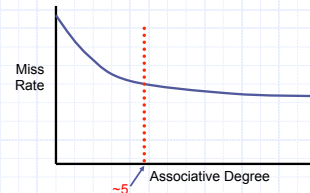
	tag (4 bits)	index (2 bits)	2 bits
Cache contents			
[0000,0100], [0010,0110], [0020,0120], [0030,0130]			
[0000,0100], [0010,0110], [0120,3020], [0030,0130]			
[0000,0100], [0010,0110], [0120,3020], [0130,3030]			
[0100,2100], [0010,0110], [0120,3020], [0130,3030]			
[0100,2100], [0110,0010], [0120,3020], [0130,3030]			
[0100,2100], [0110,0010], [3020,0020], [0130,3030]			
[0100,2100], [0110,0010], [3020,0020], [3030,0030]			
[0100,2100], [0010,0110], [3020,0020], [3030,0030]			
[2100,0100], [0010,0110], [3020,0020], [3030,0030]			
[0100,2100], [0010,0110], [3020,0020], [3030,0030]			

CS/ECE 752 (Wood): Caches

33

Increase Associativity

- Higher associative caches have better miss rates
 - However $\text{latency}_{\text{hit}}$ increases
- Diminishing returns (for a single thread)



CS/ECE 752 (Wood): Caches

34

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)**
 - Is LRU for 2-way set-associative caches
 - **pLRU (pseudo-LRU)**
 - Tree-base generalization of nMRU, uses A-1 bits
 - **Belady's**: replace block that will be used furthest in future
 - Unachievable optimum
 - Which policy is simulated in previous example?

CS/ECE 752 (Wood): Caches

35

Least-Recently Used

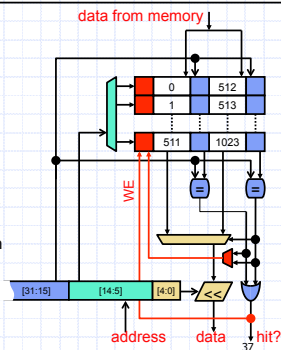
- For $a=2$
 - Single bit per set indicates LRU block
 - Set/clear on each access
- For $a>2$
 - Requires $\text{ceiling}(\log_2(a!))$
 - Simpler to use a list: $\log_2(a)$ bits/block = $a \cdot \log_2(a)$ bits
- For $a>4$, LRU is difficult/expensive
 - Timestamps? How many bits?
 - Must find min timestamp on each eviction
 - Sorted list? Re-sort on every access?

CS/ECE 752 (Wood): Caches

36

NMRU and Miss Handling

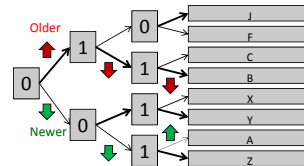
- Key idea: Don't pick the most recently used block
 - MRU data is encoded "way"
 - Hit? update MRU
- Add **MRU** field to each set
 - One bit per way, set on hit
 - Randomly select an unset way
- Alternative implementation
 - One bit per way, set on hit
 - Randomly select an unset way
- MRU/LRU bits updated on each access
- For $a=2$, same as LRU



CS/ECE 752 (Wood): Caches

Pseudo-LRU

- Rather than true LRU, use binary tree
- Each node records which half is older/newer
- Update nodes on each reference
- Follow older pointers to find LRU victim



CS/ECE 752 (Wood): Caches

38

Segmented LRU

- Partition LRU list into *filter* and *reuse* lists
 - On insert, block goes into *filter* list
 - On reuse (hit), block promoted into *reuse* list
- Provides scan & some thrash resistance
 - Blocks without reuse get evicted quickly
 - Blocks with reuse are protected from scan/thrash blocks
- No storage overhead
 - But LRU update slightly more complicated

CS/ECE 752 (Wood): Caches

39

Segmented LRU: LIP/DIP

- LIP – Simplified variant [Qureshi et al. ISCA 2007]
 - Insert new blocks into LRU position, not MRU position
 - Filter list* of size 1, *reuse list* of size $(a-1)$
- DIP — Dynamic LIP
 - Choose to insert in LRU or MRU position
 - Depends upon workload behavior
 - Choose using Set Dueling

CS/ECE 752 (Wood): Caches

40

Set dueling — Dynamic Set Sampling

- Best policy depends upon workload
 - Dynamically decide which is best
- Divide sets into N groups
 - Group 1 uses policy A
 - Group 2 uses policy B
 - Other $N-2$ groups use policy that performed the best in last interval

CS/ECE 752 (Wood): Caches

41

RRIP - [Jaleel et al. ISCA 2010]

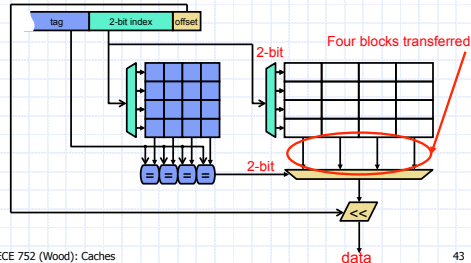
- Re-reference Interval Prediction
 - Extends NRU to multiple bits
 - Start in the middle, promote on hit, demote over time
- Can predict *near-immediate*, *intermediate*, and *distant* re-reference
 - Low overhead: 2 bits/block
 - Static and dynamic variants (like LIP/DIP)

CS/ECE 752 (Wood): Caches

42

Parallel or Serial Tag Access?

