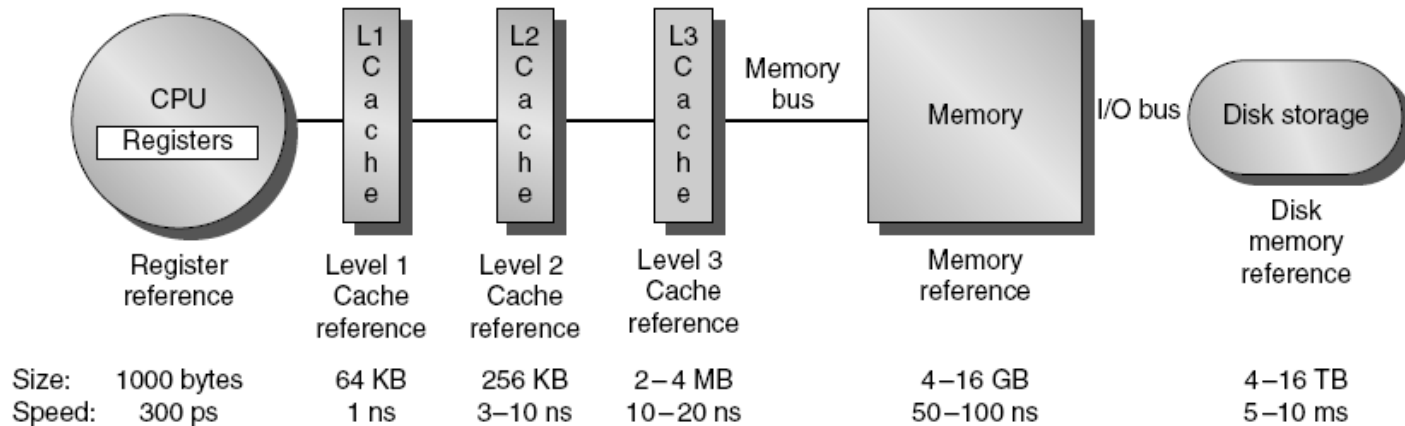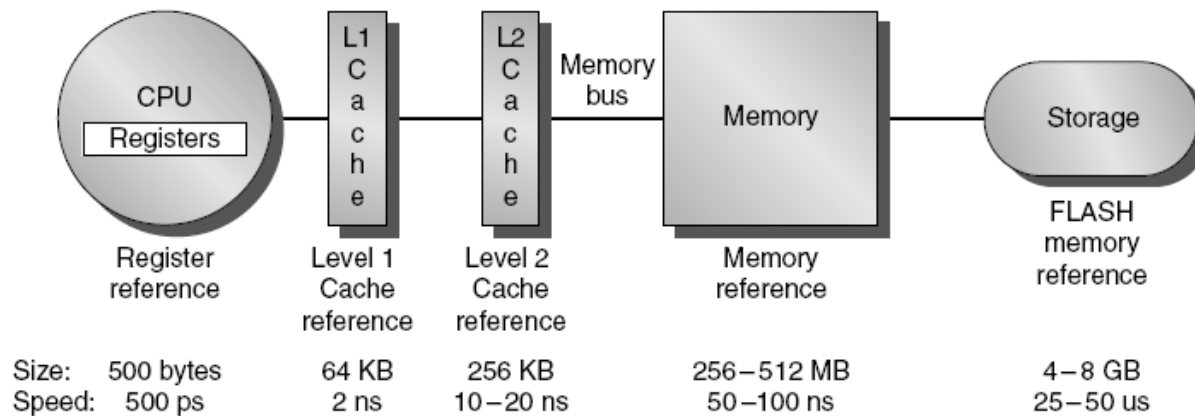# Chapter 2

# Memory Hierarchy Design

# Introduction

- Programmers want unlimited amounts of memory with low latency
- Fast memory technology is more expensive per bit than slower memory
- Solution:  organize memory system into a hierarchy
  - Entire addressable memory space available in largest, slowest memory
  - Incrementally smaller and faster memories, each containing a subset of the memory below it, proceed in steps up toward the processor
- Temporal and spatial locality insures that nearly all references can be found in smaller memories
  - Gives the allusion of a large, fast memory being presented to the processor
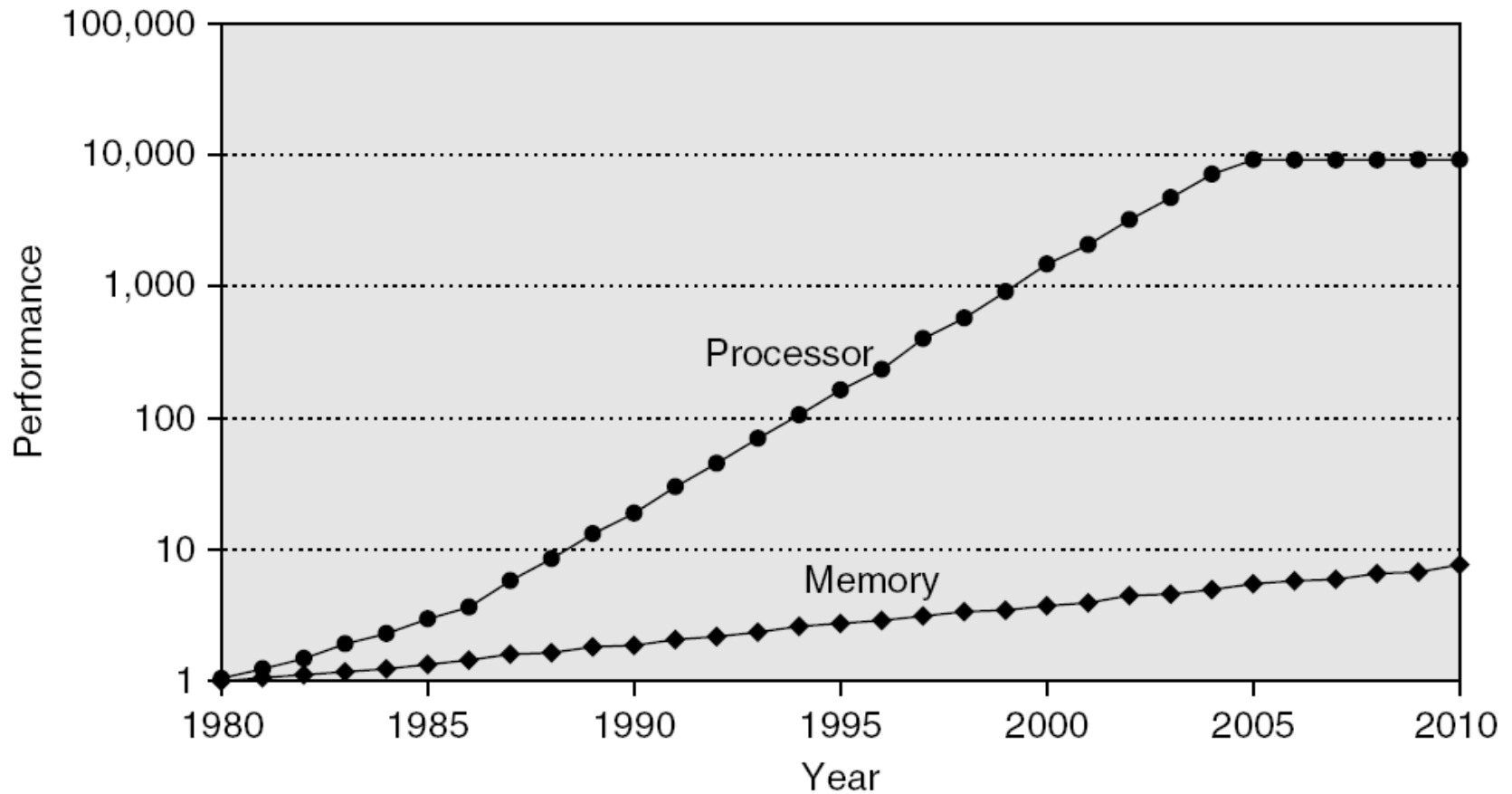
2

