# CS 110
# Computer Architecture

# *Caches Part 1*

Instructor:

**Sören Schwertfeger**
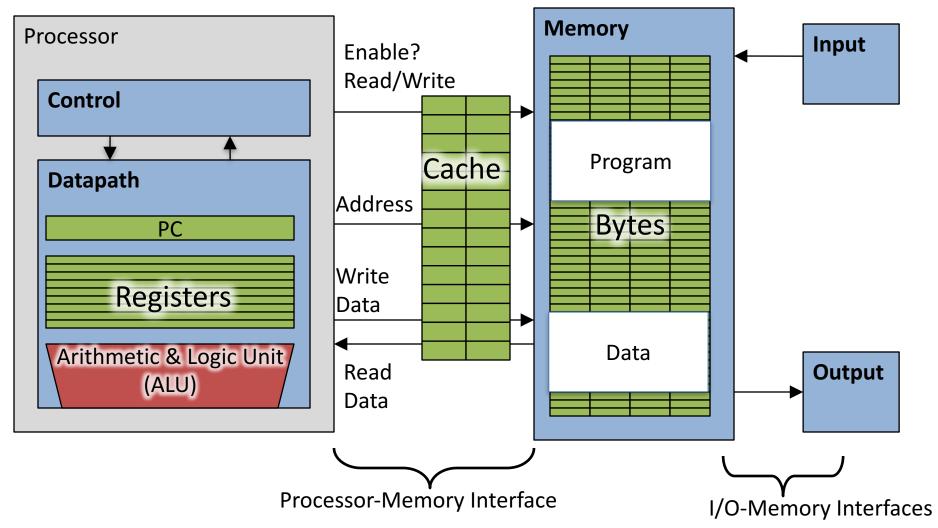
**http://shtech.org/courses/ca/**

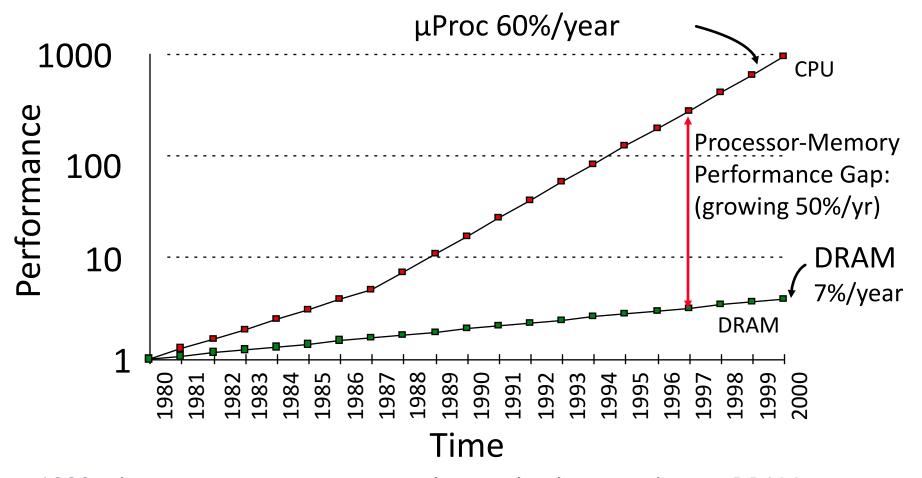**School of Information Science and Technology SIST**

**ShanghaiTech University**

**Slides based on UC Berkley's CS61C**

# Adding Cache to Computer

# Processor-DRAM Gap (latency)



μProc 60%/year

CPU

Processor-Memory
Performance Gap:
(growing 50%/yr)

DRAM
7%/year

DRAM

1000

100

10

1

Performance

1980 1981 1982 1983 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000
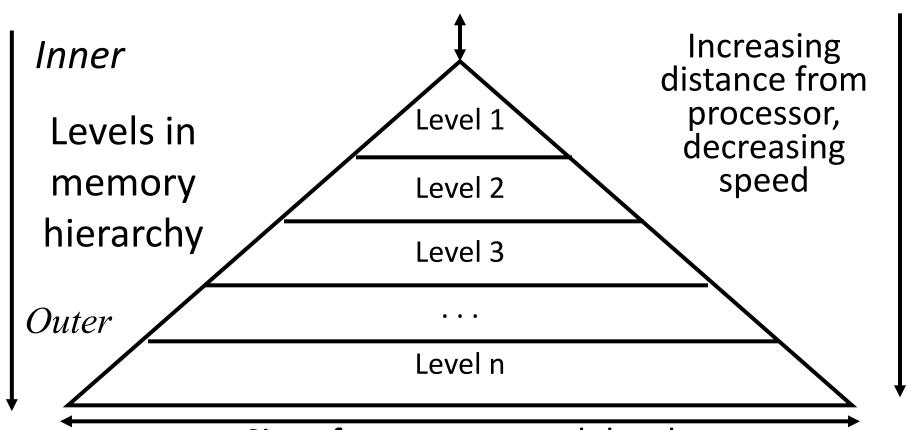
Time

1980 microprocessor executes ~one instruction in same time as DRAM access

2015 microprocessor executes ~1000 instructions in same time as DRAM access

*Slow DRAM access could have disastrous impact on CPU performance!*

3

# Big Idea: Memory Hierarchy

Processor

*Inner*

Levels in memory hierarchy

*Outer*

Increasing distance from processor, decreasing speed
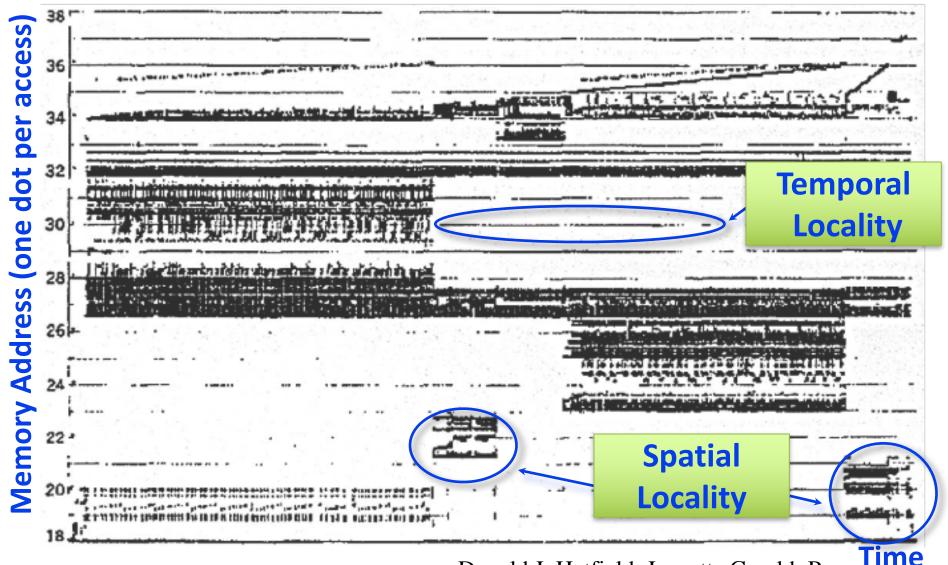
Level 1

Level 2

Level 3

. . .