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

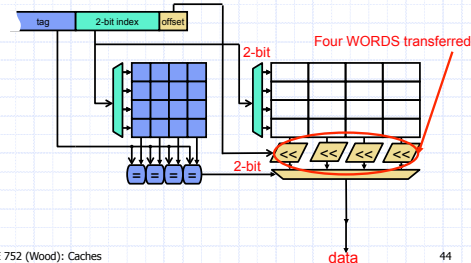


CS/ECE 752 (Wood): Caches

43

Better Parallel Access

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

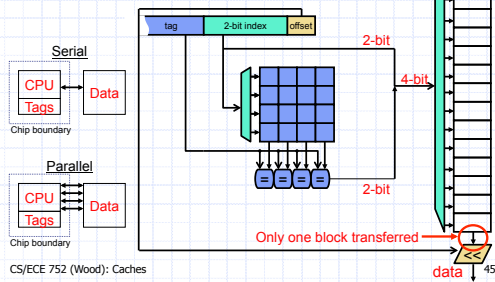


CS/ECE 752 (Wood): Caches

44

Serial Tag Access

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



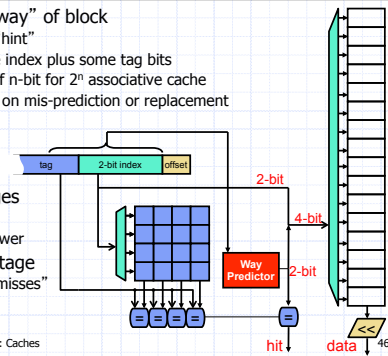
CS/ECE 752 (Wood): Caches

45

Best of Both? Way Prediction

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

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



CS/ECE 752 (Wood): Caches

46

Classifying Misses: 3(4)C Model

- Divide cache misses into three categories
 - Compulsory (cold)**: never seen this address before
 - Would miss even in infinite cache
 - Identify? easy
 - 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
- Who cares? Different techniques for attacking different misses

CS/ECE 752 (Wood): Caches

47

Cache Performance Simulation

- Parameters: 8-bit addresses, 32B cache, 4B blocks
 - Initial contents : 0000, 0010, 0020, 0030, 0100, 0110, 0120, 0130
 - Initial blocks accessed in increasing order

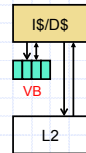
Cache contents	Address	Outcome
0000, 0010, 0020, 0030, 0100, 0110, 0120, 0130	3020	Miss (compulsory)
0000, 0010, 3020 , 0030, 0100, 0110, 0120, 0130	3030	Miss (compulsory)
0000, 0010, 3020, 3030 , 0100, 0110, 0120, 0130	2100	Miss (compulsory)
0000, 0010, 3020, 3030, 2100 , 0110, 0120, 0130	0012	Hit
0000, 0010, 3020, 3030, 2100, 0110, 0120, 0130	0020	Miss (capacity)
0000, 0010, 0020 , 3030, 2100, 0110, 0120, 0130	0030	Miss (capacity)
0000, 0010, 0020, 0030 , 2100, 0110, 0120, 0130	0110	Hit
0000, 0010, 0020, 0030, 2100, 0110, 0120, 0130	0100	Miss (capacity)
0000, 1010, 0020, 0030, 0100 , 0110, 0120, 0130	2100	Miss (conflict)
1000, 1010, 0020, 0030, 2100 , 0110, 0120, 0130	3020	Miss (conflict)

CS/ECE 752 (Wood): Caches

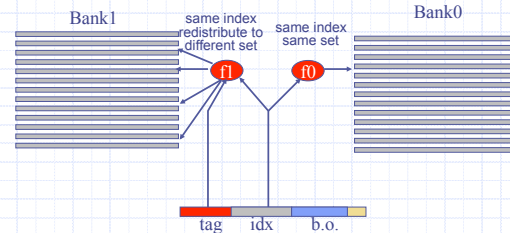
48

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 (ABC)*
- Victim buffer (VB):** small fully-associative cache
 - Sits on I\$/D\$ fill 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 for small caches
 - Does VB reduce **%miss** or **latency_{miss}**?



Seznec's Skewed-Associative Cache



Can get better utilization with less assoc?
average case? worst case? Replacement Policy?

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 $x[i+1][j]$
 - 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]; // non-contiguous accesses
```
 - Better code

```
for (i = 0; i < NROWS; i++)
    for (j = 0; j < NCOLS; j++)
        sum += x[i][j]; // contiguous accesses
```

Software Restructuring: Data

- Loop blocking:** temporal locality
 - Poor code

```
for (k=0; k < NITERATIONS; k++)
    for (i=0; i < NLEMS; 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 < NLEMS; 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

Restructuring Loops

- Loop Fusion
 - Merge two independent loops
 - Increase reuse of data
- Loop Fission
 - Split loop into independent loops
 - Reduce contention for cache resources

Fusion Example:

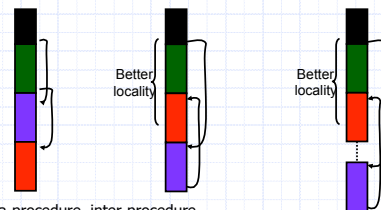
```
for (i=0; i < N; i++)
    for (j=0; j < N; j++)
        a[j][j] = 1/b[j][j]*c[j][j];
for (i=0; i < N; i++)
    for (j=0; j < N; j++)
        d[j][j] = a[j][j]+c[j][j];
```

Fused Loop:

```
for (i=0; i < N; i++)
    for (j=0; j < N; j++)
    {
        a[j][j] = 1/b[j][j]*c[j][j];
        d[j][j] = a[j][j]+c[j][j];
    }
```

Software Restructuring: Code

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



- Intra-procedure, inter-procedure
- Related to trace scheduling

Miss Cost: Critical Word First/Early Restart

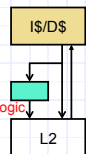
- Observation: $\text{latency}_{\text{miss}} = \text{latency}_{\text{access}} + \text{latency}_{\text{transfer}}$
 - $\text{latency}_{\text{access}}$: time to get first word
 - $\text{latency}_{\text{transfer}}$: time to get rest of block
 - Implies whole block is loaded before data returns to CPU
- Optimization
 - **Critical word first**: return requested word first
 - Must arrange for this to happen (bus, memory must cooperate)
 - **Early restart**: send requested word to CPU immediately
 - Get rest of block load into cache in parallel
 - $\text{latency}_{\text{miss}} = \text{latency}_{\text{access}}$

Miss Cost: Lockup Free Cache

- **Lockup free**: allows other accesses while miss is pending
 - Consider: Load [r1] -> r2; Load [r3] -> r4; Add r2, r4 -> r5
 - Only makes sense for...
 - Data cache
 - 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
- Simplest scenario: "hit under miss"
 - Handle hits while miss is pending
 - Easy for OoO cores
- More common: "miss under miss"
 - A little trickier, but common anyway
 - Requires split-transaction bus/interconnect
 - Requires multiple MSHRs: search to avoid frame conflicts