# Memory Hierarchy

(a) Memory hierarchy for server

| | Register reference | Level 1 Cache reference | Level 2 Cache reference | Level 3 Cache reference | Memory reference | Disk memory reference |
|---|---|---|---|---|---|---|
| Size: | 1000 bytes | 64 KB | 256 KB | 2–4 MB | 4–16 GB | 4–16 TB |
| Speed: | 300 ps | 1 ns | 3–10 ns | 10–20 ns | 50–100 ns | 5–10 ms |

(b) Memory hierarchy for a personal mobile device

| | Register reference | Level 1 Cache reference | Level 2 Cache reference | Memory reference | FLASH memory reference |
|---|---|---|---|---|---|
| Size: | 500 bytes | 64 KB | 256 KB | 256–512 MB | 4–8 GB |
| Speed: | 500 ps | 2 ns | 10–20 ns | 50–100 ns | 25–50 us |

3

# Memory Performance Gap
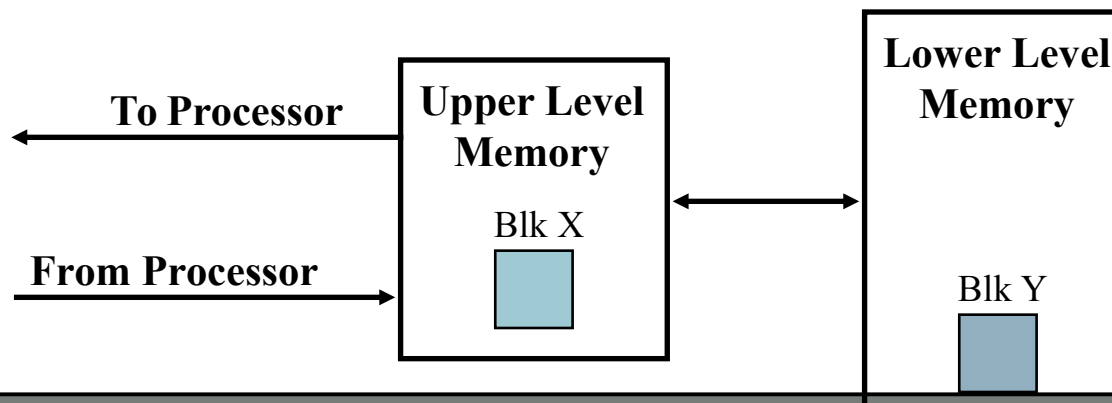
# Memory Hierarchy Design

- Memory hierarchy design becomes more crucial with recent multi-core processors:
  - Aggregate peak bandwidth grows with # cores:
    - Intel Core i7 can generate two references per core per clock
    - Four cores and 3.2 GHz clock
      - 25.6 billion 64-bit data references/second +
      - 12.8 billion 128-bit instruction references
      - = 409.6 GB/s!
  - DRAM bandwidth is only 6% of this (25 GB/s)
  - Requires:
    - Multi-port, pipelined caches
    - Two levels of cache per core
    - Shared third-level cache on chip

# Performance and Power

- High-end microprocessors have >10 MB on-chip cache
  - Consumes large amount of area and power budget

6

# Memory Hierarchy: Terminology

- Hit: data appears in some block in the upper level (example: Block X)
  - Hit Rate: the fraction of memory access found in the upper level
  - Hit Time: Time to access the upper level which consists of
    RAM access time + Time to determine hit/miss
- Miss: data needs to be retrieve from a block in the lower level (Block Y)
  - Miss Rate = 1 - (Hit Rate)
  - Miss Penalty: Time to replace a block in the upper level +
    Time to deliver the block the processor
- Hit Time << Miss Penalty (500 instructions on 21264!)

To Processor

From Processor

**Upper Level Memory**

Blk X

**Lower Level Memory**

Blk Y

# Cache Measures

- *Hit rate*: fraction found in that level
  - So high that usually talk about *Miss rate*
  - Miss rate fallacy: as MIPS to CPU performance,
    miss rate to average memory access time in memory
- Average memory-access time
      = Hit time + Miss rate x Miss penalty
      (ns or clocks)
- *Miss penalty*: time to replace a block from lower level, including time
  to replace in CPU
  - *access time*: time to lower level