Level n

Size of memory at each level

*As we move to outer levels the latency goes up and price per bit goes down. Why?*

# Big Idea: Locality

- *Temporal Locality* (locality in time)
  - Go back to same book on desktop multiple times
  - If a memory location is referenced, then it will tend to be referenced again soon

- *Spatial Locality* (locality in space)
  - When go to book shelf, pick up multiple books on J.D. Salinger since library stores related books together
  - If a memory location is referenced, the locations with nearby addresses will tend to be referenced soon

# Memory Reference Patterns
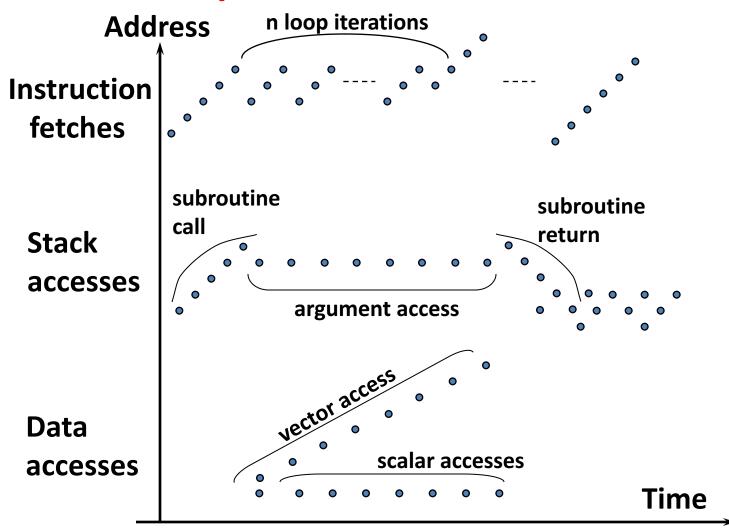


Donald J. Hatfield, Jeanette Gerald: Program Restructuring for Virtual Memory. IBM Systems Journal 10(3): 168-192 (1971)

# Principle of Locality

- *Principle of Locality*: Programs access small portion of address space at any instant of time (spatial locality) and repeatedly access that portion (temporal locality)

- What program structures lead to temporal and spatial locality in **instruction** accesses?

- In **data** accesses?

# Memory Reference Patterns

**Address**

**n loop iterations**

**Instruction fetches**

**Stack accesses**

subroutine call

subroutine return

argument access

**Data accesses**

vector access

scalar accesses

**Time**

# Cache Philosophy

- Programmer-invisible hardware mechanism to give illusion of speed of fastest memory with size of largest memory
  - Works fine even if programmer has no idea what a cache is
  - However, performance-oriented programmers today sometimes "reverse engineer" cache design to design data structures to match cache

# Memory Access with Cache

- Load word instruction: `lw $t0,0($t1)`

- $t1 contains $1022_{ten}$, Memory[1022] = 99

- With cache: Processor issues address $1022_{ten}$ to Cache

  1. Cache checks to see if has copy of data at address $1022_{ten}$

     2a. If finds a match (Hit): cache reads 99, sends to processor

     2b. No match (Miss): cache sends address 1022 to Memory

        I.   Memory reads 99 at address $1022_{ten}$
        II.  Memory sends 99 to Cache
        III. Cache replaces word with new 99
        IV.  Cache sends 99 to processor
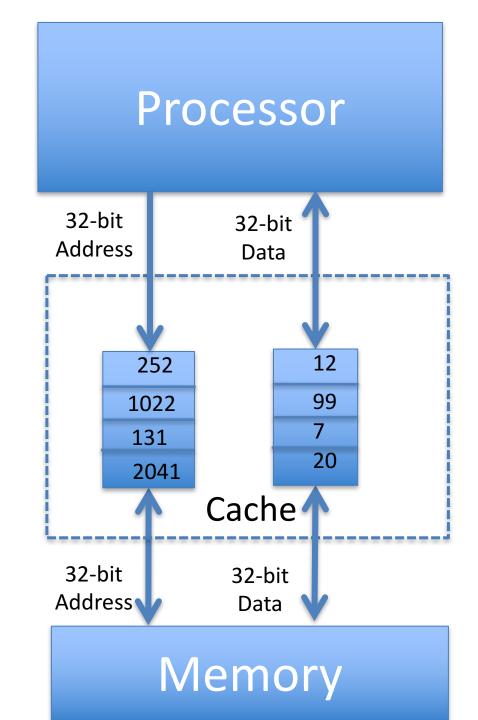
  2. Processor loads 99 into register $t0

# Cache "Tags"

- Need way to tell if have copy of location in memory so that can decide on hit or miss

- On cache miss, put memory address of block in "tag address" of cache block

  1022 placed in tag next to data from memory (99)

| Tag | Data |
|------|------|
| 252 | 12 |
| 1022 | 99 |
| 131 | 7 |
| 2041 | 20 |

From earlier instructions

# Anatomy of a 16 Byte Cache, 4 Byte Block

- Operations:
  1. Cache Hit
  2. Cache Miss
  3. Refill cache from memory

- Cache needs Address Tags to decide if Processor Address is a Cache Hit or Cache Miss
  - Compares all 4 tags

