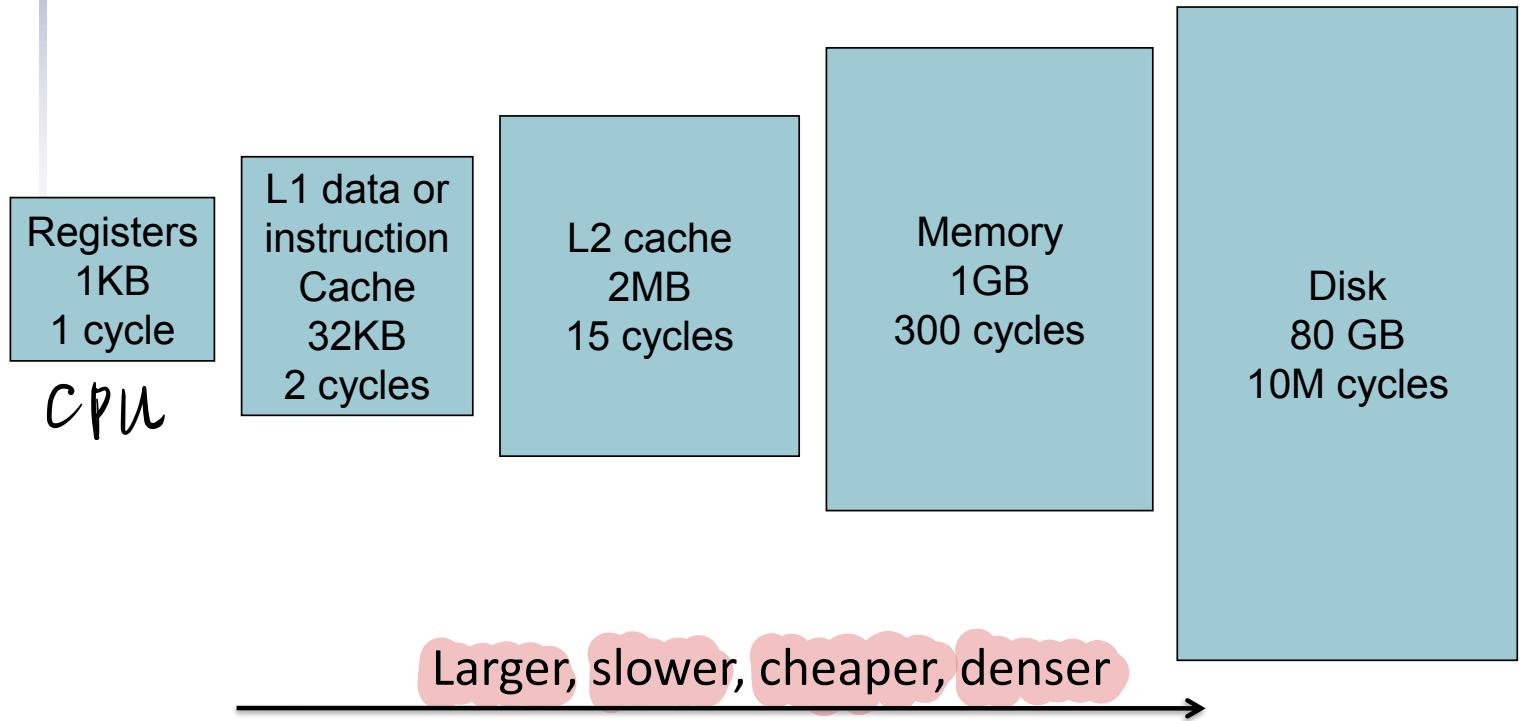


Chapter 5

**Large and Fast: Exploiting
Memory Hierarchy**

Memory Hierarchy

- Various storage devices in computers:



Memory Technology

- Access time and price per bit vary widely among different technologies

Memory technology	Typical access time	\$ per GiB in 2012
SRAM semiconductor memory	0.5–2.5 ns	\$500–\$1000
DRAM semiconductor memory	50–70 ns	\$10–\$20
Flash semiconductor memory	5,000–50,000 ns	\$0.75–\$1.00
Magnetic disk	5,000,000–20,000,000 ns	\$0.05–\$0.10

Data in 2012

- Ideal memory
 - ◆ Access time of Cache
 - ◆ Capacity and cost/GB of disk

Outline

- Cache (CPU $\leftarrow\rightarrow$ memory)
 - ◆ Direct mapped cache
 - ◆ Set associative cache
 - ◆ Multi-level cache
- Virtual memory (memory $\leftarrow\rightarrow$ disk)
- Dependable memory
- Real examples

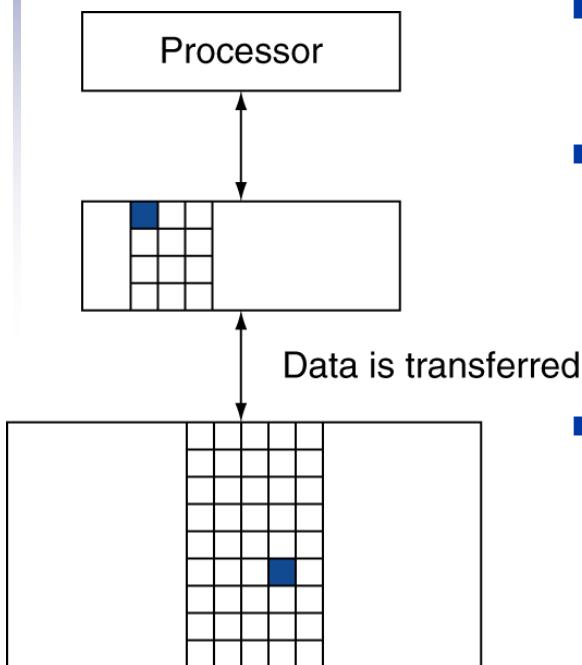
Cache Hierarchies

- Data and instructions are stored on DRAM chips
 - ◆ DRAM is a technology that has **high bit density**, but relatively **poor latency**
 - ◆ an access to data in memory can take as many as 300 cycles!
- Hence, some data is stored on the processor in a structure called the **cache**
 - ◆ caches employ SRAM technology, which is **faster**, but has lower bit density ~~same size less data~~
- Internet browsers also cache web pages – same concept

Memory hierarchy

- Store everything on disk
 - Copy recently accessed (and nearby) items from disk to smaller DRAM memory
 - ◆ Main memory
 - Copy more recently accessed (and nearby) items from DRAM to smaller SRAM memory
 - ◆ Cache memory attached to CPU
- Register*

Memory Hierarchy Levels



- **Block** (also called line): unit of copying
 - ◆ May be multiple words
- The memory in upper level is originally empty
- If accessed data is absent
 - ◆ **Miss**: block copied from lower level
 - Time taken: **miss penalty**
 - Miss ratio: **misses/Accesses**
- If accessed data is present in upper level
 - ◆ **Hit**: access satisfied by upper level
 - Hit ratio: $\text{hits}/\text{Accesses} = 1 - \text{miss ratio}$
 - ◆ Then accessed data supplied from upper level