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** -> anticipate miss on **X+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



Prefetching Characterization

- Useless prefetch
 - Prefetch brings in data that is never used
- Harmful prefetch
 - Prefetch brings in data that displaces a block that WOULD have been reused.
 - So-called (negative) interference
 - Positive interference is possible in multicores
- Performance v. Energy
 - Original characterization deals with performance
 - Useless prefetches unnecessarily use energy, thus are harmful

Software Prefetching

- Software prefetching: two kinds
 - **Binding**: prefetch into register (e.g., software pipelining)
 - + No ISA support needed, use normal loads (non-blocking cache)
 - Need more registers, and what about faults?
 - **Non-binding**: prefetch into cache (or other buffer) only
 - Need ISA support: non-binding, non-faulting loads
 - + Simpler semantics
- Example


```
for (i = 0; i < NROWS; i++)
  for (j = 0; j < NCOLS; j += BLOCK_SIZE) {
    prefetch(&X[i][j] + BLOCK_SIZE);
    for (jj = j; jj < j + BLOCK_SIZE - 1; jj++)
      sum += x[i][jj];
  }
```

Hardware Prefetching

- What to prefetch?
 - One block ahead
 - How much latency do we need to hide (Little's Law)?
 - Can also do N blocks ahead to hide more latency
 - + Simple, works for sequential things: insns, array data
 - **Address-prediction**
 - Needed for non-sequential data: lists, trees, etc.
- When to prefetch?
 - On every reference?
 - On every miss?
 - + Works better than doubling the block size
 - Ideally: when resident block becomes dead (avoid useful evictions)
 - How to know when that is? ["Dead-Block Prediction", ISCA'01]

Address Prediction for Prefetching

- “Next-block” prefetching is easy, what about other options?
- **Correlating predictor**
 - Large table stores (miss-addr \rightarrow 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
 - Can do in hardware too
- An active area of research

CS/ECE 752 (Wood): Caches

61

Write Issues

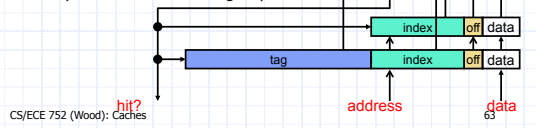
- So far we have looked at reading from cache (loads)
- What about writing into cache (stores)?
- Several new issues
 - Tag/data access
 - Write-through vs. write-back
 - Write-allocate vs. write-not-allocate
- Buffers
 - Store buffers (queues)
 - Write buffers
 - Writeback buffers

CS/ECE 752 (Wood): Caches

62

Tag/Data Access

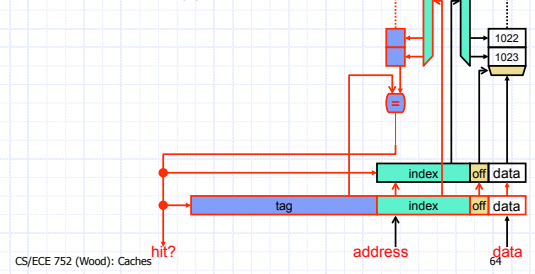
- Reads: read tag and data in parallel
 - Tag mis-match \rightarrow data is garbage (OK)
- Writes: read tag, write data in parallel?
 - Tag mis-match \rightarrow clobbered data (oops)
 - For associative cache, which way is written?
- Writes are a pipelined 2 cycle process
 - Cycle 1: match tag
 - Cycle 2: write to matching way



CS/ECE 752 (Wood): Caches

Tag/Data Access

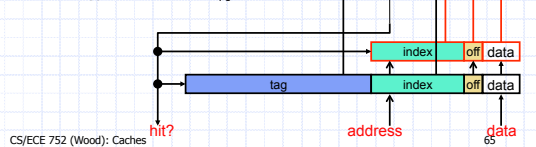
- Cycle 1: check tag
 - Hit? Advance “store pipeline”
 - Miss? Stall “store pipeline”



CS/ECE 752 (Wood): Caches

Tag/Data Access

- Cycle 2: write data
- Advanced Technique
 - Decouple write pipeline
 - In the same cycle
 - Check tag of store_i
 - Write data of store_{i-1}
 - Bypass data of store_{i-1} to loads



CS/ECE 752 (Wood): Caches

Write-Through vs. Write-Back

- When to propagate new value to (lower level) memory?
 - **Write-through:** immediately
 - + Conceptually simpler
 - + Uniform latency on misses
 - Requires additional bus bandwidth
 - **Write-back:** when block is replaced
 - + Requires additional “dirty” bit per block
 - + Lower bus bandwidth for large caches
 - Only writeback dirty blocks
 - Non-uniform miss latency
 - Clean miss: one transaction with lower level (fill)
 - Dirty miss: two transactions (writeback + fill)
 - Writeback buffer: fill, then writeback (later)
- Common design: Write through L1, write-back L2/L3

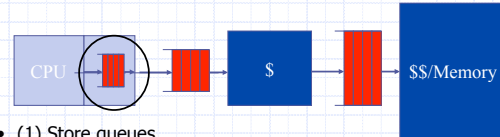
CS/ECE 752 (Wood): Caches

66

Write-allocate vs. Write-non-allocate

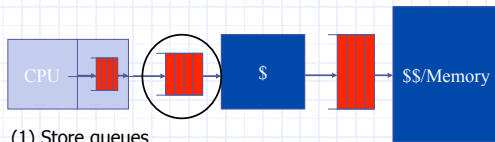
- What to do on a write miss?
 - Write-allocate:** read block from lower level, write value into it
 - + Decreases read misses
 - Requires additional bandwidth
 - Used mostly with write-back
 - Write-non-allocate:** just write to next level
 - Potentially more read misses
 - + Uses less bandwidth
 - Used mostly with write-through
- Write allocate is common for write-back
 - Write-non-allocate for write through