    = f(latency to lower level)
  - *transfer time*: time to transfer block

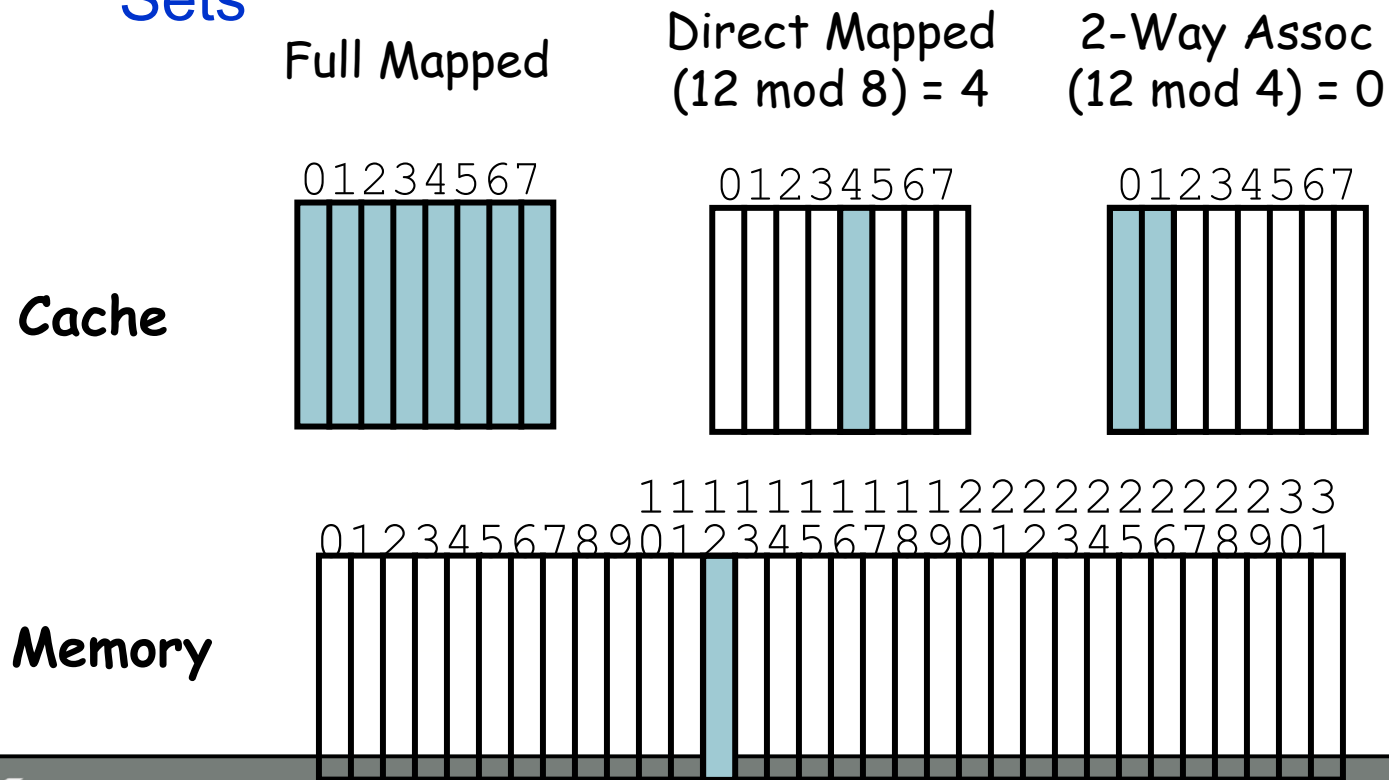    =f(BW between upper & lower levels)

# 4 Questions for Memory Hierarchy

- Q1: Where can a block be placed in the upper level?     *(Block placement)*

- Q2: How is a block found if it is in the upper level?     *(Block identification)*

- Q3: Which block should be replaced on a miss?     *(Block replacement)*

- Q4: What happens on a write?     *(Write strategy)*

# Q1: Where can a block be placed in the upper level?

- Block 12 placed in 8 block cache:
    - Fully associative, direct mapped, 2-way set associative
    - S.A. Mapping = Block Number Modulo Number Sets

Full Mapped

Direct Mapped
(12 mod 8) = 4

2-Way Assoc
(12 mod 4) = 0

01234567

01234567

01234567

**Cache**

**Memory**

1111111111222222222233
01234567890123456789012345678901

# Q2: How is a block found if it is in the upper level?

- Tag on each block
  - No need to check index or block offset
- Increasing associativity shrinks index, expands tag

| Block Address | | Block Offset |
|---|---|---|
| Tag | Index | |

# **Example**

- Suppose we have a 16KB of data in a direct-mapped cache with 4 word blocks

- Determine the size of the tag, index and offset fields if we're using a 32-bit architecture

- Offset

  - need to specify correct byte within a block

  - block contains       4 words
                        16 bytes
                        $2^4$ bytes

  - need 4 bits to specify correct byte

# Example [contd…]

- Index: (~index into an "array of blocks")
  - need to specify correct row in cache
  - cache contains 16 KB = $2^{14}$ bytes
  - block contains $2^4$ bytes (4 words)
  - # rows/cache = # blocks/cache (since there's one block/row)

$$= \frac{\text{bytes/cache}}{\text{bytes/row}} = \frac{2^{14} \text{ bytes/cache}}{2^4 \text{ bytes/row}}$$

$$= 2^{10} \text{ rows/cache}$$

  - need 10 bits to specify this many rows

# Example [contd…]

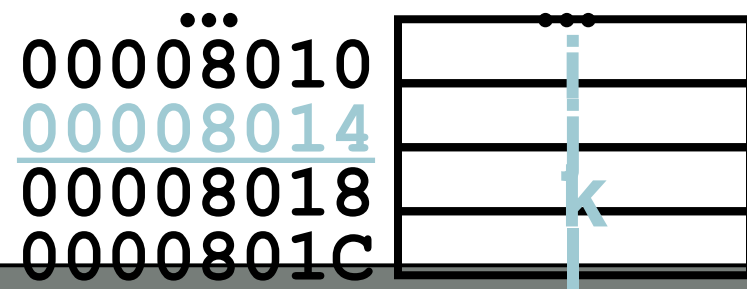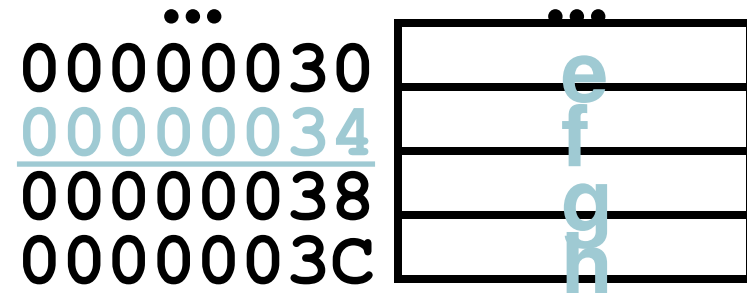- Tag: use remaining bits as tag

  – tag length = mem addr length

              - offset

              - index

           = 32 - 4 - 10 bits

           = 18 bits

  – so tag is leftmost 18 bits of memory address

# Accessing data in cache

## Address (hex)

- Ex.: 16KB of data, direct-mapped,
  4 word blocks

- Read 4 addresses
  - 0x00000014,
    0x0000001C,
    0x00000034, 0x00008014

- Memory values on right:
  - only cache/memory level
    of hierarchy

...
00000010        a
00000014        b
00000018        c
0000001C        d

...
00000030        e
00000034        f
00000038        g
0000003C        h

...
00008010        i
00008014        j
00008018        k
0000801C        l

...                     ...

# Accessing data in cache [contd…]

- 4 Addresses:
  - 0x00000014, 0x0000001C, 0x00000034, 0x00008014

- 4 Addresses divided (for convenience) into Tag, Index, Byte Offset fields

```
0000000000000000 000000001 0100
0000000000000000 000000001 1100
0000000000000000 000000011 0100
0000000000000010 000000001 0100
```

       Tag                     Index    Offset

# 16 KB Direct Mapped Cache, 16B blocks

- <u>Valid bit:</u> determines whether anything is stored in that row (when computer initially turned on, all entries are invalid)

**Valid**　　　　　　　　　　　　**Example Block**

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|-------|-----|-------|-------|-------|-------|
| 0 | 0 | | | | |
| 1 | 0 | | | | |
| 2 | 0 | | | | |
| 3 | 0 | | | | |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

**...**　　　　　　　　**...**

| Index | Tag | | | | |
|-------|-----|--|--|--|--|
| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# Read 0x00000014 = 0…00 0..001 0100