Locality

- Why do caches work?
 - ◆ **Temporal locality** if you used some data recently, you will likely use it again **时间局部性**
 - ◆ **Spatial locality** if you used some data recently, you will likely access its neighbors **空间局部性** \Leftarrow the size of the block
- No hierarchy:
 - ◆ average access time for data = 300 cycles (*memory*)
- 32KB 1-cycle L1 cache that has a hit rate of 95%:
 - ◆ average access time = $0.95 \uparrow \times 1 + 0.05 \downarrow \times (301) = 16$ cycles

SRAM Technology

random access 与 Disk 中的连续访问相对应

- Static RAM

- ◆ Memory arrays with a single read/write port

- It's Volatile (SRAM/DRAM)

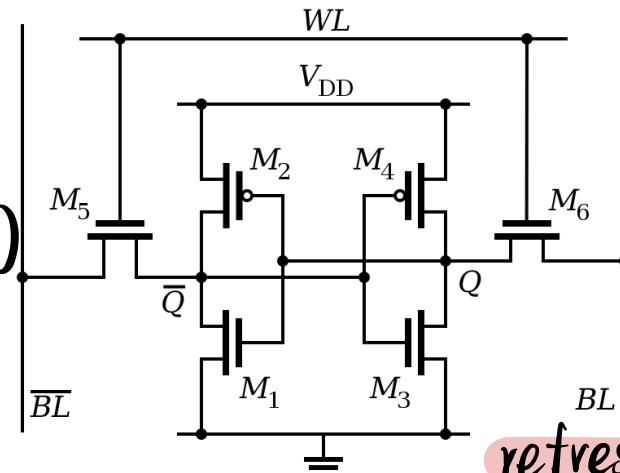
- ◆ The data will lost when SRAM is not powered

- Compared with DRAM

DRAM is easy to lost, need to refresh

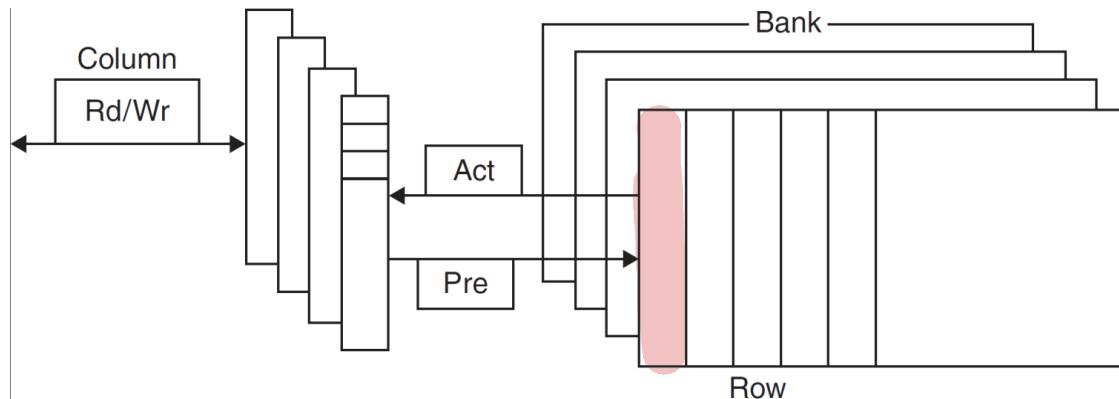
- ◆ Don't need to refresh, use 6-8 transistors to install a bit

- Used in CPU cache, integrated onto the processor chip



DRAM Technology

- Data stored as a charge in a capacitor
 - ◆ Single transistor used to access the charge
 - ◆ Must periodically be refreshed
 - Read contents and write back
 - Performed on a DRAM “row”

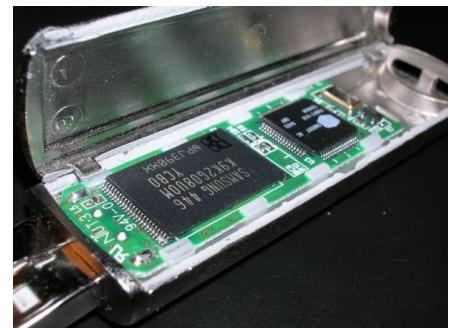


Advanced DRAM Organization

- Bits in a DRAM are organized as a rectangular array
 - ◆ DRAM accesses an entire row
 - ◆ Burst mode: supply successive words from a row with reduced latency
- Synchronous DRAM
 - ◆ A clock is added, the memory and processor are synchronized
 - ◆ Allows for consecutive accesses in bursts without needing to send each address
 - ◆ Improves bandwidth
- Double data rate (DDR) DRAM
 - ◆ Transfer on rising and falling clock edges
 - ◆ DDR4-3200 DRAM: 3200M times of transfer per second

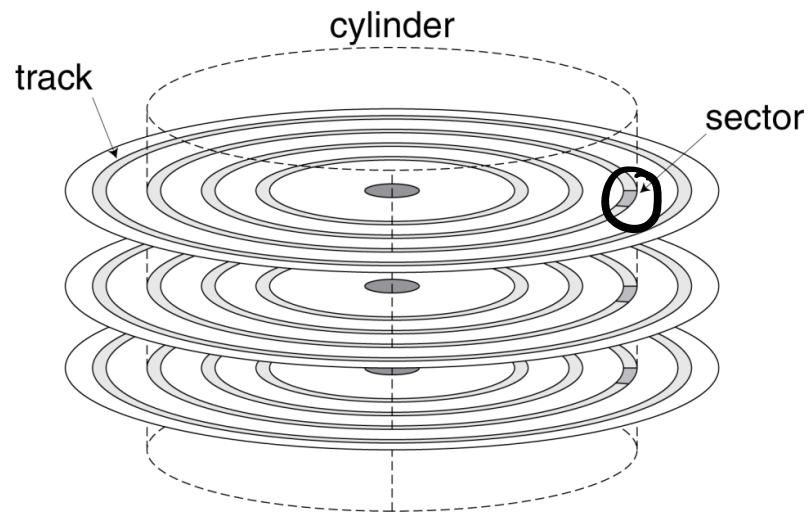
Flash Storage

- Nonvolatile semiconductor storage
 - ◆ 100× – 1000× faster than disk
 - ◆ Smaller, lower power, more robust
 - ◆ But more \$/GB (between disk and DRAM)
- Flash bits wears out after 1000's of accesses
 - ◆ Not suitable for direct RAM or disk replacement
 - ◆ Wear leveling: remap data to less used blocks



Disk Storage

- Nonvolatile, rotating magnetic storage



Disk Sectors and Access

- Each sector records
 - ◆ Sector ID
 - ◆ Data (512 bytes, 4096 bytes proposed)
 - ◆ Error correcting code (ECC)
 - Used to hide defects and recording errors
- Access to a sector involves
 - ◆ Queuing delay if other accesses are pending
 - ◆ Seek: move the heads
 - ◆ Rotational latency
 - ◆ Data transfer
 - ◆ Controller overhead

Cache Memory

- Cache memory
 - ◆ The level of the memory hierarchy closest to the CPU
- Given accesses X_1, \dots, X_{n-1}, X_n

X_4
X_1
X_{n-2}
X_{n-1}
X_2
X_3

a. Before the reference to X_n

X_4
X_1
X_{n-2}
X_{n-1}
x_2
x_n
X_3

b. After the reference to X_n

- How do we know if the data is present?
- Where do we look?