Buffering Writes 1 of 3: Store Queues



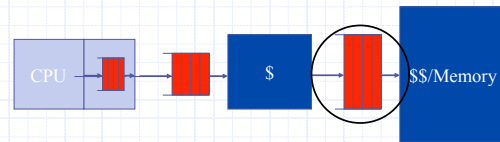
- (1) Store queues
 - Part of speculative processor; transparent to architecture
 - Hold speculatively executed stores
 - May rollback store if earlier exception occurs
 - Used to track load/store dependencies
- (2) Write buffers
- (3) Writeback buffers

Buffering Writes 2 of 3: Write Buffer



- (1) Store queues
- (2) Write buffers
 - Holds committed architectural state
 - Transparent to single thread
 - May affect memory consistency model
 - Hides latency of memory access or cache miss
 - May bypass values to later loads (or stall)
 - Store queue & write buffer may be in same physical structure
- (3) Writeback buffers

Buffering Writes 3 of 3: Writeback Buffer



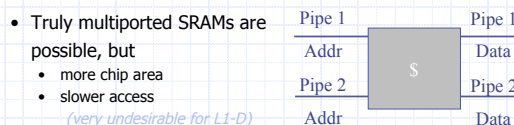
- (1) Store queues
- (2) Write buffers
- (3) Writeback buffers (Special case of Victim Buffer)
 - Transparent to architecture
 - Holds victim block(s) so miss/prefetch can start immediately
 - (Logically part of cache for multiprocessor coherence)

Increasing Cache Bandwidth

- What if we want to access the cache twice per cycle?
- Option #1: multi-ported cache
 - Same number of six-transistor cells
 - Double the decoder logic, bitlines, wordlines
 - Areas becomes "wire dominated" -> slow
 - **OR**, time multiplex the wires
- Option #2: banked cache
 - Split cache into two smaller "banks"
 - Can do two parallel access to different parts of the cache
 - Bank conflict occurs when two requests access the same bank
- Option #3: replication
 - Make two copies (2x area overhead)
 - Writes both replicas (does not improve write bandwidth)
 - Independent reads
 - No bank conflicts, but lots of area
 - Split instruction/data caches is a special case of this approach

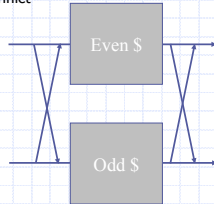
Multi-Port Caches

- Superscalar processors requires multiple data references per cycle
- Time-multiplex a single port (double pump)
 - need cache access to be faster than datapath clock
 - not scalable



Multi-Banking (Interleaving) Caches

- Address space is statically partitioned and assigned to different caches *Which addr bit to use for partitioning?*
- A compromise (e.g. Intel P6, MIPS R10K)
 - multiple references per cyc. if no conflict
 - only one reference goes through if conflicts are detected
 - the rest are deferred
(bad news for scheduling logic)
- Most helpful if compiler knows about the interleaving rules

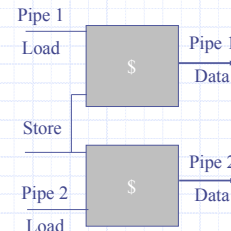


CS/ECE 752 (Wood): Caches

73

Multiple Cache Copies: e.g. Alpha 21164

- Independent fast load paths
- Single shared store path
- Not a scalable solution
 - Store is a bottleneck
 - Doubles area



CS/ECE 752 (Wood): Caches

74

Evaluation Methods

- The three system evaluation methodologies
 - Hardware counters
 - Analytic modeling
 - Software simulation
 - Hardware prototyping and measurement

CS/ECE 752 (Wood): Caches

75

Methods: Hardware Counters

- See Clark, TOCS 1983
 - ✓ accurate
 - ✓ realistic workloads, system + user + others
 - ✗ difficult, why?
 - ✗ must first have the machine
 - ✗ hard to vary cache parameters
 - ✗ experiments not deterministic
 - ✗ use statistics!
 - ✗ take multiple measurements
 - ✗ compute mean and confidence measures
- Most modern processors have built-in hardware counters

CS/ECE 752 (Wood): Caches

76

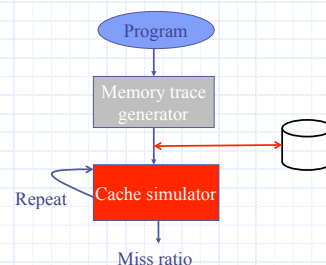
Methods: Analytic Models

- Mathematical expressions
 - ✓ insightful: can vary parameters
 - ✓ fast
 - ✗ absolute accuracy suspect for models with few parameters
 - ✗ hard to determine parameter values
 - ✗ difficult to evaluate cache interaction with system
 - ✗ bursty behavior hard to evaluate

CS/ECE 752 (Wood): Caches

77

Methods: Trace-Driven Simulation



CS/ECE 752 (Wood): Caches

78

Methods: Trace-Driven Simulation

- ✓ experiments repeatable
- ✓ can be accurate
- ✓ much recent progress
- ✗ reasonable traces are very large (gigabytes?)
- ✗ simulation is time consuming
- ✗ hard to say if traces are representative
- ✗ don't directly capture speculative execution
- ✗ don't model interaction with system

- ✗ Widely used in industry

Methods: Execution-Driven Simulation

- Simulate the program execution
 - simulates each instruction's execution on the computer
 - model processor, memory hierarchy, peripherals, etc.
- ✓ reports execution time
 - ✓ accounts for all system interactions
- ✓ no need to generate/store trace
- ✗ much more complicated simulation model
- ✗ time-consuming but good programming can help
- ✗ multi-threaded programs exhibit variability
- ✗ Very common in academia and industry today

- ✗ Watch out for repeatability in multithreaded workloads

Low-Power Caches

- Caches consume significant power
 - 15% in Pentium4
 - 45% in StrongARM
- Three techniques
 - Way prediction (already talked about)
 - Dynamic resizing
 - Drowsy caches

Low-Power Access: Dynamic Resizing

- **Dynamic cache resizing**
 - Observation I: data, tag arrays implemented as many small arrays
 - Observation II: many programs don't fully utilize caches
- Idea: dynamically turn off unused arrays
 - Turn off means disconnect power (V_{DD}) plane
 - + Helps with both dynamic and static power
- There are always tradeoffs
 - Flush dirty lines before powering down → costs power↑
 - Cache-size↓ → %_{miss}↑ → power↑, execution time↑

Dynamic Resizing: When to Resize

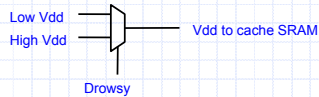
- Use %_{miss} feedback
 - %_{miss} near zero? Make cache smaller (if possible)
 - %_{miss} above some threshold? Make cache bigger (if possible)
- Aside: how to track miss-rate in hardware?
 - Hard, easier to track miss-rate vs. some threshold
 - Example: is %_{miss} higher than 5%?
 - N-bit counter (N = 8, say)
 - Hit? counter -= 1
 - Miss? counter += 19
 - Counter positive? More than 1 miss per 19 hits (%_{miss} > 5%)

Dynamic Resizing: How to Resize?