**Processor**

32-bit Address     32-bit Data

| 252 | | 12 |
|-----|---|-----|
| 1022 | | 99 |
| 131 | | 7 |
| 2041 | | 20 |

Cache

32-bit Address     32-bit Data

**Memory**

# Cache Replacement

- Suppose processor now requests location 511, which contains 11?
- Doesn't match any cache block, so must "evict" one resident block to make room
  - Which block to evict?
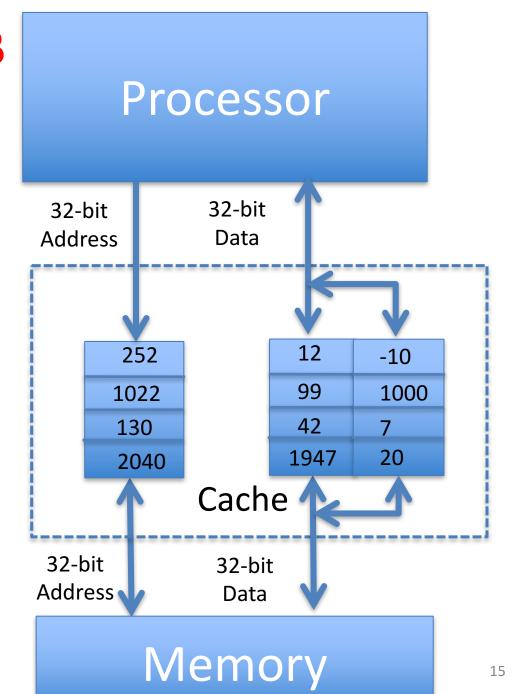- Replace "victim" with new memory block at address 511

| Tag | Data |
|-----|------|
| 252 | 12 |
| 1022 | 99 |
| 511 | 11 |
| 2041 | 20 |

# Block Must be Aligned in Memory

- Word blocks are aligned, so binary address of all words in cache always ends in $00_{two}$

- How to take advantage of this to save hardware and energy?

- Don't need to compare last 2 bits of 32-bit byte address (comparator can be narrower)

=> Don't need to store last 2 bits of 32-bit byte address in Cache Tag (Tag can be narrower)

# Anatomy of a 32B Cache, 8B Block

- Blocks must be aligned in pairs, otherwise could get same word twice in cache

- ➢ Tags only have even-numbered words

- ➢ Last 3 bits of address always $000_{two}$

- ➢ Tags, comparators can be narrower

- Can get hit for either word in block

## Processor

32-bit Address

32-bit Data

### Cache

| | | |
|---|---|---|
| 252 | 12 | -10 |
| 1022 | 99 | 1000 |
| 130 | 42 | 7 |
| 2040 | 1947 | 20 |

32-bit Address

32-bit Data

## Memory

# Hardware Cost of Cache

- Need to compare every tag to the Processor address

- Comparators are expensive

- Optimization: use 2 "sets" => ½ comparators

- 1 Address bit selects which set

- Compare only tags from selected set

- Generalize to more sets

# Processor Address Fields used by Cache Controller

- Block Offset: Byte address within block

- Set Index: Selects which set

- Tag: Remaining portion of processor address

Processor Address (32-bits total)

| Tag | Set Index | Block offset |
|-----|-----------|--------------|

- Size of Index = log2 (number of sets)

- Size of Tag = Address size – Size of Index – log2 (number of bytes/block)

# What is limit to number of sets?

- For a given total number of blocks, we can save more comparators if have more than 2 sets

- Limit: As Many Sets as Cache Blocks => only one block per set – only needs one comparator!

- Called "Direct-Mapped" Design

| Tag | Index | Block offset |
|-----|-------|--------------|

# Direct Mapped Cache Ex: Mapping a 6-bit Memory Address

| 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|
| *Tag* | | *Index* | | *Byte Offset* | |

Mem Block Within $ Block

Block Within $

Byte Within Block

- In example, block size is 4 bytes/1 word
- Memory and cache blocks always the same size, unit of transfer between memory and cache
- # Memory blocks >> # Cache blocks
  - 16 Memory blocks = 16 words = 64 bytes => 6 bits to address all bytes
  - 4 Cache blocks, 4 bytes (1 word) per block
  - 4 Memory blocks map to each cache block
- Memory block to cache block, aka *index*: middle two bits
- Which memory block is in a given cache block, aka *tag*: top two bits

# One More Detail: Valid Bit

- When start a new program, cache does not have valid information for this program

- Need an indicator whether this tag entry is valid for this program

- Add a "valid bit" to the cache tag entry

  0 => cache miss, even if by chance, address = tag

  1 => cache hit, if processor address = tag

# Caching:  A Simple First Example

**Cache**

Index   Valid   Tag        Data

00

01

10

11

Q: Is the memory block in cache?

Compare the cache tag to the high-order 2 memory address bits to tell if the memory block is in the cache (provided valid bit is set)

**Main Memory**

0000xx
0001xx
0010xx
0011xx
0100xx
0101xx
0110xx
0111xx
1000xx
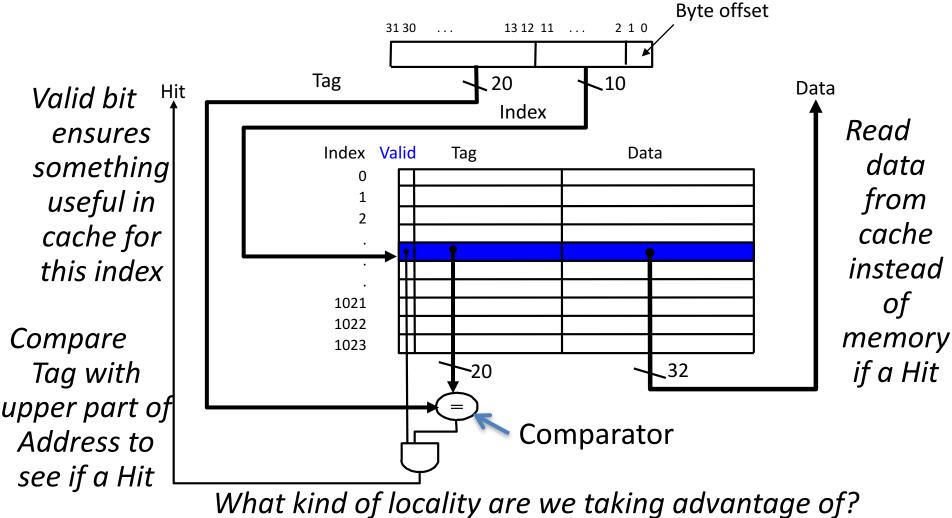1001xx
1010xx
1011xx
1100xx
1101xx
1110xx
1111xx

One word blocks
Two low order bits (xx) define the byte in the block (32b words)

Q: Where in the cache is the mem block?