- 00000000000000000 0000000001 0100

**Tag field**      **Index field**    **Offset**

**Valid**

| Index | Valid | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|-------|-------|-----|-------|-------|-------|-------|
| 0 | 0 | | | | | |
| 1 | 0 | | | | | |
| 2 | 0 | | | | | |
| 3 | 0 | | | | | |
| 4 | 0 | | | | | |
| 5 | 0 | | | | | |
| 6 | 0 | | | | | |
| 7 | 0 | | | | | |

...      ...

| 1022 | | | | | | |
| 1023 | 0 | | | | | |

0

# So we read block 1 (0000000001)

- 0000000000000000000 0000000001 0100

**Tag field**        **Index field**   **Offset**

**Valid**

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|
| 0 | 0 | | | | |
| 1 | 0 | | | | |
| 2 | 0 | | | | |
| 3 | 0 | | | | |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...         ...

| | | | | | |
|---|---|---|---|---|---|
| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# No valid data

- 00000000000000000 0000000001 0100

**Tag field**      **Index field**   **Offset**

| Index | Valid Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|-------|-----------|-------|-------|-------|-------|
| 0 | 0 | | | | |
| 1 | 0 | | | | |
| 2 | 0 | | | | |
| 3 | 0 | | | | |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# So load that data into cache, setting tag, valid

- <span style="color:red">0000000000000000</span> 0000000001 0100

**Tag field**  **Index field**  **Offset**

**Valid**

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|
| 0 | 0 | | | | |
| 1 | 1   0 | a | b | c | d |
| 2 | 0 | | | | |
| 3 | 0 | | | | |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...  ...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# Read from cache at offset, return word b

- 0000000000000000000 0000000001 0100

Tag field            Index field    Offset

**Valid**

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|-------|-----|-------|-------|-------|-------|
| 0 | 0 | . | | | |
| 1 | 1   0 | a | b | c | d |
| 2 | 0 | | | | |
| 3 | 0 | | | | |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...         ...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# Read 0x0000001C = 0…00 0..001 1100

- 00000000000000000  0000000001  1100

| | | Tag field | Index field | Offset |

Valid

| Index | | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|-------|---|-----|-------|-------|-------|-------|
| 0 | 0 | | | | | |
| 1 | 1 | 0 | a | b | c | d |
| 2 | 0 | | | | | |
| 3 | 0 | | | | | |
| 4 | 0 | | | | | |
| 5 | 0 | | | | | |
| 6 | 0 | | | | | |
| 7 | 0 | | | | | |

…                    …

| 1022 | 0 | | | | | |
| 1023 | 0 | | | | | |

# Data valid, tag OK, so read offset return word d

- <u>00000000000000000</u>  <u>0000000001</u>  <u>1100</u>

**Valid**

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | <u>**0xc-f**</u> |
|-------|-----|-------|-------|-------|--------|
| 0 | 0 |  |  |  | . |
| <u>**1**</u> **1** | <u>**0**</u> | **a** | **b** | **c** | **d** |
| 2 | 0 |  |  |  |  |
| 3 | 0 |  |  |  |  |
| 4 | 0 |  |  |  |  |
| 5 | 0 |  |  |  |  |
| 6 | 0 |  |  |  |  |
| 7 | 0 |  |  |  |  |

...          ...

| 1022 | 0 |  |  |  |  |
| 1023 | 0 |  |  |  |  |

# Read 0x00000034 = 0…00 0..011 0100

- 00000000000000000000 0000000011 0100

| | Valid<br>Tag field | | Index field | Offset |
|---|---|---|---|---|

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|
| 0 | 0 | | | | |
| 1 | **1**   **0** | **a** | **b** | **c** | **d** |
| 2 | 0 | | | | |
| 3 | 0 | | | | |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...     ...

| 1022 | 0 | | | | |
|---|---|---|---|---|---|
| 1023 | 0 | | | | |

# So read block 3

- 0000000000000000000 <span style="color:red">0000000011</span> 0100

| | | **Tag field** | | **Index field** | **Offset** |
|---|---|---|---|---|---|

**Valid**

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|
| 0 | 0 | | | | |
| 1 | **1** **0** | **a** | **b** | **c** | **d** |
| 2 | 0 | | | | |
| 3 | 0 | | | | |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...          ...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# No valid data

- 00000000000000000000 0000000011 0100

| | | **Tag field** | **Index field** | **Offset** |

**Valid**

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|
| 0 | 0 | | | | |
| 1 | **1** | **0** | **a** | **b** | **c** | **d** |
| 2 | 0 | | | | |
| 3 | 0 | | | | |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...

...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# Load that cache block, return word f

- 000000000000000000 0000000011 0100

**Tag field**      **Index field**    **Offset**

**Valid**

| Index | Tag | `0x0-3` | `0x4-7` | `0x8-b` | `0xc-f` |
|-------|-----|---------|---------|---------|---------|
| 0 | 0 | | | | |
| 1 | **1**   **0** | **a** | **b** | **c** | **d** |
| 2 | 0 | . | | | |
| 3 | **1**   **0** | **e** | **f** | . **g** | **h** |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...       ...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# Read 0x00008014 = 0…10 0..001 0100

- 00000000000000010 0000000001 0100

**Valid**

**Tag field**     **Index field**     **Offset**

| Index | Valid | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|---|
| 0 | 0 | | | | | |
| 1 | 1 | 0 | a | b | c | d |
| 2 | 0 | | | | | |
| 3 | 1 | 0 | e | f | g | h |
| 4 | 0 | | | | | |
| 5 | 0 | | | | | |
| 6 | 0 | | | | | |
| 7 | 0 | | | | | |

...           ...

| 1022 | 0 | | | | | |
| 1023 | 0 | | | | | |

- 00000000000000010 0000000001 0100

**Tag field**      **Index field**   **Offset**

| Index | Valid | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|-------|-------|-----|-------|-------|-------|-------|
| 0 | 0 | | | | | |
| 1 | 1 | 0 | a | b | c | d |
| 2 | 0 | | | | | |
| 3 | 1 | 0 | e | f | g | h |
| 4 | 0 | | | | | |
| 5 | 0 | | | | | |
| 6 | 0 | | | | | |
| 7 | 0 | | | | | |

...       ...

| 1022 | 0 | | | | | |
| 1023 | 0 | | | | | |

# Cache Block 1 Tag does not match (0 != 2)

- 00000000000000010 0000000001 0100

**Tag field**    **Index field**    **Offset**

| Valid Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|
| 0 | 0 | | | | |
| 1 | **1** **0** | **a** | **b** | **c** | **d** |
| 2 | 0 | | | | |
| 3 | **1** **0** | **e** | **f** | **g** | **h** |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...    ...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# Miss, so replace block 1 with new data & tag

- 00000000000000010  0000000001  0100

| | Tag field | | Index field | Offset |
|---|---|---|---|---|

| Valid Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|
| 0 | 0 | | | | |
| 1 | **1** **2** | **i** | **j** | **k** | **l** |
| 2 | 0 | | | | |
| 3 | **1** **0** | **e** | **f** | **g** | **h** |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...  ...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