Memory Structure

■ Address and data

- ◆ Address is the index, are not stored in memory
 - Address can be in **unit of byte** or **in unit of word**
- ◆ Only data is stored in memory

address

000

001

010

011

100

101

110

111

data

Byte 1

Byte 2

Byte 3

...

Word1

address

000

001

010

011

100

101

110

111

data

Word1

Word2

Word3

Word4

...

in unit of byte

in unit of word

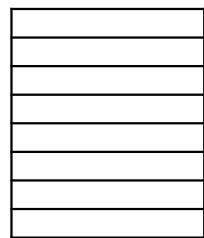
Direct Mapped Cache

- Memory size: 32 words, cache size: 8 words
- The address is in unit of word

Last three bits is better



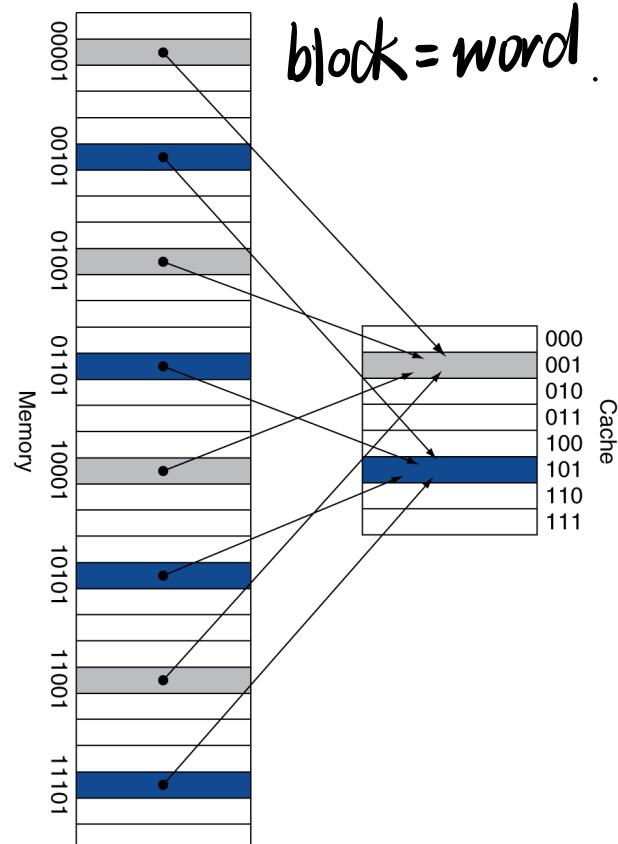
due to locality.



Direct Mapped Cache

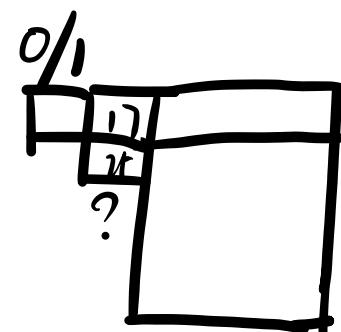
- Memory size: 32 words, cache size: 8 words
- The address is in unit of word
- Direct mapped cache: (直接映射)

 - ◆ Location determined by address
 - ◆ One data in memory is mapped to only one location in cache
 - ◆ Use low-order address bits or high-order bits?
 - ◆ The lower bits defines the address of the cache



Tags and Valid Bits

- How do we know which particular block is stored in a cache location?
 - ◆ Store block address as well as the data
 - ◆ Actually, only need the high-order bits
 - ◆ Called the **tag**
- What if there is no data in a location?
 - ◆ **Valid bit:** 1 = present, 0 = not present
 - ◆ **Initially 0**



Cache Example

- 8-blocks, 1 word/block, direct mapped
- Initial state

Index	V	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	N		
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Miss	110

Index	V	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	N		
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Miss	110

Index	V	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
26	11 010	Miss	010

Index	V	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
26	11 010	Miss	010

Index	V	Tag	Data
000	N		
001	N		
010	Y	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Hit	110
26	11 010	Hit	010

Index	V	Tag	Data
000	N		
001	N		
010	Y	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
16	10 000	Miss	000
3	00 011	Miss	011
16	10 000	Hit	000

Index	V	Tag	Data
000	N		
001	N		
010	Y	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
16	10 000	Miss	000
3	00 011	Miss	011
16	10 000	Hit	000

Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
010	Y	11	Mem[11010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
18	10 010	Miss	010

Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
010	Y	11	Mem[11010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