- **Reduce ways**
 - ["Selective Cache Ways", Albonesi, ISCA-98]
 - + Resizing doesn't change mapping of blocks to sets → simple
 - Lose associativity
- **Reduce sets**
 - ["Resizable Cache Design", Yang+, HPCA-02]
 - Resizing changes mapping of blocks to sets → tricky
 - When cache made bigger, need to relocate some blocks
 - Actually, just flush them
 - Why would anyone choose this way?
 - + More flexibility: number of ways typically small
 - + Lower %_{miss}: for fixed capacity, higher associativity better

Drowsy Caches

- Circuit technique to reduce leakage power
 - Lower Vdd → Much lower leakage
 - But too low Vdd → Unreliable read/destructive read
- Key: Drowsy state (low Vdd) to hold value w/ low leakage
- Key: Wake up to normal state (high Vdd) to access
 - 1-3 cycle additional latency



CS/ECE 752 (Wood): Caches

85

Memory Hierarchy Design

- Important: design hierarchy components together
- **I\$, D\$:** optimized for latency_{hit} and parallel access
 - Insns/data in separate caches (3-4 cycle latency)
 - Capacity: 8-64KB, block size: 16-64B, associativity: 1-4
 - Power: parallel tag/data access, way prediction?
 - Bandwidth: banking or multi-porting/replication
 - Other: write-through or write-back
- **L2:** optimized for %_{miss} power (latency_{hit}: 6-12 cycles)
 - Insns and data in one cache (for higher utilization, %_{miss})
 - Capacity: 128KB-2MB, block size: 64-256B, associativity: 4-16
 - Power: parallel or serial tag/data access, banking
 - Bandwidth: banking
 - Other: write-back
- **L3:** starting to appear (latency_{hit} = 24-30 cycles)

CS/ECE 752 (Wood): Caches

86

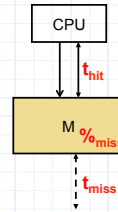
Hierarchy: Inclusion versus Exclusion

- 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

CS/ECE 752 (Wood): Caches

87

Memory Performance Equation



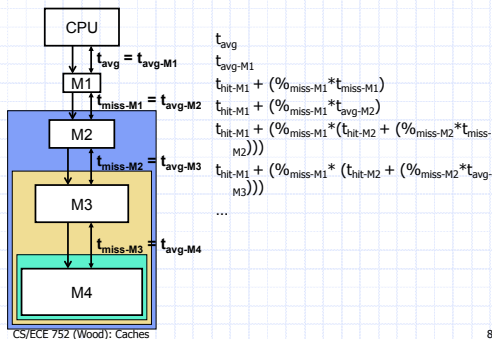
- For memory component M
 - **Access:** read or write to M
 - **Hit:** desired data found in M
 - **Miss:** desired data not found in M
 - Must get from another (slower) component
 - **Fill:** action of placing data in M
- %_{miss} (miss-rate): #misses / #accesses
- t_{hit}: time to read data from (write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

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

CS/ECE 752 (Wood): Caches

88

Hierarchy Performance



CS/ECE 752 (Wood): Caches

89

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%

CS/ECE 752 (Wood): Caches

90

Performance Calculation I

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

Performance Calculation II

- In a pipelined processor, I\$/D\$ t_{hit} is “built in” (effectively 0)
- Parameters
 - Base pipeline CPI = 1
 - Instruction mix: 30% loads/stores
 - I\$: $\%_{miss} = 2\%$, $t_{miss} = 10$ cycles
 - D\$: $\%_{miss} = 10\%$, $t_{miss} = 10$ cycles
- What is new CPI?
 - $CPI_{IS} = \%_{missIS} * t_{miss} = 0.02 * 10 \text{ cycles} = 0.2 \text{ cycle}$
 - $CPI_{DS} = \%_{memory} * \%_{missD\$} * t_{missD\$} = 0.30 * 0.10 * 10 \text{ cycles} = 0.3 \text{ cycle}$
 - $CPI_{new} = CPI + CPI_{IS} + CPI_{DS} = 1 + 0.2 + 0.3 = 1.5$

An Energy Calculation

- Parameters
 - 2-way SA D\$
 - 10% miss rate
 - 5 μ W/access tag way, 10 μ W/access data way
- What is power/access of parallel tag/data design?
 - Parallel: each access reads both tag ways, both data ways
 - Misses write additional tag way, data way (for fill)
 - $[2 * 5\mu W + 2 * 10\mu W] + [0.1 * (5\mu W + 10\mu W)] = 31.5 \mu W/\text{access}$
- What is power/access of serial tag/data design?
 - Serial: each access reads both tag ways, one data way
 - Misses write additional tag way (actually...)
 - $[2 * 5\mu W + 0.9 * 10\mu W] + [0.1 * (5\mu W + 10\mu W)] = 20.5 \mu W/\text{access}$

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
- Power optimizations**: way prediction, dynamic resizing
- Write issues**
 - Write-back vs. write-through/write-allocate vs. write-no-allocate

Backups

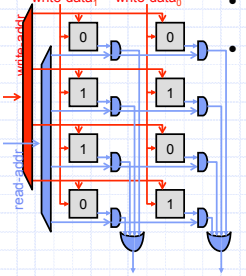
SRAM Technology

- SRAM**: static RAM
 - Static**: bits directly connected to power/ground
 - Naturally/continuously “refreshed”, never decay
 - Designed for speed
- Implements all storage arrays in real processors
 - Register file, caches, branch predictor, etc.
 - Everything except pipeline latches
- Latches vs. SRAM
 - Latches: singleton word, always read/write same one
 - SRAM: array of words, always read/write different one
 - Address indicates which one