Use next 2 low-order memory address bits – the index – to determine which cache block (i.e., modulo the number of blocks in the cache)

# Direct-Mapped Cache Example

- One word blocks, cache size = 1K words (or 4KB)

Byte offset

31 30   . . .   13 12 11   . . .   2 1 0

Tag

Hit

*Valid bit ensures something useful in cache for this index*

*Compare Tag with upper part of Address to see if a Hit*

20

10

Index

Data

*Read data from cache instead of memory if a Hit*

| Index | Valid | Tag | Data |
|-------|-------|-----|------|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| . | | | |
| . | | | |
| . | | | |
| 1021 | | | |
| 1022 | | | |
| 1023 | | | |

20

32

=

Comparator

*What kind of locality are we taking advantage of?*

# Multiword-Block Direct-Mapped Cache

- Four words/block, cache size = 1K words

Hit

Byte offset

Data

31 30 . . . 13 12 11 . . . 4 3 2 1 0

Tag

20

8

2

Word offset

Index

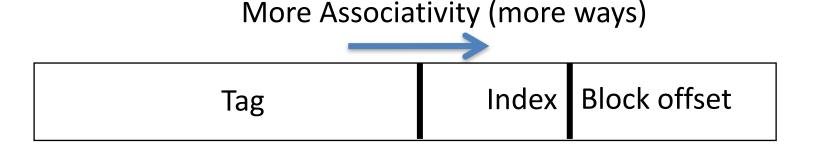Index  Valid  Tag

Data

0
1
2
.
.
.
253
254
255

20

=

32

*What kind of locality are we taking advantage of?*

# Cache Names for Each Organization

- "Fully Associative": Line can go anywhere
  - First design in lecture
  - Note: No Index field, but 1 comparator/ line
- "Direct Mapped": Line goes one place
  - Note: Only 1 comparator
  - Number of sets = number blocks
- "N-way Set Associative": N places for a line
  - Number of sets = number of lines/ N
  - N comparators
  - *Fully Associative: N = number of lines*
  - *Direct Mapped: N = 1*
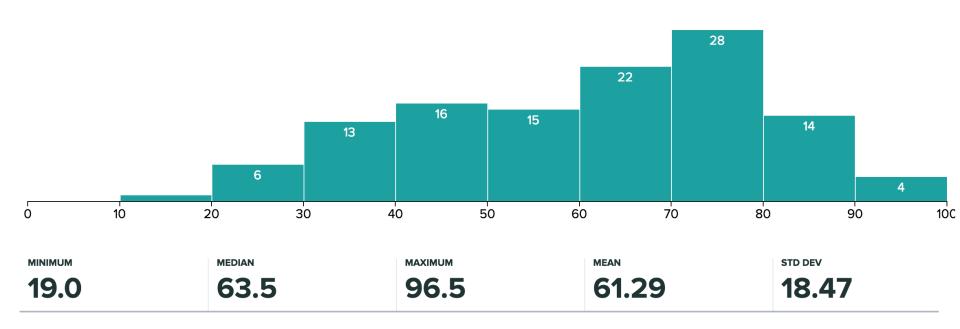
# Range of Set-Associative Caches

- For a fixed-size cache, and a given block size, each increase by a factor of 2 in associativity doubles the number of blocks per set (i.e., the number of "ways") and halves the number of sets –
  - decreases the size of the index by 1 bit and increases the size of the tag by 1 bit

More Associativity (more ways)

| Tag | Index | Block offset |
|-----|-------|--------------|

# Question

- For a cache with constant total capacity, if we increase the number of ways by a factor of 2, which statement is false:

- A: The number of sets could be doubled

- B: The tag width could decrease

- C: The block size could stay the same

- D: The block size could be halved

- E: Tag width must increase

# Midterm I



| MINIMUM | MEDIAN | MAXIMUM | MEAN | STD DEV |
|---------|--------|---------|------|---------|
| **19.0** | **63.5** | **96.5** | **61.29** | **18.47** |

# Total Cash Capacity =

## Associativity  *  # of sets  *  block_size

*Bytes = blocks/set  *  sets  *  Bytes/block*

*C = N  *  S  *  B*

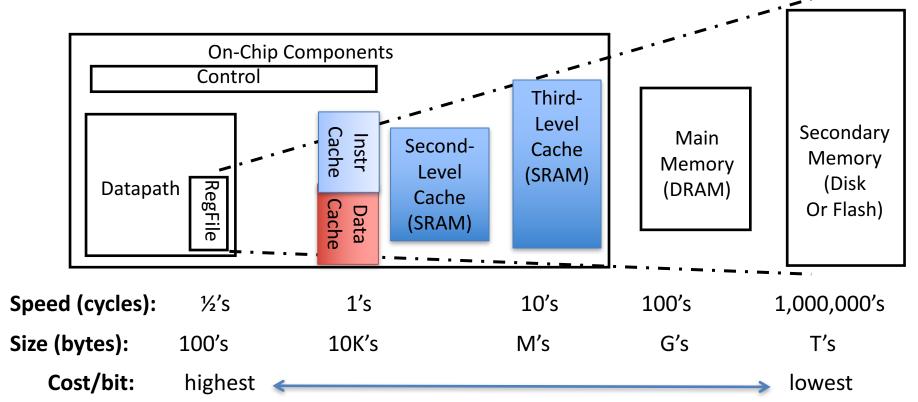| Tag | Index | Byte Offset |
|-----|-------|-------------|

address_size = tag_size + index_size + offset_size
= tag_size + log2(S) + log2(B)

Clicker Question:  C remains constant, S and/or B can change such that
C = 2N * (SB)' => (SB)' = SB/2
Tag_size = address_size – (log2(S) + log2(B)) = address_size – log2(SB)

# Typical Memory Hierarchy



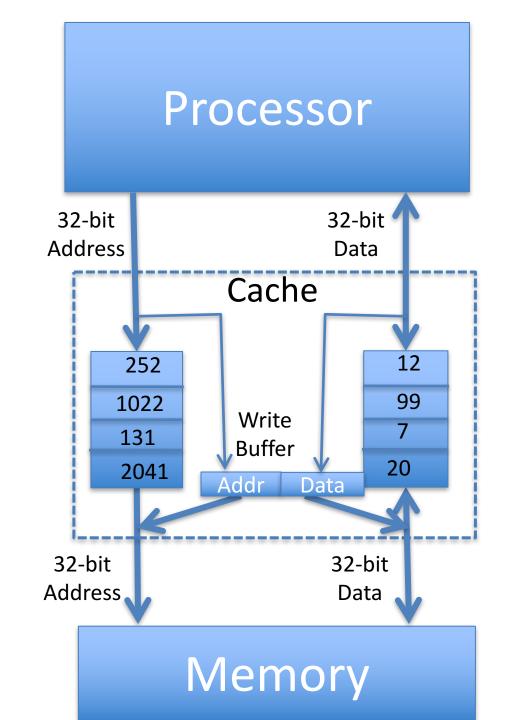| Speed (cycles): | ½'s | 1's | 10's | 100's | 1,000,000's |
|---|---|---|---|---|---|
| Size (bytes): | 100's | 10K's | M's | G's | T's |
| Cost/bit: | highest | | ← | → | lowest |