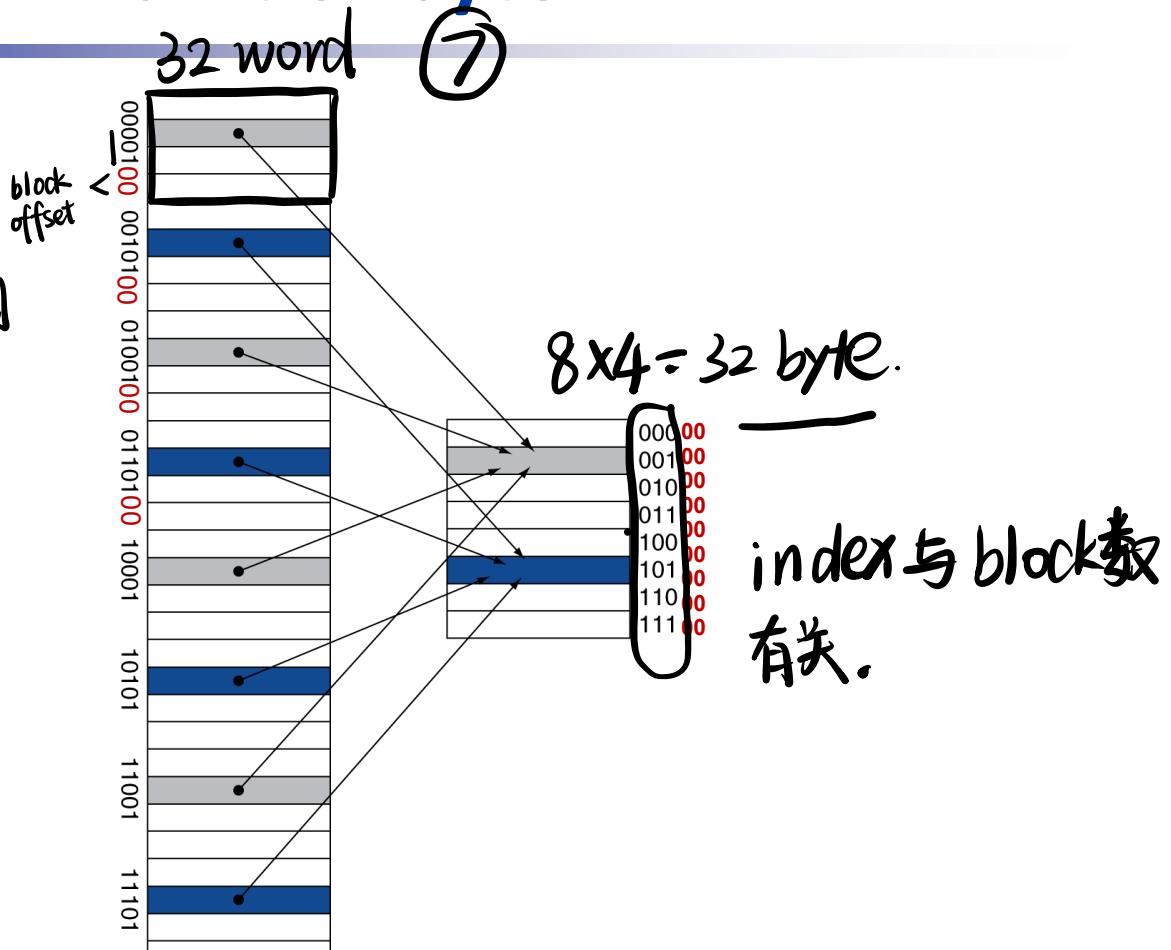
Cache Example

Word addr	Binary addr	Hit/miss	Cache block
18	10 010	Miss	010

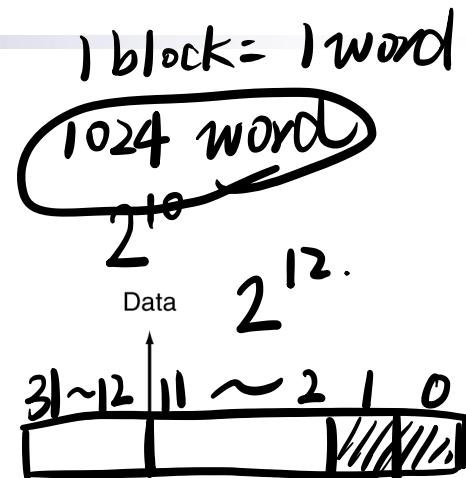
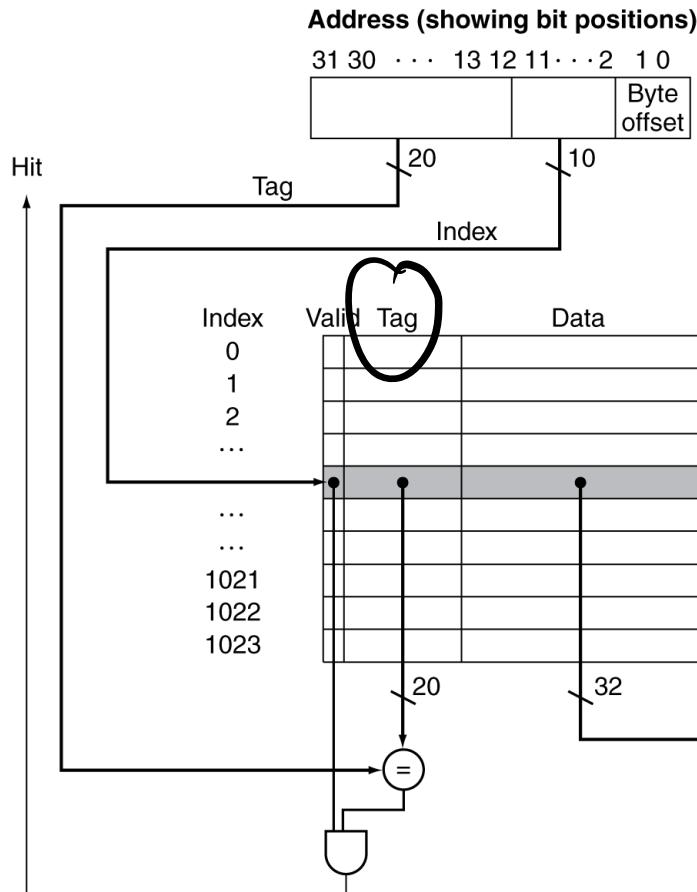
Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
010	Y	10	Mem[10010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Memory in unit of byte

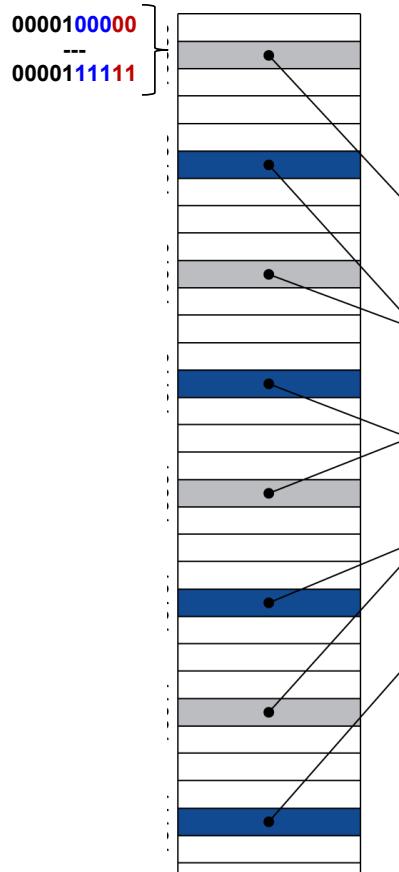
如果访问
第2



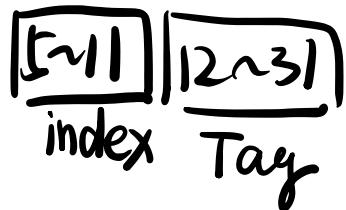
Address Subdivision



Larger Block Size



cache: 1024 word,
1 block = 8 word.
↓
2⁷ index.



Larger Block Size

- Assume:

- 32-bit address

- Direct mapped cache

- 2^n number of blocks, so n bit for index

- Block size: 2^m words, so m bit for the word within the block

- Calculate:

- Size of tag field: $32 - (n + m + 2)$

- Size of cache: $2^n * (\text{block size} + \text{tag size} + \text{valid field size})$

$$= 2^n * [2^m * 32 + (32 - n - m - 2) + 1]$$

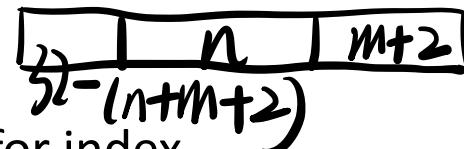
word in block

Tag



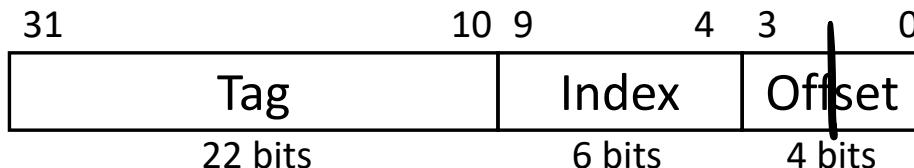
bit to 位.