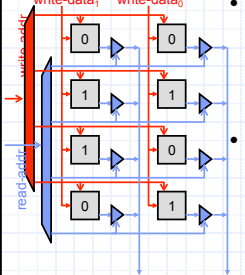
(CMOS) Memory Components

- Interface
 - N-bit **address** bus (on N-bit machine)
 - Data** bus
 - Typically read/write on same data bus
 - Can have multiple **ports**: address/data bus pairs
 - Can be **synchronous**: read/write on clock edges
 - Can be **asynchronous**: untimed "handshake"

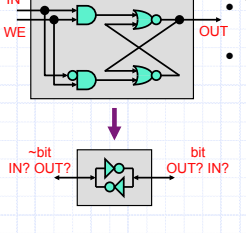
SRAM: First Cut

- 
- 4x2 (4 2-bit words) RAM
 - 2-bit addr
 - First cut: bits are D-Latches
 - Write port**
 - Addr **decodes** to enable signals
 - Read port**
 - Addr **decodes** to mux selectors
 - 1024 input OR gate?
 - Physical layout of output wires
 - RAM width \propto M
 - Wire delay \propto wire length

SRAM: Second Cut

- 
- Second cut: tri-state wired-OR
 - Read mux using **tri-states**
 - Scalable, distributed "muxes"
 - Better layout of output wires
 - RAM width independent of M
 - Standard RAM**
 - Bits in word connected by **wordline**
 - 1-hot decode address
 - Bits in position connected by **bitline**
 - Shared input/output wires
 - Port**: one set of wordlines/bitlines
 - Grid-like design

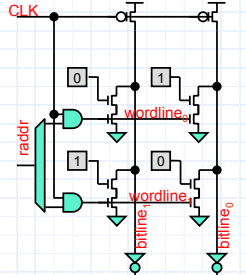
SRAM: Third Cut

- 
- Third cut: replace latches with...
 - 28 transistors per bit
 - Cross-coupled inverters (CCI)**
 - 4 transistors
 - Convention
 - Right node is **bit**, left is **~bit**
 - Non-digital interface
 - What is the input and output?
 - Where is write enable?
 - Implement ports in "analog" way
 - Transistors, not full gates

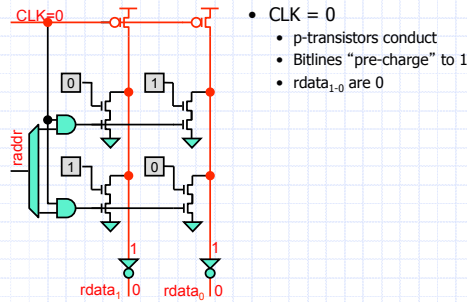
SRAM: Register Files and Caches

- Two different SRAM port styles
 - Regfile style**
 - Modest size: <4KB
 - Many ports: some read-only, some write-only
 - Write and read both take half a cycle (write first, read second)
 - Cache style**
 - Larger size: >8KB
 - Few ports: read/write in a single port
 - Write and read can both take full cycle

Regfile-Style Read Port

- 
- Two phase read
 - Phase I: $\text{clk} = 0$
 - Pre-charge bitlines to 1
 - Negated bitlines are 0
 - Phase II: $\text{clk} = 1$
 - One wordline goes high
 - All "1" bits in that row discharge their bitlines to 0
 - Negated bitlines go to 1