- **Principle of locality + memory hierarchy** presents programmer with ≈ as much memory as is available in the *cheapest* technology at the ≈ speed offered by the *fastest* technology

# Handling Stores with Write-Through

- Store instructions write to memory, changing values

- Need to make sure cache and memory have same values on writes: 2 policies

1) Write-Through Policy: write cache and write *through* the cache to memory

  – Every write eventually gets to memory

  – Too slow, so include Write Buffer to allow processor to continue once data in Buffer

  – Buffer updates memory in parallel to processor

# Write-Through Cache

- Write both values in cache and in memory
- Write buffer stops CPU from stalling if memory cannot keep up
- Write buffer may have multiple entries to absorb bursts of writes
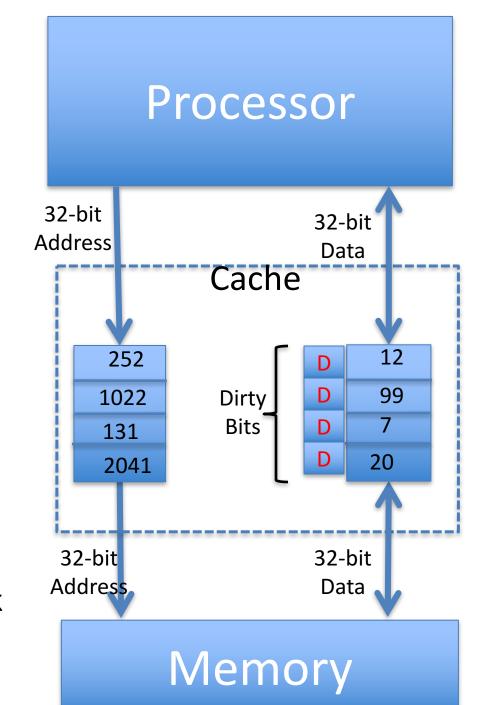- What if store misses in cache?

# Handling Stores with Write-Back

2) Write-Back Policy: write only to cache and then write cache block *back* to memory when evict block from cache

– Writes collected in cache, only single write to memory per block

– Include bit to see if wrote to block or not, and then only write back if bit is set

• Called "Dirty" bit (writing makes it "dirty")

# Write-Back Cache

- Store/cache hit, write data in cache *only* & set dirty bit
  - Memory has stale value
- Store/cache miss, read data from memory, then update and set dirty bit
  - "Write-allocate" policy
- Load/cache hit, use value from cache
- On any miss, write back evicted block, only if dirty. Update cache with new block and clear dirty bit.

**Processor**

32-bit Address

32-bit Data

Cache

| | |
|---|---|
| 252 | |
| 1022 | |
| 131 | |
| 2041 | |

Dirty Bits

| D | 12 |
|---|---|
| D | 99 |
| D | 7 |
| D | 20 |

32-bit Address

32-bit Data

**Memory**

# Write-Through vs. Write-Back

- Write-Through:
  - Simpler control logic
  - More predictable timing simplifies processor control logic
  - Easier to make reliable, since memory always has copy of data (big idea: Redundancy!)

- Write-Back
  - More complex control logic
  - More variable timing (0,1,2 memory accesses per cache access)
  - Usually reduces write traffic
  - Harder to make reliable, sometimes cache has only copy of data

# Cache (*Performance)* Terms

- Hit rate: fraction of accesses that hit in the cache

- Miss rate: 1 − Hit rate

- Miss penalty: time to replace a line/ block from lower level in memory hierarchy to cache

- Hit time: time to access cache memory (including tag comparison)

- Abbreviation: "$" = cache ( cash … )

# Average Memory Access Time (AMAT)

- Average Memory Access Time (AMAT) is the average time to access memory considering both hits and misses in the cache

AMAT =    Time for a hit
                            +  Miss rate × Miss penalty

# Question
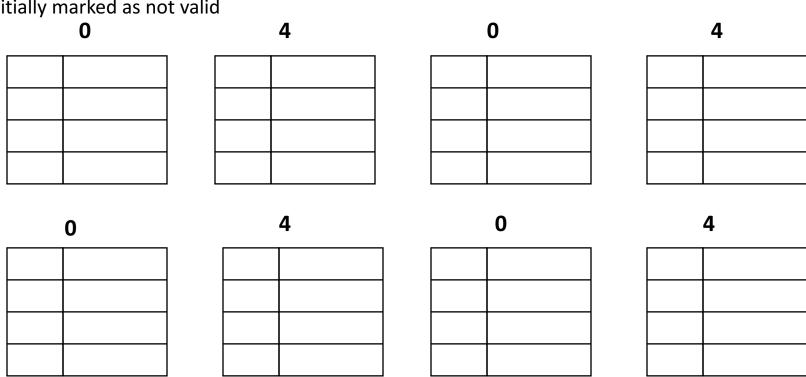
AMAT =  Time for a hit  +  Miss rate x Miss penalty

Given a 200 psec clock, a miss penalty of 50 clock cycles, a miss rate of 0.02 misses per instruction and a cache hit time of 1 clock cycle, what is AMAT?

☐ A:  ≤200 psec

☐ B:  400 psec

☐ C:  600 psec

☐ D:  ≥ 800 psec

# Example: Direct-Mapped Cache
## with 4 Single-Word Blocks, Worst-Case Reference String

- Consider the main memory address (words) reference string of word numbers:        0   4   0   4   0   4   0   4