number of index: # of block in the cache
 2^n
each block is 2^m word



Example: Larger Block Size

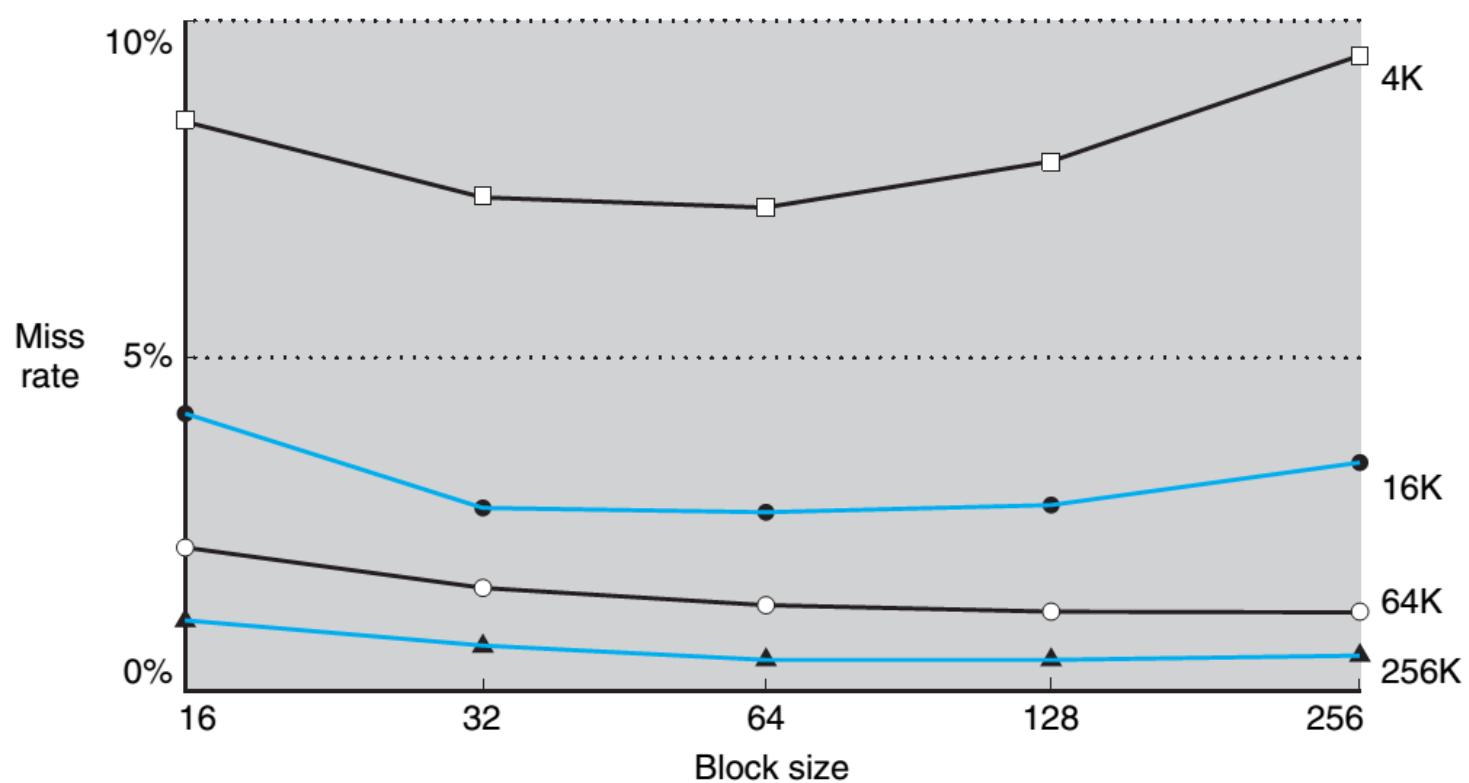
- 64 blocks, 2^4 blocks, 16 bytes/block
 - To what block number does address 1200 map?
- Block address = $\lfloor 1200/16 \rfloor = 75$
- Block number = 75 modulo 64 = 11



Block Size Considerations

- Larger blocks should reduce miss rate
 - ◆ Due to spatial locality
- But in a fixed-sized cache
 - ◆ Larger blocks \Rightarrow fewer of them
 - More competition \Rightarrow increased miss rate
 - ◆ Larger blocks \Rightarrow pollution 全部重写
- Larger miss penalty (transfer time)
 - ◆ Can override benefit of reduced miss rate
 - ◆ Early restart and critical-word-first can help

Block Size Considerations



Cache Misses

- On cache hit, CPU proceeds normally
- On cache miss
 - ◆ Read miss vs. write miss
 - ◆ Stall the CPU pipeline
 - ◆ Fetch block from next level of hierarchy
 - ◆ Instruction cache miss
 - Restart instruction fetch
 - ◆ Data cache miss
 - Complete data access

Write-Through

- On data-write hit, could just update the block in cache
 - ◆ But then cache and memory would be inconsistent
- Write through: also update memory
- But makes writes take longer
 - ◆ e.g., if base CPI = 1, 10% of instructions are stores, write to memory takes 100 cycles
 - Effective CPI = $1 + 0.1 \times 100 = 11$
- Solution: write buffer
$$1 + 10\% \times 100$$
 - ◆ Holds data waiting to be written to memory
 - ◆ CPU continues immediately
 - Only stalls on write if write buffer is already full

Write-Back

- Alternative: On data-write hit, just update the block in cache
 - ◆ Keep track of whether each block is dirty \Rightarrow dirty bit
- When a dirty block is replaced
 - ◆ Write it back to memory
 - ◆ Can use a write buffer to allow replacing block to be read first

Write Allocation

- What should happen on a write miss?
- Alternatives for write-through
 - ◆ Allocate on miss: fetch the block
 - ◆ Write around: don't fetch the block
 - Since programs often write a whole block before reading it (e.g., initialization)
- For write-back
 - ◆ Usually fetch the block