- 00000000000000010 0000000001 0100

| | | Tag field | | index field | Offset |
|---|---|---|---|---|---|

Valid

| Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|---|---|---|---|---|---|
| 0 | 0 | | . | | |
| 1 | 1 2 | i | j | k | l |
| 2 | 0 | | . | | |
| 3 | 1 0 | e | f | g | h |
| 4 | 0 | | | | |
| 5 | 0 | | | | |
| 6 | 0 | | | | |
| 7 | 0 | | | | |

...  ...

| 1022 | 0 | | | | |
| 1023 | 0 | | | | |

# Q3: Which block should be replaced on a miss?

- Easy for Direct Mapped
- Set Associative or Fully Associative:
  - Random
  - LRU (Least Recently Used)

| Assoc: | 2-way | | 4-way | | 8-way | |
|---|---|---|---|---|---|---|
| Size | LRU | Ran | LRU | Ran | LRU | Ran |
| 16 KB | 5.2% | 5.7% | 4.7% | 5.3% | 4.4% | 5.0% |
| 64 KB | 1.9% | 2.0% | 1.5% | 1.7% | 1.4% | 1.5% |
| 256 KB | 1.15% | 1.17% | 1.13% | 1.13% | 1.12% | 1.12% |

# Q3: After a cache read miss, if there are no empty cache blocks, which block should be removed from the cache?

**The Least Recently Used (LRU) block? Appealing, but hard to implement for high associativity**

**A randomly chosen block? Easy to implement, how well does it work?**

## Miss Rate for 2-way Set Associative Cache

| Size | Random | LRU |
|---|---|---|
| 16 KB | 5.7% | 5.2% |
| 64 KB | 2.0% | 1.9% |
| 256 KB | 1.17% | 1.15% |

**Also, try other LRU approx.**

# Q4: What happens on a write?

| | Write-Through | Write-Back |
|---|---|---|
| **Policy** | **Data written to cache block** **also written to lower-level memory** | **Write data only to the cache** **Update lower level when a block falls out of the cache** |
| **Debug** | **Easy** | **Hard** |
| **Do read misses produce writes?** | **No** | **Yes** |
| **Do repeated writes make it to lower level?** | **Yes** | **No** |

**Additional option -- let writes to an un-cached address allocate a new cache line ("write-allocate").**

# Write Buffers for Write-Through Caches



**Holds data awaiting write-through to lower level memory**

**Q. Why a write buffer ?**

**A. So CPU doesn't stall**

**Q. Why a buffer, why not just one register ?**

**A. Bursts of writes are common.**

**Q. Are Read After Write (RAW) hazards an issue for write buffer?**

**A. Yes! Drain buffer before next read, or send read 1st after check write buffers.**

# Memory Hierarchy Basics

- When a word is not found in the cache, a *miss* occurs:
  - Fetch word from lower level in hierarchy, requiring a higher latency reference
  - Lower level may be another cache or the main memory
  - Also fetch the other words contained within the *block*
    - Takes advantage of spatial locality
  - Place block into cache in any location within its *set*, determined by address
    - block address MOD number of sets

# **Memory Hierarchy Basics**

- *n* sets => *n-way set associative*
  - *Direct-mapped cache* => one block per set
  - *Fully associative* => one set

- Writing to cache:  two strategies
  - *Write-through*
    - Immediately update lower levels of hierarchy
  - *Write-back*
    - Only update lower levels of hierarchy when an updated block is replaced
  - Both strategies use *write buffer* to make writes asynchronous

MORGAN KAUFMANN

# **Memory Hierarchy Basics**

- Miss rate
  - Fraction of cache access that result in a miss

- Causes of misses
  - Compulsory
    - First reference to a block
  - Capacity
    - Blocks discarded and later retrieved
  - Conflict
    - Program makes repeated references to multiple addresses from different blocks that map to the same location in the cache

# Memory Hierarchy Basics

$$\frac{\text{Misses}}{\text{Instruction}} = \frac{\text{Miss rate} \times \text{Memory accesses}}{\text{Instruction count}} = \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}}$$

$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

- Note that speculative and multithreaded processors may execute other instructions during a miss
  - Reduces performance impact of misses

# Memory Hierarchy Basics

- Six basic cache optimizations:
  - Larger block size
    - Reduces compulsory misses
    - Increases capacity and conflict misses, increases miss penalty
  - Larger total cache capacity to reduce miss rate
    - Increases hit time, increases power consumption
  - Higher associativity
    - Reduces conflict misses
    - Increases hit time, increases power consumption
  - Higher number of cache levels
    - Reduces overall memory access time
  - Giving priority to read misses over writes
    - Reduces miss penalty
  - Avoiding address translation in cache indexing
    - Reduces hit time

# Ten Advanced Optimizations

- Small and simple first level caches
  - Critical timing path:
    - addressing tag memory, then
    - comparing tags, then
    - selecting correct set
  - Direct-mapped caches can overlap tag compare and transmission of data
  - Lower associativity reduces power because fewer cache lines are accessed

# L1 Size and Associativity

Access time vs. size and associativity

# L1 Size and Associativity

Energy per read vs. size and associativity

# Way Prediction

- To improve hit time, predict the way to pre-set mux
  - Mis-prediction gives longer hit time
  - Prediction accuracy
    - > 90% for two-way
    - > 80% for four-way
    - I-cache has better accuracy than D-cache
  - First used on MIPS R10000 in mid-90s
  - Used on ARM Cortex-A8
- Extend to predict block as well
  - "Way selection"
  - Increases mis-prediction penalty

# Pipelining Cache

- Pipeline cache access to improve bandwidth
  - Examples:
    - Pentium:  1 cycle
    - Pentium Pro – Pentium III:  2 cycles
    - Pentium 4 – Core i7:  4 cycles

- Increases branch mis-prediction penalty
- Makes it easier to increase associativity

# Nonblocking Caches

- Allow hits before previous misses complete
  - "Hit under miss"
  - "Hit under multiple miss"
- L2 must support this
- In general, processors can hide L1 miss penalty but not L2 miss penalty

# Multibanked Caches

- Organize cache as independent banks to support simultaneous access
    - ARM Cortex-A8 supports 1-4 banks for L2
    - Intel i7 supports 4 banks for L1 and 8 banks for L2

- Interleave banks according to block address



**Figure 2.6** Four-way interleaved cache banks using block addressing. Assuming 64 bytes per blocks, each of these addresses would be multiplied by 64 to get byte addressing.