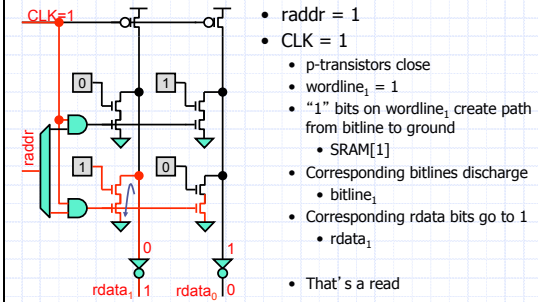
Read Port In Action: Phase I



CS/ECE 752 (Wood): Caches

103

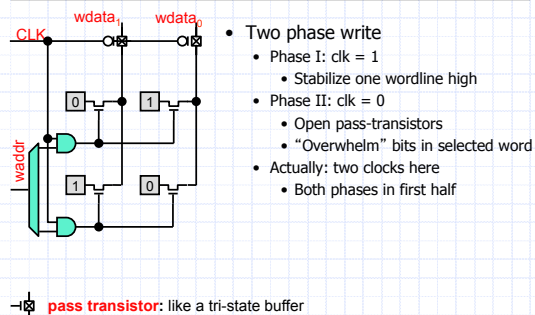
Read Port In Action: Phase II



CS/ECE 752 (Wood): Caches

104

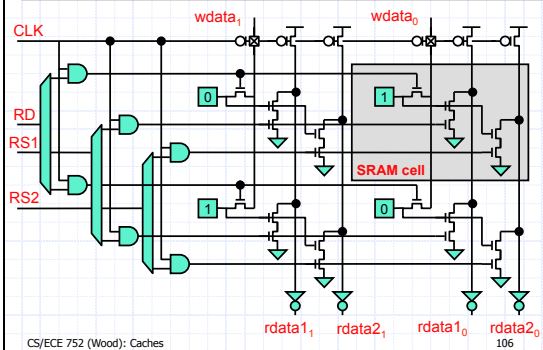
Regfile-Style Write Port



CS/ECE 752 (Wood): Caches

105

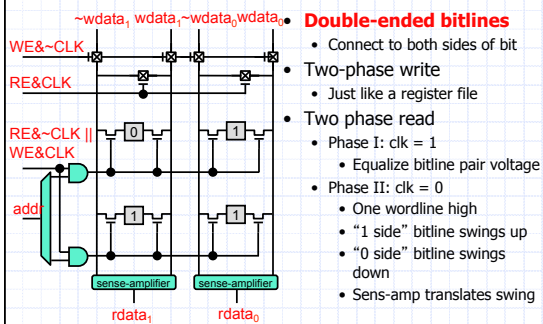
A 2-Read Port 1-Write Port Regfile



CS/ECE 752 (Wood): Caches

106

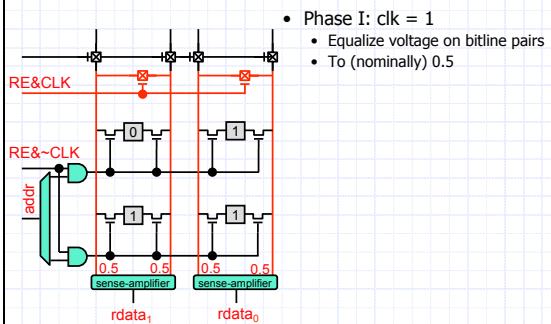
Cache-Style Read/Write Port



CS/ECE 752 (Wood): Caches

107

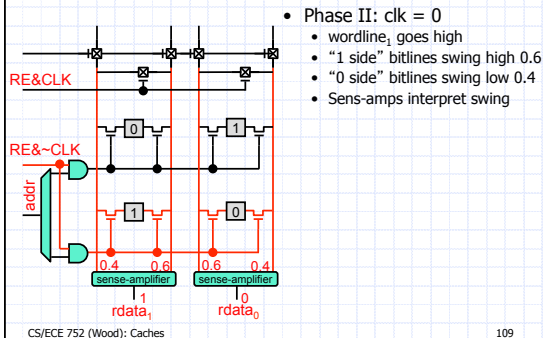
Read/Write Port in Read Action: Phase I



CS/ECE 752 (Wood): Caches

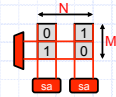
108

Read/Write Port in Read Action: Phase II



Cache-Style SRAM Latency

- Assume
 - M N-bit words
 - Some minimum wire spacing L
 - CCIs occupy no space
- 4 major latency components: taken in series
 - Decoder:** $\propto \log_2 M$
 - Wordlines:** $\propto 2NL$ (cross 2N bitlines)
 - Bitlines:** $\propto ML$ (cross M wordlines)
 - Muxes + sens-amps:** constant
 - 32KB SRAM: red components contribute about equally
- Latency: $\propto (2N+M)L$
 - Make SRAMs as square as possible: minimize $2N+M$
- Latency: $\propto \sqrt{\# \text{bits}}$

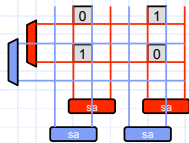


CS/ECE 752 (Wood): Caches

110

Multi-Ported Cache-Style SRAM Latency

- Previous calculation had hidden constant
 - Number of ports **P**
- Recalculate latency components
 - Decoder: $\propto \log_2 M$ (unchanged)
 - Wordlines: $\propto 2NL$ (cross 2N bitlines)
 - Bitlines: $\propto MLP$ (cross MP wordlines)
 - Muxes + sens-amps: constant (unchanged)
- Latency: $\propto (2N+M)L$
- Latency: $\propto \sqrt{\# \text{bits}} * \# \text{ports}$**
- How does latency scale?

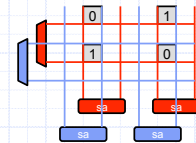


CS/ECE 752 (Wood): Caches

111

Multi-Ported Cache-Style SRAM Power

- Same four components for power
 - $P_{\text{dynamic}} = C * V_{DD}^2 * f$, what is C?
 - Decoder: $\propto \log_2 M$
 - Wordlines: $\propto 2NLP$
 - Huge C per wordline (drives 2N gates)
 - But only one ever high at any time (overall consumption low)
 - Bitlines: $\propto MLP$
 - C lower than wordlines, but large
 - $+ V_{\text{swing}} \ll V_{DD}$ ($C * V_{\text{swing}}^2 * f$)
 - Muxes + **sens-amps:** constant
 - 32KB SRAM: sens-amps are 60–70%
- How does power scale?



CS/ECE 752 (Wood): Caches

112

Multi-Porting an SRAM

- Why multi-porting?
 - Multiple accesses per cycle
- True multi-porting** (physically adding a port) not good
 - Any combination of accesses will work
 - Increases access latency, energy $\propto P$, area $\propto P^2$
- Another option: **pipelining**
 - Timeshare single port on clock edges (wave pipelining: no latches)
 - Negligible area, latency, energy increase
 - Not scalable beyond 2 ports
- Yet another option: **replication**
 - Don't laugh: used for register files, even caches (Alpha 21164)
 - Smaller and faster than true multi-porting $2*P^2 < (2*P)^2$
 - Adds read bandwidth, any combination of reads will work
 - Doesn't add write bandwidth, not really scalable beyond 2 ports

CS/ECE 752 (Wood): Caches

113

Banking an SRAM

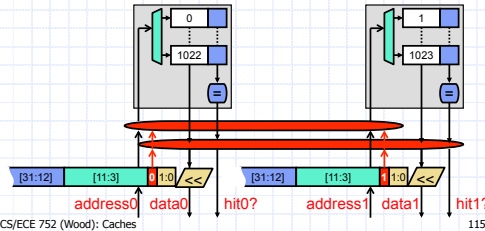
- Still yet another option: **banking (inter-leaving)**
 - Divide SRAM into banks
 - Allow parallel access to different banks
 - Two accesses to same bank? **bank-conflict**, one waits
 - Low area, latency overhead for routing requests to banks
 - Few bank conflicts given sufficient number of banks
 - Rule of thumb: N simultaneous accesses \rightarrow 2N banks
- How to divide words among banks?
 - Round robin:** using address LSB (least significant bits)
 - Example: 16 word RAM divided into 4 banks
 - b0:** 0,4,8,12; **b1:** 1,5,9,13; **b2:** 2,6,10,14; **b3:** 3,7,11,15
 - Why? Spatial locality

CS/ECE 752 (Wood): Caches

114

A Banked Cache

- **Banking** a cache
 - Simple: bank SRAMs
 - Which address bits determine bank? LSB of index
 - **Bank network** assigns accesses to banks, resolves conflicts
 - Adds some latency too