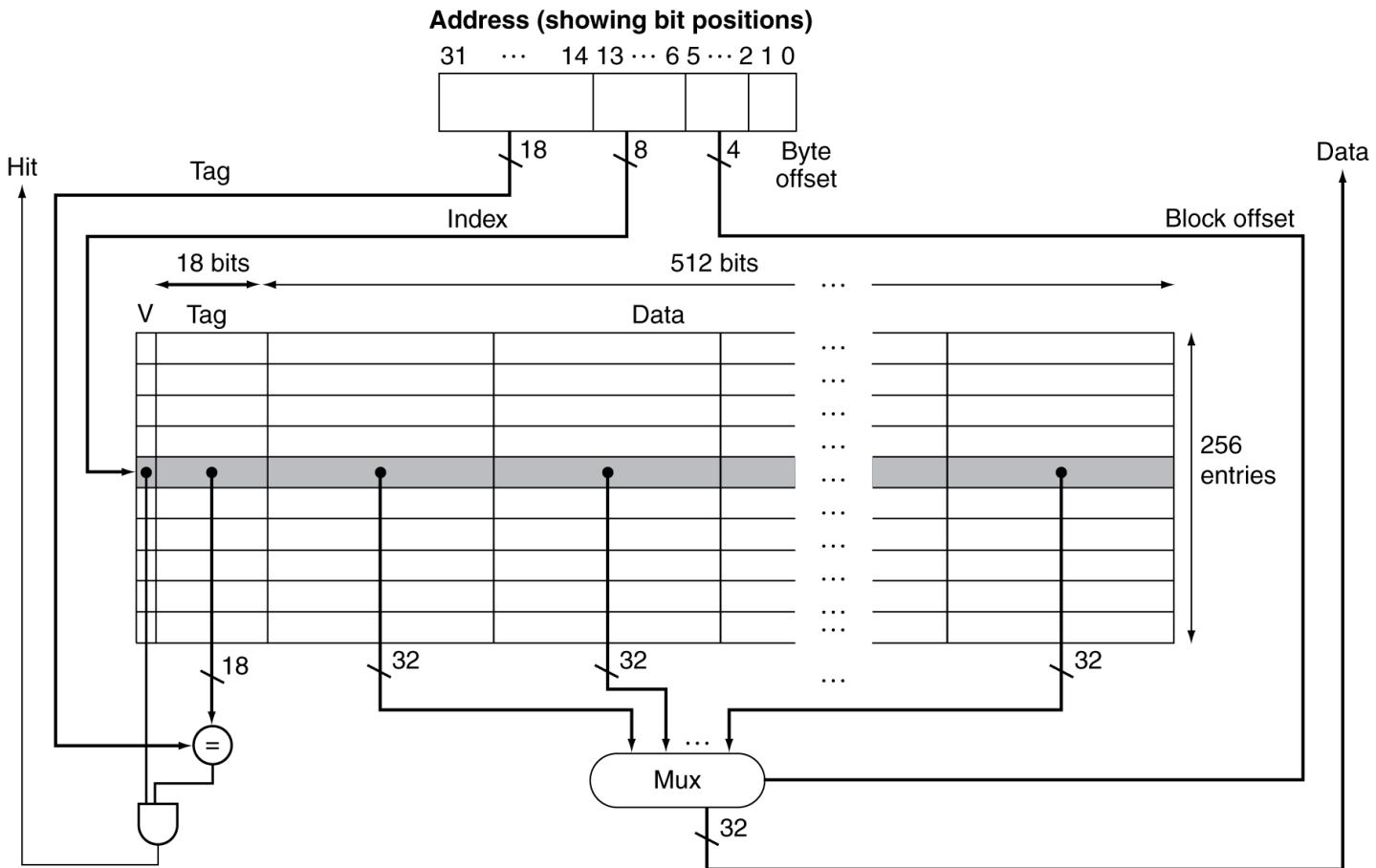
Write Policies Summary

- If that memory location is in the cache?
 - ◆ Send it to the cache
 - ◆ Should we also send it to memory right away?
(write-through policy)
 - ◆ Wait until we kick the block out *(write-back policy)*
- If it is not in the cache?
 - ◆ Allocate the line (put it in the cache)?
(write allocate policy)
 - ◆ Write it directly to memory without allocation?
(no write allocate policy)

Example: Intrinsity FastMATH

- Embedded MIPS processor
 - ◆ 12-stage pipeline
 - ◆ Instruction and data access on each cycle
- Split cache: separate I-cache and D-cache
 - ◆ Each 16KB: $256 \text{ blocks} \times 16 \text{ words/block}$
 - ◆ D-cache: write-through or write-back
- SPEC2000 miss rates
 - ◆ I-cache: 0.4%
 - ◆ D-cache: 11.4%
 - ◆ Weighted average: 3.2%

Example: Intrinsity FastMATH



Measuring Cache Performance

- Components of CPU time
 - ◆ Program execution cycles
 - Includes cache hit time
 - ◆ Memory stall cycles
 - Mainly from cache misses
- With simplifying assumptions:

Memory stall cycles

$$= \frac{\text{Memory accesses}}{\text{Program}} \times \text{Miss rate} \times \text{Miss penalty}$$

$$= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Misses}}{\text{Instruction}} \times \text{Miss penalty}$$

Cache Performance Example

- Calculate actual CPI, given that
 - ◆ I-cache miss rate = 2%
 - ◆ D-cache miss rate = 4%
 - ◆ Miss penalty = 100 cycles
 - ◆ Base CPI (ideal cache) = 2
 - ◆ Load & stores are 36% of instructions
- Miss cycles per instruction (assume N ins. In total)
 - ◆ I-cache: $N \times 0.02 \times 100/N = 2$
 - ◆ D-cache: $N \times 0.36 \times 0.04 \times 100/N = 1.44$
- Actual CPI = $2 + 2 + 1.44 = 5.44$
 - ◆ Ideal CPU is $5.44/2 = 2.72$ times faster

Average Access Time

- Hit time is also important for performance
- Average memory access time (AMAT)
 - ◆ $\text{AMAT} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$
- Example
 - ◆ CPU with 1ns clock, hit time = 1 cycle, miss penalty = 20 cycles, I-cache miss rate = 5%
 - ◆ $\text{AMAT} = 1 + 0.05 \times 20 = 2\text{ns}$
 - 2 cycles per instruction

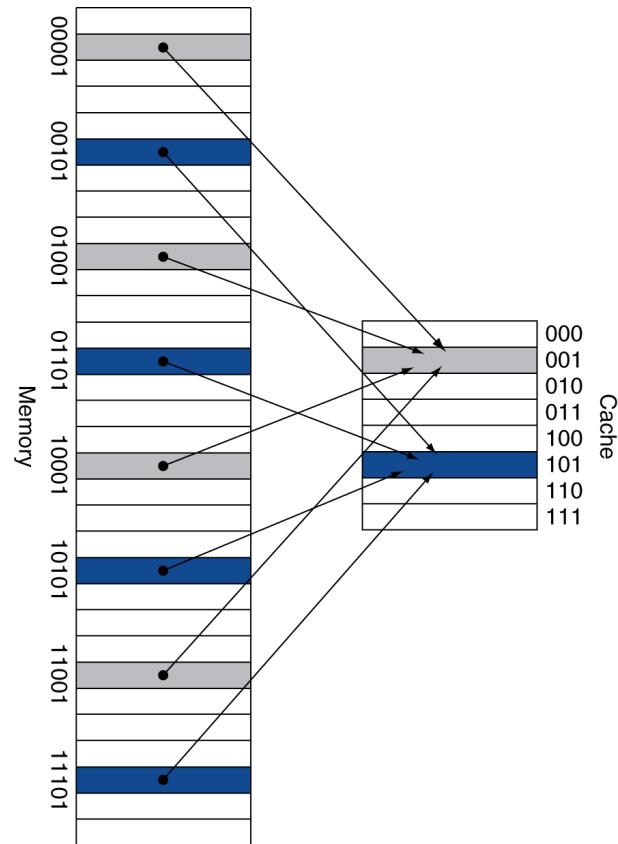
$$1 + 5\% \times 20 = 2\text{ns}$$

Performance Summary

- When CPU performance increased
 - ◆ Miss penalty becomes more significant
 - ◆ CPI=2, Miss=3.44, % of memory stall: $3.44/5.44=63\%$
 - ◆ CPI=1, Miss=3.44, % of memory stall: $3.44/4.44=77\%$
- Decreasing base CPI
 - ◆ Greater proportion of time spent on memory stalls
- Increasing clock rate
 - ◆ Memory stalls account for more CPU cycles
- Can't neglect cache behavior when evaluating system performance

Recall: Direct Mapped Cache

- Location determined by address
- Direct mapped: only one choice
 - ◆ Capacity of cache is not fully exploited
 - ◆ Miss rate is high



Cache Example

Word addr	Binary addr	Hit/miss	Cache block
16	10 000	Miss	000
3	00 011	Miss	011
16	10 000	Hit	000

Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
010	Y	11	Mem[11010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