# **Critical Word First, Early Restart**

- Critical word first
    - Request missed word from memory first
    - Send it to the processor as soon as it arrives
- Early restart
    - Request words in normal order
    - Send missed work to the processor as soon as it arrives

- Effectiveness of these strategies depends on block size and likelihood of another access to the portion of the block that has not yet been fetched

# Merging Write Buffer

- When storing to a block that is already pending in the write buffer, update write buffer
- Reduces stalls due to full write buffer
- Do not apply to I/O addresses

| Write address | V | | V | | V | | V | |
|---|---|---|---|---|---|---|---|---|
| 100 | 1 | Mem[100] | 0 | | 0 | | 0 | |
| 108 | 1 | Mem[108] | 0 | | 0 | | 0 | |
| 116 | 1 | Mem[116] | 0 | | 0 | | 0 | |
| 124 | 1 | Mem[124] | 0 | | 0 | | 0 | |

No write buffering

| Write address | V | | V | | V | | V | |
|---|---|---|---|---|---|---|---|---|
| 100 | 1 | Mem[100] | 1 | Mem[108] | 1 | Mem[116] | 1 | Mem[124] |
| | 0 | | 0 | | 0 | | 0 | |
| | 0 | | 0 | | 0 | | 0 | |
| | 0 | | 0 | | 0 | | 0 | |

Write buffering

# Compiler Optimizations

- McFarling [1989] reduced caches misses by 75% on 8KB direct mapped cache, 4B blocks in software
- Instructions
  - Reorder procedures in memory so as to reduce conflict misses
  - Profiling to look at conflicts (using tools they developed)
- Data
  - *Merging Arrays*: improve spatial locality by single array of compound elements vs. 2 arrays
  - *Loop Interchange*: change nesting of loops to access data in order stored in memory
  - *Loop Fusion*: Combine 2 independent loops that have same looping and some variables overlap
  - *Blocking*: Improve temporal locality by accessing "blocks" of data repeatedly vs. going down whole columns or rows

# Merging Arrays Example

```c
/* Before: 2 sequential arrays */
int val[SIZE];
int key[SIZE];

/* After: 1 array of stuctures */
struct merge {
  int val;
  int key;
};
struct merge merged_array[SIZE];
```

Reducing conflicts between val & key;
  improve spatial locality

# Loop Interchange Example

```
/* Before */
for (k = 0; k < 100; k = k+1)
  for (j = 0; j < 100; j = j+1)
  for (i = 0; i < 5000; i = i+1)
  x[i][j] = 2 * x[i][j];
/* After */
for (k = 0; k < 100; k = k+1)
  for (i = 0; i < 5000; i = i+1)
  for (j = 0; j < 100; j = j+1)
  x[i][j] = 2 * x[i][j];
```

Sequential accesses instead of striding through memory every 100 words; improved spatial locality

# Loop Fusion Example

```
/* Before */
for (i = 0; i < N; i = i+1)
   for (j = 0; j < N; j = j+1)
    a[i][j] = 1/b[i][j] * c[i][j];
for (i = 0; i < N; i = i+1)
   for (j = 0; j < N; j = j+1)
   d[i][j] = a[i][j] + c[i][j];
/* After */
for (i = 0; i < N; i = i+1)
   for (j = 0; j < N; j = j+1)
   {    a[i][j] = 1/b[i][j] * c[i][j];
   d[i][j] = a[i][j] + c[i][j];}
```
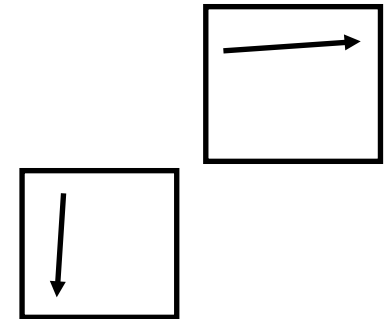
2 misses per access to `a` & `c` vs. one miss per access; improve spatial locality

# Blocking Example

```
/* Before */
for (i = 0; i < N; i = i+1)
   for (j = 0; j < N; j = j+1)
      {r = 0;
       for (k = 0; k < N; k = k+1){
          r = r + y[i][k]*z[k][j];};
       x[i][j] = r;
      };
```

- Two Inner Loops:
  - Read all NxN elements of z[]
  - Read N elements of 1 row of y[] repeatedly
  - Write N elements of 1 row of x[]
- Capacity Misses a function of N & Cache Size:
  - $2N^3 + N^2$ => (assuming no conflict; otherwise …)
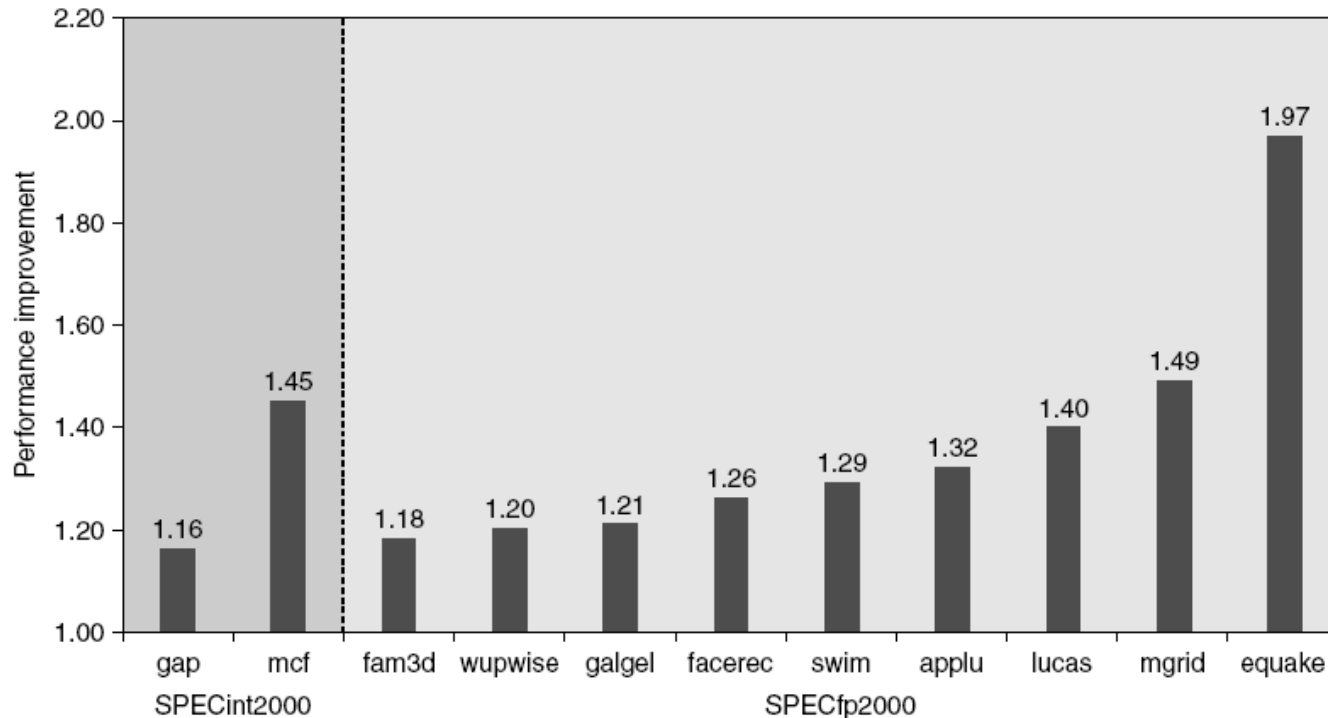- Idea: compute on BxB submatrix that fits in cache