Start with an empty cache - all blocks
initially marked as not valid

**0**           **4**           **0**           **4**

**0**           **4**           **0**           **4**

# Example: Direct-Mapped Cache
## with 4 Single-Word Blocks, Worst-Case Reference String

- Consider the main memory address (words) reference string of word numbers:  0  4  0  4  0  4  0  4
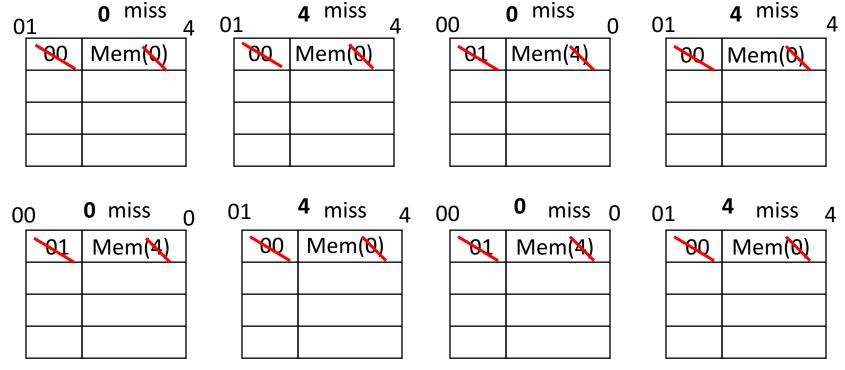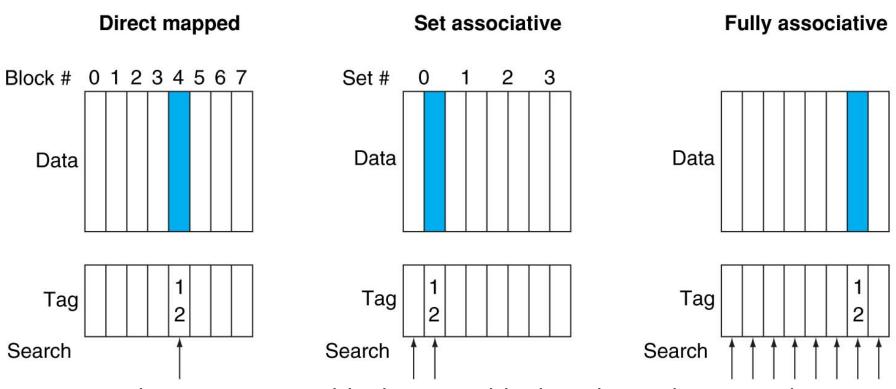
Start with an empty cache - all blocks
initially marked as not valid

| 01 | **0** miss | 4 |
|----|----|----|
| 00 | Mem(0) | |
| | | |
| | | |
| | | |

| 01 | **4** miss | 4 |
|----|----|----|
| 00 | Mem(0) | |
| | | |
| | | |
| | | |

| 00 | **0** miss | 0 |
|----|----|----|
| 01 | Mem(4) | |
| | | |
| | | |
| | | |

| 01 | **4** miss | 4 |
|----|----|----|
| 00 | Mem(0) | |
| | | |
| | | |
| | | |

| 00 | **0** miss | 0 |
|----|----|----|
| 01 | Mem(4) | |
| | | |
| | | |
| | | |

| 01 | **4** miss | 4 |
|----|----|----|
| 00 | Mem(0) | |
| | | |
| | | |
| | | |

| 00 | **0** miss | 0 |
|----|----|----|
| 01 | Mem(4) | |
| | | |
| | | |
| | | |

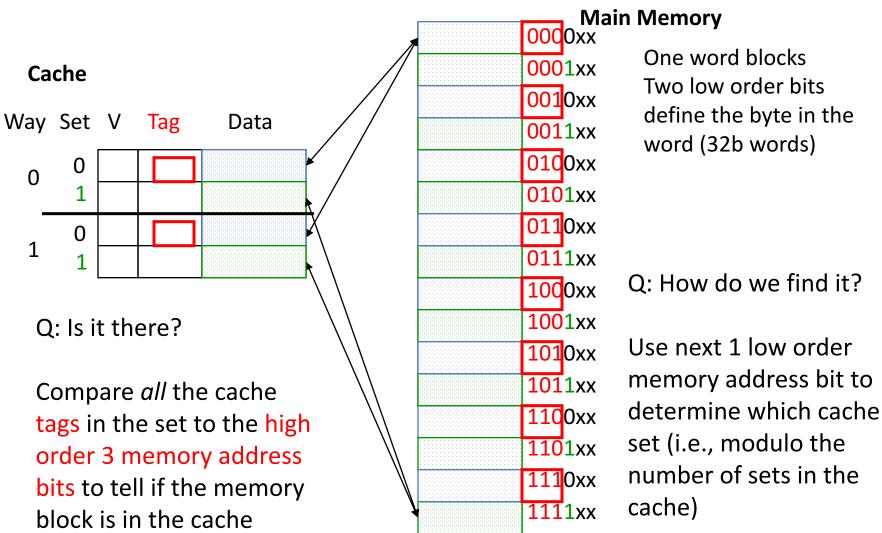| 01 | **4** miss | 4 |
|----|----|----|
| 00 | Mem(0) | |
| | | |
| | | |
| | | |

- 8 requests, 8 misses
- Ping-pong effect due to conflict misses - two memory locations that map into the same cache block

# Alternative Block Placement Schemes

**Direct mapped**  **Set associative**  **Fully associative**

Block #  0 1 2 3 4 5 6 7    Set #  0  1  2  3

Data

Tag  1 2

Search

- DM placement: mem block 12 in 8 block cache: only one cache block where mem block 12 can be found—(12 modulo 8) = 4
- SA placement: four sets x 2-ways (8 cache blocks), memory block 12 in set (12 mod 4) = 0; either element of the set
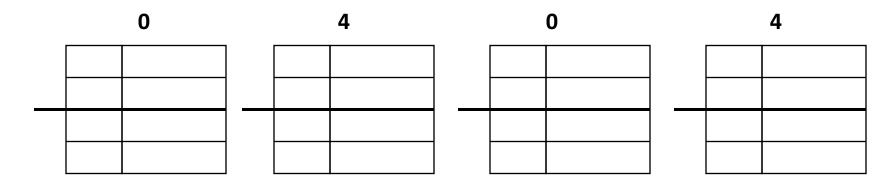- FA placement: mem block 12 can appear in any cache blocks

# Example: 2-Way Set Associative $
## (4 words = 2 sets x 2 ways per set)

**Cache**

| Way | Set | V | Tag | Data |
|-----|-----|---|-----|------|
| 0 | 0 | | ▭ | |
| | 1 | | | |
| 1 | 0 | | ▭ | |
| | 1 | | | |

Q: Is it there?

Compare *all* the cache tags in the set to the high order 3 memory address bits to tell if the memory block is in the cache

**Main Memory**

| | |
|---|---|
| 000 | 0xx |
| 0001 | xx |
| 001 | 0xx |
| 0011 | xx |
| 010 | 0xx |
| 0101 | xx |
| 011 | 0xx |
| 0111 | xx |
| 100 | 0xx |
| 1001 | xx |
| 101 | 0xx |
| 1011 | xx |
| 110 | 0xx |
| 1101 | xx |
| 111 | 0xx |
| 1111 | xx |

One word blocks
Two low order bits define the byte in the word (32b words)

Q: How do we find it?