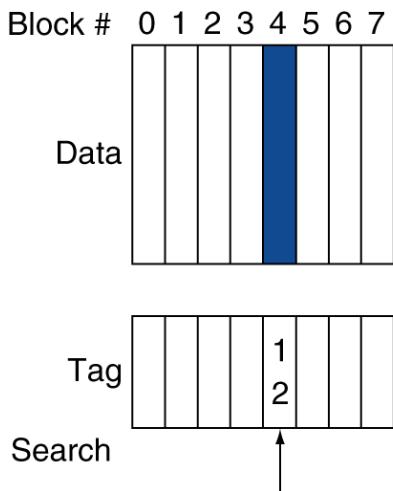
Cache Example

Word addr	Binary addr	Hit/miss	Cache block
18	10 010	Miss	010

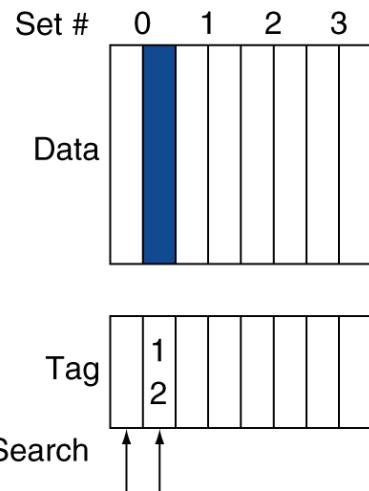
Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
010	Y	10	Mem[10010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Associative Cache Example

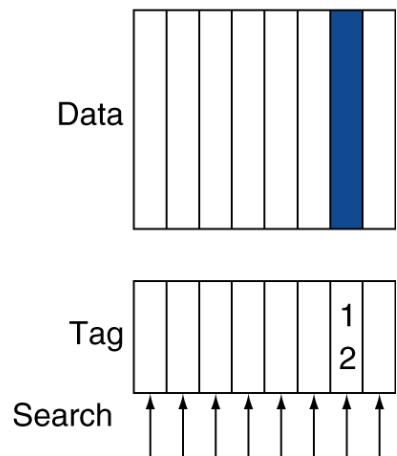
Direct mapped



Set associative



Fully associative



Associative Caches

- Fully associative
 - ◆ Allow a given block to go in any cache entry
 - ◆ Requires all entries to be searched at once
 - ◆ Comparator per entry (expensive)
- n -way set associative
 - ◆ Each set contains n entries
 - ◆ Block number determines which set
 - $(\text{Block number}) \bmod (\#\text{Sets in cache})$
 - ◆ Search all entries in a given set at once
 - ◆ n comparators (less expensive)

Spectrum of Associativity

- For a cache with 8 blocks

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				

Four-way set associative

Set	Tag	Data	Tag	Data	Tag	Data	Tag	Data
0								
1								

Eight-way set associative (fully associative)

Associativity Example

- Compare 4-block caches
 - ◆ Direct mapped, 2-way set associative, fully associative
 - ◆ Block access sequence: [0, 8, 0, 6, 8]
- Direct mapped

Block address	Cache index	Hit/miss	Cache content after access			
			0	1	2	3
0	0	miss	Mem[0]			
8	0	miss	Mem[8]			
0	0	miss	Mem[0]			
6	2	miss	Mem[0]		Mem[6]	
8	0	miss	Mem[8]		Mem[6]	

Associativity Example

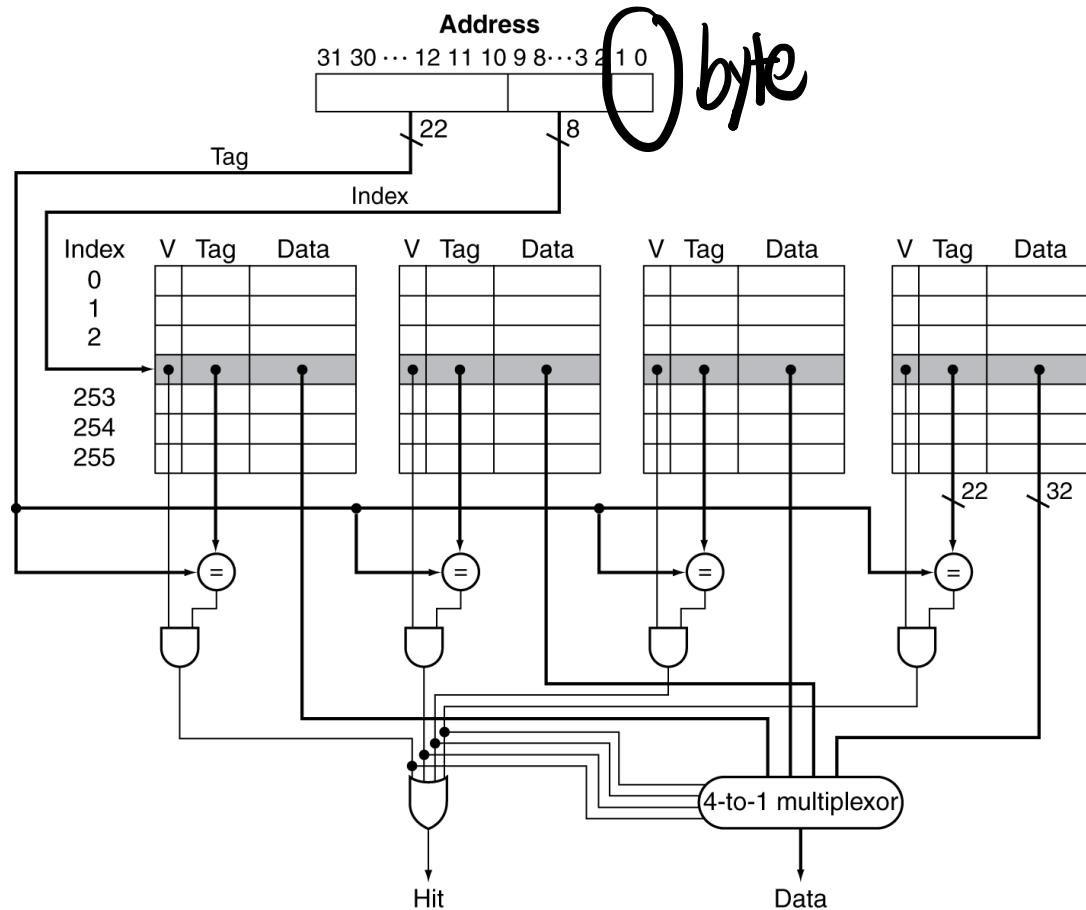
- 2-way set associative

Block address	Cache index	Hit/miss	Cache content after access			
			Set 0		Set 1	
0	0	miss	Mem[0]			
8	0	miss	Mem[0]	Mem[8]		
0	0	hit	Mem[0]	Mem[8]		
6	0	miss	Mem[0]	Mem[6]		
8	0	miss	Mem[8]	Mem[6]		

- Fully associative

Block address		Hit/miss	Cache content after access			
			Mem[0]	Mem[8]	Mem[6]	
0		miss	Mem[0]			
8		miss	Mem[0]	Mem[8]		
0		hit	Mem[0]	Mem[8]		
6		miss	Mem[0]	Mem[8]	Mem[6]	
8		hit	Mem[0]	Mem[8]	Mem[6]	