# Blocking Example

```
/* After */
for (jj = 0; jj < N; jj = jj+B)
for (kk = 0; kk < N; kk = kk+B)
for (i = 0; i < N; i = i+1)
    for (j = jj; j < min(jj+B-1,N); j = j+1)
  {r = 0;
    for (k = kk; k < min(kk+B-1,N); k = k+1) {
  r = r + y[i][k]*z[k][j];};
   x[i][j] = x[i][j] + r;
   };
```

- B called *Blocking Factor*
- Capacity Misses from $2N^3 + N^2$ to $2N^3/B + N^2$

# Hardware Prefetching

■ Fetch two blocks on miss (include next sequential block)



Pentium 4 Pre-fetching

# Compiler Prefetching

- Insert prefetch instructions before data is needed
- Non-faulting: prefetch doesn't cause exceptions

- Register prefetch
  - Loads data into register
- Cache prefetch
  - Loads data into cache

- Combine with loop unrolling and software pipelining

# Summary

| Technique | Hit time | Band-width | Miss penalty | Miss rate | Power consumption | Hardware cost/ complexity | Comment |
|---|---|---|---|---|---|---|---|
| Small and simple caches | + | | | − | + | 0 | Trivial; widely used |
| Way-predicting caches | + | | | | + | 1 | Used in Pentium 4 |
| Pipelined cache access | − | + | | | | 1 | Widely used |
| Nonblocking caches | | + | + | | | 3 | Widely used |
| Banked caches | | + | | | + | 1 | Used in L2 of both i7 and Cortex-A8 |
| Critical word first and early restart | | | + | | | 2 | Widely used |
| Merging write buffer | | | + | | | 1 | Widely used with write through |
| Compiler techniques to reduce cache misses | | | | + | | 0 | Software is a challenge, but many compilers handle common linear algebra calculations |
| Hardware prefetching of instructions and data | | | + | + | − | 2 instr., 3 data | Most provide prefetch instructions; modern high-end processors also automatically prefetch in hardware. |
| Compiler-controlled prefetching | | | + | + | | 3 | Needs nonblocking cache; possible instruction overhead; in many CPUs |

**Figure 2.11** Summary of 10 advanced cache optimizations showing impact on cache performance, power consumption, and complexity. Although generally a technique helps only one factor, prefetching can reduce misses if done sufficiently early; if not, it can reduce miss penalty. + means that the technique improves the factor, − means it hurts that factor, and blank means it has no impact. The complexity measure is subjective, with 0 being the easiest and 3 being a challenge.

# Memory Technology

- Performance metrics
  - Latency is concern of cache
  - Bandwidth is concern of multiprocessors and I/O
  - Access time
    - Time between read request and when desired word arrives
  - Cycle time
    - Minimum time between unrelated requests to memory

- DRAM used for main memory, SRAM used for cache

# Memory Technology

- ## SRAM
  - Requires low power to retain bit
  - Requires 6 transistors/bit

- ## DRAM
  - Must be re-written after being read
  - Must also be periodically refeshed
    - Every ~ 8 ms
    - Each row can be refreshed simultaneously
  - One transistor/bit
  - Address lines are multiplexed:
    - Upper half of address:  row access strobe (RAS)
    - Lower half of address:  column access strobe (CAS)

# Memory Technology

- Amdahl:
  - Memory capacity should grow linearly with processor speed
  - Unfortunately, memory capacity and speed has not kept pace with processors

- Some optimizations:
  - Multiple accesses to same row
  - Synchronous DRAM
    - Added clock to DRAM interface
    - Burst mode with critical word first
  - Wider interfaces
  - Double data rate (DDR)
  - Multiple banks on each DRAM device

# Memory Optimizations

| Production year | Chip size | DRAM Type | Row access strobe (RAS) | | Column access strobe (CAS)/ data transfer time (ns) | Cycle time (ns) |
|---|---|---|---|---|---|---|
| | | | Slowest DRAM (ns) | Fastest DRAM (ns) | | |
| 1980 | 64K bit | DRAM | 180 | 150 | 75 | 250 |
| 1983 | 256K bit | DRAM | 150 | 120 | 50 | 220 |
| 1986 | 1M bit | DRAM | 120 | 100 | 25 | 190 |
| 1989 | 4M bit | DRAM | 100 | 80 | 20 | 165 |
| 1992 | 16M bit | DRAM | 80 | 60 | 15 | 120 |
| 1996 | 64M bit | SDRAM | 70 | 50 | 12 | 110 |
| 1998 | 128M bit | SDRAM | 70 | 50 | 10 | 100 |
| 2000 | 256M bit | DDR1 | 65 | 45 | 7 | 90 |
| 2002 | 512M bit | DDR1 | 60 | 40 | 5 | 80 |
| 2004 | 1G bit | DDR2 | 55 | 35 | 5 | 70 |
| 2006 | 2G bit | DDR2 | 50 | 30 | 2.5 | 60 |
| 2010 | 4G bit | DDR3 | 36 | 28 | 1 | 37 |
| 2012 | 8G bit | DDR3 | 30 | 24 | 0.5 | 31 |

**Figure 2.13** Times of fast and slow DRAMs vary with each generation. (Cycle time is defined on page 95.) Performance improvement of row access time is about 5% per year. The improvement by a factor of 2 in column access in 1986 accompanied the switch from NMOS DRAMs to CMOS DRAMs. The introduction of various burst transfer modes in the mid-1990s and SDRAMs in the late 1990s has significantly complicated the calculation of access time for blocks of data; we discuss this later in this section when we talk about SDRAM access time and power. The DDR4 designs are due for introduction in mid- to late 2012. We discuss these various forms of DRAMs in the next few pages.

# Memory Optimizations

| Standard | Clock rate (MHz) | M transfers per second | DRAM name | MB/sec /DIMM | DIMM name |
|---|---|---|---|---|---|
| DDR | 133 | 266 | DDR266 | 2128 | PC2100 |
| DDR | 150 | 300 | DDR300 | 2400 | PC2400 |
| DDR | 200 | 400 | DDR400 | 3200 | PC3200 |
| DDR2 | 266 | 533 | DDR2-533 | 4264 | PC4300 |
| DDR2 | 333 | 667 | DDR2-667 | 5336 | PC5300 |
| DDR2 | 400 | 800 | DDR2-800 | 6400 | PC6400 |
| DDR3 | 533 | 1066 | DDR3-1066 | 8528 | PC8500 |
| DDR3 | 666 | 1333 | DDR3-1333 | 10,664 | PC10700 |
| DDR3 | 800 | 1600 | DDR3-1600 | 12,800 | PC12800 |
| DDR4 | 1066–1600 | 2133–3200 | DDR4-3200 | 17,056–25,600 | PC25600 |