Use next 1 low order memory address bit to determine which cache set (i.e., modulo the number of sets in the cache)

41

# Example: 4-Word 2-Way SA $ Same Reference String

- Consider the main memory address (word) reference string

Start with an empty cache - all blocks initially marked as not valid

0  4  0  4  0  4  0  4

# Example: 4-Word 2-Way SA $ Same Reference String

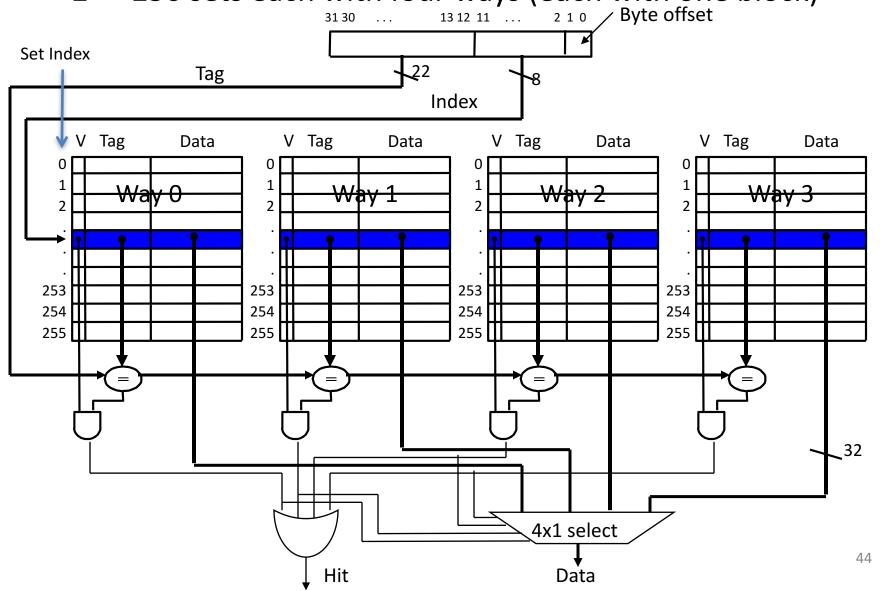- Consider the main memory address (word) reference string

Start with an empty cache - all blocks initially marked as not valid

0 4 0 4 0 4 0 4

**0** miss

| 000 | Mem(0) |
|-----|--------|
|     |        |
|     |        |
|     |        |

**4** miss

| 000 | Mem(0) |
|-----|--------|
|     |        |
| 010 | Mem(4) |
|     |        |

**0** hit

| 000 | Mem(0) |
|-----|--------|
|     |        |
| 010 | Mem(4) |
|     |        |

**4** hit

| 000 | Mem(0) |
|-----|--------|
|     |        |
| 010 | Mem(4) |
|     |        |

- 8 requests, 2 misses

- Solves the ping-pong effect in a direct-mapped cache due to conflict misses since now two memory locations that map into the same cache set can co-exist!

# Four-Way Set-Associative Cache

- $2^8$ = 256 sets each with four ways (each with one block)

# Different Organizations of an Eight-Block Cache

**One-way set associative (direct mapped)**

| Block | Tag | Data |
|-------|-----|------|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |

**Two-way set associative**

| Set | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

Total size of $ in blocks is equal to *number of sets × associativity*. For fixed $ size and fixed block size, increasing associativity decreases number of sets while increasing number of elements per set. With eight blocks, an 8-way set-associative $ is same as a fully associative $.
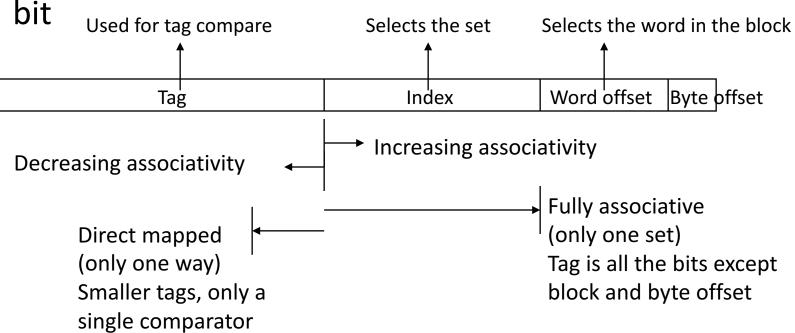
**Four-way set associative**

| Set | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|-----|------|-----|------|
| 0 | | | | | | | | |
| 1 | | | | | | | | |

**Eight-way set associative (fully associative)**

| Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|
| | | | | | | | | | | | | | | | |

# Range of Set-Associative Caches

- For a *fixed-size* cache and fixed block size, each increase by a factor of two in associativity doubles the number of blocks per set (i.e., the number or ways) and halves the number of sets – decreases the size of the index by 1 bit and increases the size of the tag by 1 bit
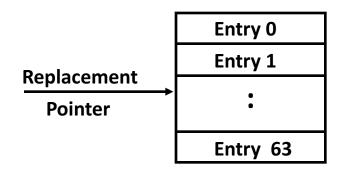
Used for tag compare          Selects the set          Selects the word in the block

| Tag | Index | Word offset | Byte offset |

Decreasing associativity

Increasing associativity

Direct mapped
(only one way)
Smaller tags, only a
single comparator

Fully associative
(only one set)
Tag is all the bits except
block and byte offset

# Costs of Set-Associative Caches

- N-way set-associative cache costs
  - N comparators (delay and area)
  - MUX delay (set selection) before data is available
  - Data available after set selection (and Hit/Miss decision). DM $: block is available before the Hit/Miss decision
    - In Set-Associative, not possible to just assume a hit and continue and recover later if it was a miss
- When miss occurs, which way's block selected for replacement?
  - Least Recently Used (LRU): one that has been unused the longest (principle of temporal locality)
    - Must track when each way's block was used relative to other blocks in the set
    - For 2-way SA $, one bit per set → set to 1 when a block is referenced; reset the other way's bit (i.e., "last used")
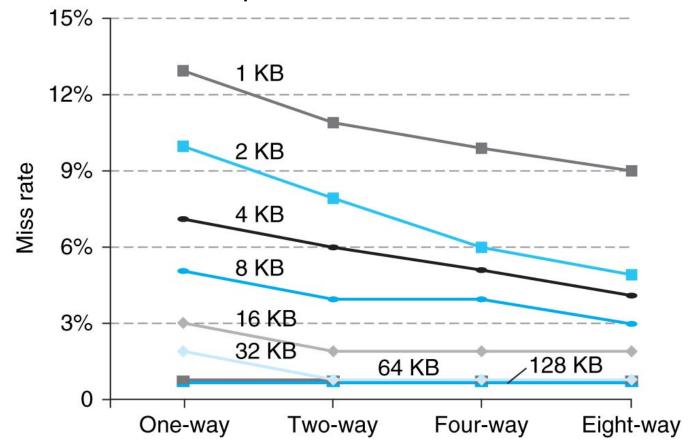
# Cache Replacement Policies