CS/ECE 752 (Wood): Caches

SRAM Summary

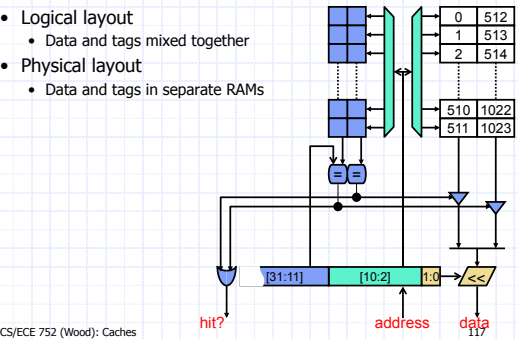
- Large storage arrays are not implemented “digitally”
- SRAM implementation exploits analog transistor properties
 - Inverter pair bits much smaller than latch/flip-flop bits
 - Wordline/bitline arrangement gives simple “grid-like” routing
 - Basic understanding of read, write, read/write ports
 - Wordlines select words
 - Overwhelm inverter-pair to write
 - Drain pre-charged line or swing voltage to read
 - Latency proportional to $\sqrt{\text{#bits} * \text{\#ports}}$

CS/ECE 752 (Wood): Caches

116

Aside: Physical Cache Layout I

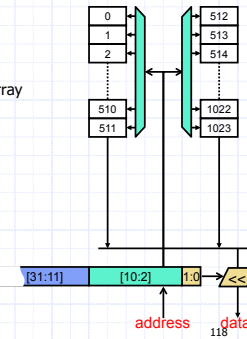
- Logical layout
 - Data and tags mixed together
- Physical layout
 - Data and tags in separate RAMs



CS/ECE 752 (Wood): Caches

Physical Cache Layout II

- Logical layout
 - Data array is monolithic
- Physical layout
 - Each data “way” in separate array

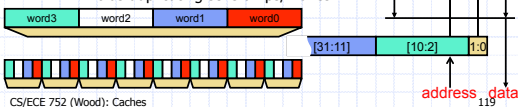


CS/ECE 752 (Wood): Caches

118

Physical Cache Layout III

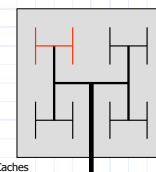
- Logical layout
 - Data blocks are contiguous
- Physical layout
 - Only if full block needed on read
 - E.g., I\$ (read consecutive words)
 - E.g., L2 (read block to fill D\$, I\$)
 - For D\$ (access size is 1 word)...
 - Words in same data blocks are bit-interleaved
 - Word₀.bit₀ adjacent to word₁.bit₀
 - + Builds word selection logic into array
 - + Avoids duplicating sens-amps/muxes



CS/ECE 752 (Wood): Caches

Physical Cache Layout IV

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

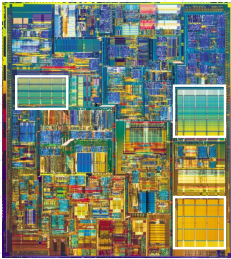


CS/ECE 752 (Wood): Caches

120

Physical Cache Layout

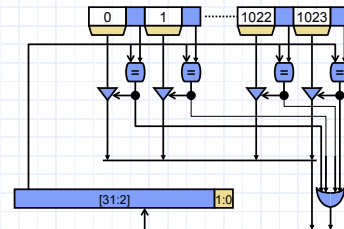
- Arrays and h-trees make caches easy to spot in μ graphs



CS/ECE 752 (Wood): Caches

121

Full-Associativity



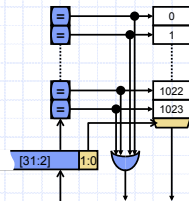
- How to implement full (or at least high) associativity?
 - 1K tag matches? unavoidable, but at least tags are small
 - 1K data reads? Terribly inefficient

CS/ECE 752 (Wood): Caches

122

Full-Associativity with CAMs

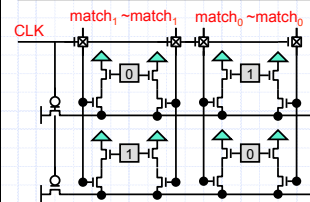
- CAM**: content associative memory
 - Array of words with built-in comparators
 - Matchlines instead of bitlines
 - Output is "one-hot" encoding of match
- FA cache?
 - Tags as CAM
 - Data as RAM
- Hardware is not software**
 - No such thing as software CAM



CS/ECE 752 (Wood): Caches

123

CAM Circuit

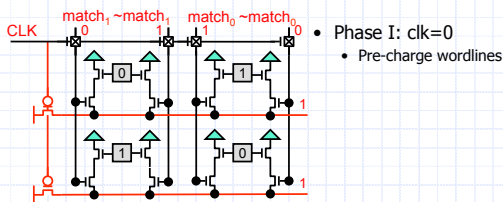


- CAM: reverse RAM
 - Bitlines are inputs
 - Called **matchlines**
 - Wordlines are outputs
- Two phase match
 - Phase I: $\text{clk}=0$
 - Pre-charge wordlines
 - Phase II: $\text{clk}=1$
 - Enable matchlines
 - Non-matching bits discharge wordlines

CS/ECE 752 (Wood): Caches

124

CAM Circuit In Action: Phase I

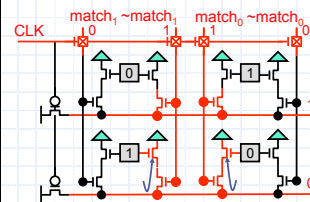


- Phase I: $\text{clk}=0$
 - Pre-charge wordlines

CS/ECE 752 (Wood): Caches

125

CAM Circuit In Action: Phase II



- Phase II: $\text{clk}=1$
 - Enable matchlines
 - Note: bits flipped
 - Non-matching bit discharges wordline
 - ANDs matches
 - NORs non-matches
- Similar technique for doing a fast OR for hit detection

CS/ECE 752 (Wood): Caches

126

CAM Upshot

- CAMs: effective but expensive
 - Matchlines are very expensive (for nasty EE reasons)
- Used but only for 16 or 32 way (max) associativity
- Not for 1024-way associativity
 - No good way of doing something like that
 - + No real need for it, either