**Figure 2.14** Clock rates, bandwidth, and names of DDR DRAMS and DIMMs in 2010. Note the numerical relationship between the columns. The third column is twice the second, and the fourth uses the number from the third column in the name of the DRAM chip. The fifth column is eight times the third column, and a rounded version of this number is used in the name of the DIMM. Although not shown in this figure, DDRs also specify latency in clock cycles as four numbers, which are specified by the DDR standard. For example, DDR3-2000 CL 9 has latencies of 9-9-9-28. What does this mean? With a 1 ns clock (clock cycle is one-half the transfer rate), this indicate 9 ns for row to columns address (RAS time), 9 ns for column access to data (CAS time), and a minimum read time of 28 ns. Closing the row takes 9 ns for precharge but happens only when the reads from that row are finished. In burst mode, transfers occur on every clock on both edges, when the first RAS and CAS times have elapsed. Furthermore, the precharge in not needed until the entire row is read. DDR4 will be produced in 2012 and is expected to reach clock rates of 1600 MHz in 2014, when DDR5 is expected to take over. The exercises explore these details further.
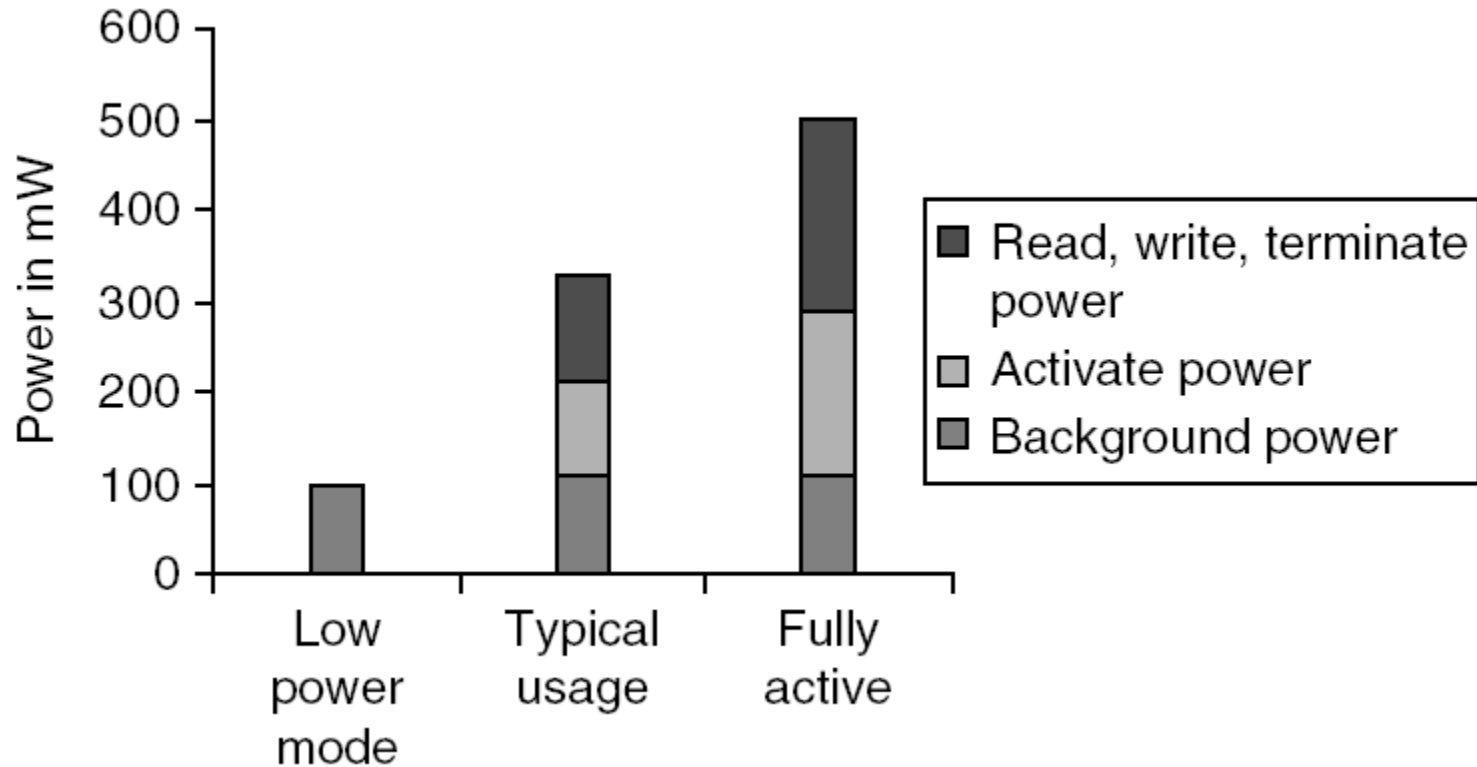
# Memory Optimizations

- DDR:
  - DDR2
    - Lower power (2.5 V -> 1.8 V)
    - Higher clock rates (266 MHz, 333 MHz, 400 MHz)
  - DDR3
    - 1.5 V
    - 800 MHz
  - DDR4
    - 1-1.2 V
    - 1600 MHz

- GDDR5 is graphics memory based on DDR3

# Memory Optimizations

- Graphics memory:
  - Achieve 2-5 X bandwidth per DRAM vs. DDR3
    - Wider interfaces (32 vs. 16 bit)
    - Higher clock rate
      - Possible because they are attached via soldering instead of socketted DIMM modules

- Reducing power in SDRAMs:
  - Lower voltage
  - Low power mode (ignores clock, continues to refresh)

# Memory Power Consumption

# Flash Memory

- Type of EEPROM
- Must be erased (in blocks) before being overwritten
- Non volatile
- Limited number of write cycles
- Cheaper than SDRAM, more expensive than disk
- Slower than SRAM, faster than disk

# **Memory Dependability**

- Memory is susceptible to cosmic rays
- *Soft errors*:  dynamic errors
  - Detected and fixed by error correcting codes (ECC)
- *Hard errors*:  permanent errors
  - Use sparse rows to replace defective rows

- Chipkill:  a RAID-like error recovery technique

# Virtual Memory

- Protection via virtual memory
  - Keeps processes in their own memory space

- Role of architecture:
  - Provide user mode and supervisor mode
  - Protect certain aspects of CPU state
  - Provide mechanisms for switching between user mode and supervisor mode
  - Provide mechanisms to limit memory accesses
  - Provide TLB to translate addresses

# Virtual Machines

- Supports isolation and security
- Sharing a computer among many unrelated users
- Enabled by raw speed of processors, making the overhead more acceptable

- Allows different ISAs and operating systems to be presented to user programs
  - "System Virtual Machines"
  - SVM software is called "virtual machine monitor" or "hypervisor"
  - Individual virtual machines run under the monitor are called "guest VMs"

# Impact of VMs on Virtual Memory

- Each guest OS maintains its own set of page tables
  - VMM adds a level of memory between physical and virtual memory called "real memory"
  - VMM maintains shadow page table that maps guest virtual addresses to physical addresses
    - Requires VMM to detect guest's changes to its own page table
    - Occurs naturally if accessing the page table pointer is a privileged operation