- Random Replacement
  - Hardware randomly selects a cache evict
- Least-Recently Used
  - Hardware keeps track of access history
  - Replace the entry that has not been used for the longest time
  - For 2-way set-associative cache, need one bit for LRU replacement
- Example of a Simple "Pseudo" LRU Implementation
  - Assume 64 Fully Associative entries
  - Hardware replacement pointer points to one cache entry
  - Whenever access is made to the entry the pointer points to:
    - Move the pointer to the next entry
  - Otherwise: do not move the pointer
  - (example of "not-most-recently used" replacement policy)

**Replacement**

**Pointer**

| Entry 0 |
| Entry 1 |
| : |
| Entry 63 |

# Benefits of Set-Associative Caches

- Choice of DM $ versus SA $ depends on the cost of a miss versus the cost of implementation



- Largest gains are in going from direct mapped to 2-way (20%+ reduction in miss rate)

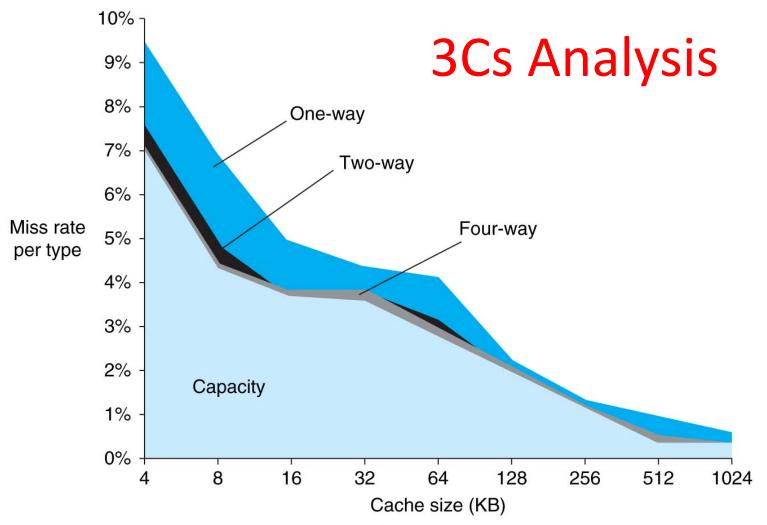# Understanding Cache Misses: The 3Cs

- **Compulsory** (cold start or process migration, 1$^{st}$ reference):
  - First access to block impossible to avoid; small effect for long running programs
  - Solution: increase block size (increases miss penalty; very large blocks could increase miss rate)
- **Capacity**:
  - Cache cannot contain all blocks accessed by the program
  - Solution: increase cache size (may increase access time)
- *Conflict (collision):*
  - *Multiple memory locations mapped to the same cache location*
  - *Solution 1: increase cache size*
  - *Solution 2: increase associativity (may increase access time)*

# How to Calculate 3C's using Cache Simulator

1. *Compulsory*: set cache size to infinity and fully associative, and count number of misses

2. *Capacity*: Change cache size from infinity, usually in powers of 2, and count misses for each reduction in size
   – 16 MB, 8 MB, 4 MB, … 128 KB, 64 KB, 16 KB

3. *Conflict*: Change from fully associative to n-way set associative while counting misses
   – Fully associative, 16-way, 8-way, 4-way, 2-way, 1-way

# 3Cs Analysis



- Three sources of misses (SPEC2000 integer and floating-point benchmarks)
    - Compulsory misses 0.006%; not visible
    - Capacity misses, function of cache size
    - Conflict portion depends on associativity and cache size

# Improving Cache Performance

AMAT = Time for a hit + Miss rate x Miss penalty

- Reduce the time to hit in the cache
  - E.g., Smaller cache
- Reduce the miss rate
  - E.g., Bigger cache
- Reduce the miss penalty
  - E.g., Use multiple cache levels

# Impact of Larger Cache on AMAT?

- 1) Reduces misses (what kind(s)?)
- 2) Longer Access time (Hit time): smaller is faster
  - Increase in hit time will likely add another stage to the pipeline
- At some point, increase in hit time for a larger cache may overcome the improvement in hit rate, yielding a decrease in performance
- Computer architects expend considerable effort optimizing organization of cache hierarchy – big impact on performance and power!

# Questions: Impact of longer cache blocks on misses?

- For fixed total cache capacity and associativity, what is effect of longer blocks on each type of miss:

  - A: Decrease, B: Unchanged, C: Increase

- Compulsory?

- Capacity?

- Conflict?

# Questions: Impact of longer blocks on AMAT

- For fixed total cache capacity and associativity, what is effect of longer blocks on each component of AMAT:

  - A: Decrease, B: Unchanged, C: Increase

- Hit Time?

- Miss Rate?

- Miss Penalty?

# Question:

For fixed capacity and fixed block size, how does increasing associativity effect AMAT?
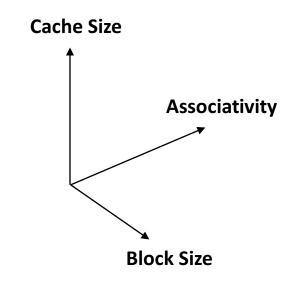
A: Increases hit time, decreases miss rate
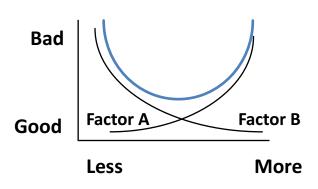B: Decreases hit time, decreases miss rate
C: Increases hit time, increases miss rate
D: Decreases hit time, increases miss rate

# Cache Design Space

- Several interacting dimensions
  - Cache size
  - Block size
  - Associativity
  - Replacement policy
  - Write-through vs. write-back
  - Write allocation
- Optimal choice is a compromise
  - Depends on access characteristics
    - Workload
    - Use (I-cache, D-cache)
  - Depends on technology / cost
- Simplicity often wins

**Cache Size**

**Associativity**

**Block Size**

**Bad**

**Good**

Factor A    Factor B

**Less**    **More**

# And In Conclusion, …

- Principle of Locality for Libraries /Computer Memory
- Hierarchy of Memories (speed/size/cost per bit) to Exploit Locality
- Cache – copy of data lower level in memory hierarchy
- Direct Mapped to find block in cache using Tag field and Valid bit for Hit
- Cache design choice:
    - Write-Through vs. Write-Back