Set Associative Cache Organization



How Much Associativity

- Increased associativity decreases miss rate
 - ◆ But with diminishing returns
- Simulation of a system with 64KB D-cache, 16-word blocks, SPEC2000
 - ◆ 1-way: 10.3%
 - ◆ 2-way: 8.6%
 - ◆ 4-way: 8.3%
 - ◆ 8-way: 8.1%

Replacement Policy

- Direct mapped: no choice
- Set associative
 - ◆ Prefer non-valid entry, if there is one
 - ◆ Otherwise, choose among entries in the set
- Least-recently used (LRU)
 - ◆ Choose the one unused for the longest time
 - Simple for 2-way, manageable for 4-way, too hard beyond that
- Random
 - ◆ Gives approximately the same performance as LRU for high associativity

Multilevel Caches

- Primary cache attached to CPU
 - ◆ Small, but fast
- Level-2 cache services misses from primary cache
 - ◆ Larger, slower, but still faster than main memory
- Main memory services L-2 cache misses
- Some high-end systems include L-3 cache

Multilevel Cache Example

- Given

- CPU base CPI = 1, clock rate = 4GHz
- Miss rate/instruction = 2%
- Main memory access time = 100ns

$$1 + 2\% \times 400 = 9$$

- With just primary cache

- Miss penalty = $100\text{ns}/0.25\text{ns} = 400 \text{ cycles}$
- Effective CPI = $1 + 0.02 \times 400 = 9$

Cycle per Instruction

Example (cont.)

- Now add L-2 cache
 - ◆ Access time = 5ns
 - ◆ Global miss rate to main memory = 0.5%
- Primary miss with L-2 hit
 - ◆ Penalty = $5\text{ns}/0.25\text{ns} = 20 \text{ cycles}$
- Primary miss with L-2 miss
 - ◆ Extra penalty = 400 cycles
- CPI = $1 + 0.02 \times 20 + 0.005 \times 400 = 3.4$
- Performance ratio = $9/3.4 = 2.6$

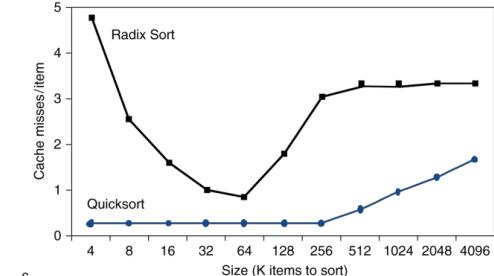
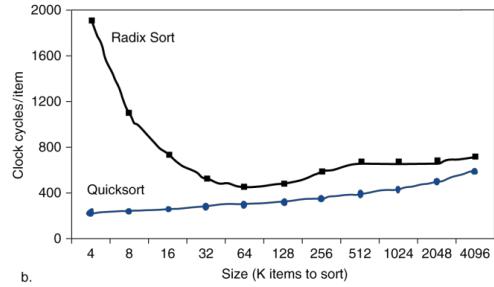
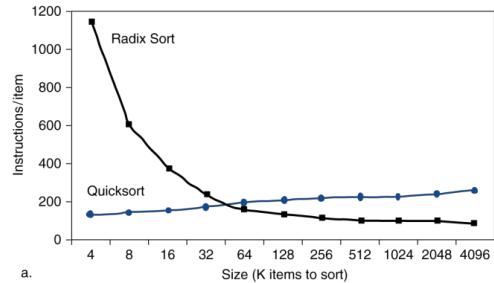
$$1 + 2\% \times 20 + 0.005 \times 400$$

Multilevel Cache Considerations

- Primary cache
 - ◆ Focus on minimal hit time
- L-2 cache
 - ◆ Focus on low miss rate to avoid main memory access
 - ◆ Hit time has less overall impact
- Results
 - ◆ L-1 cache usually smaller than a single cache
 - ◆ L-1 block size smaller than L-2 block size

Interactions with Software

- Compare two algorithms:
Radix sort & Quicksort
- When size is large,
 - ◆ Radix sort has less instructions
 - ◆ But quicksort has less clock cycles
 - ◆ Because miss rate of radix sort is higher



Software Optimization via Blocking

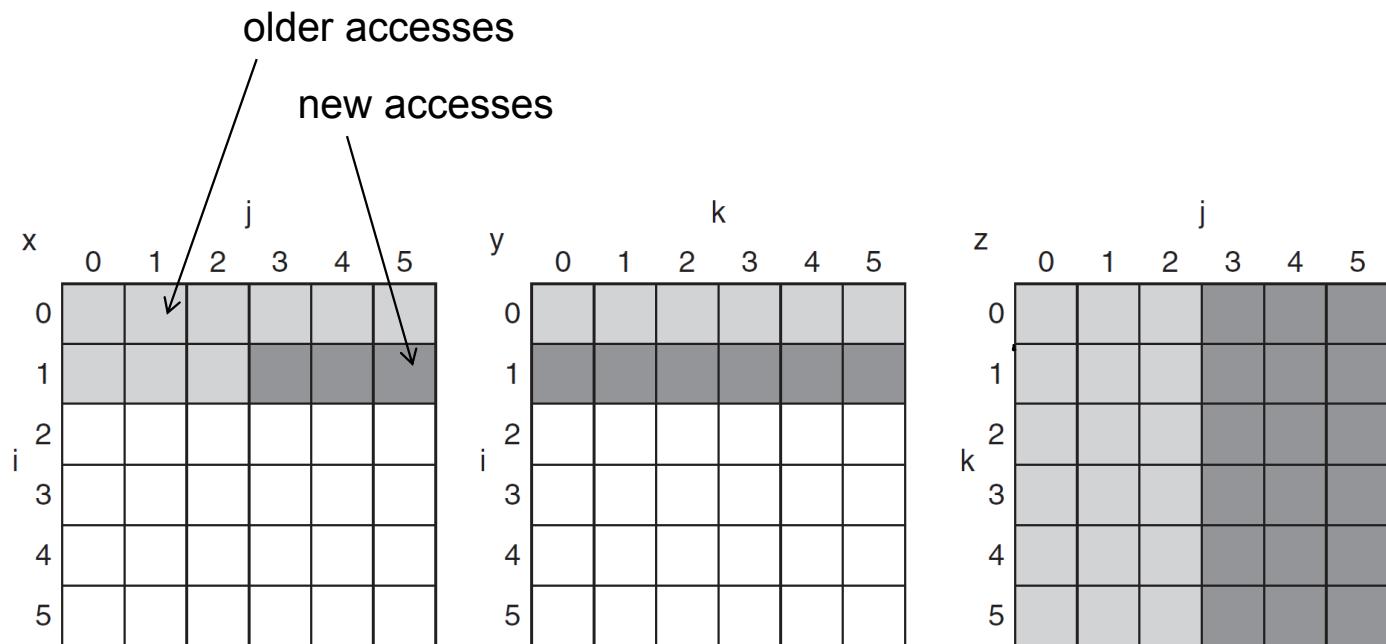
- Goal: maximize accesses to data before it is replaced
- Consider inner loops of DGEMM:

```
for (int j = 0; j < n; ++j)
{
    double cij = C[i+j*n];
    for( int k = 0; k < n; k++ )
        cij += A[i+k*n] * B[k+j*n];
    C[i+j*n] = cij;
}
```

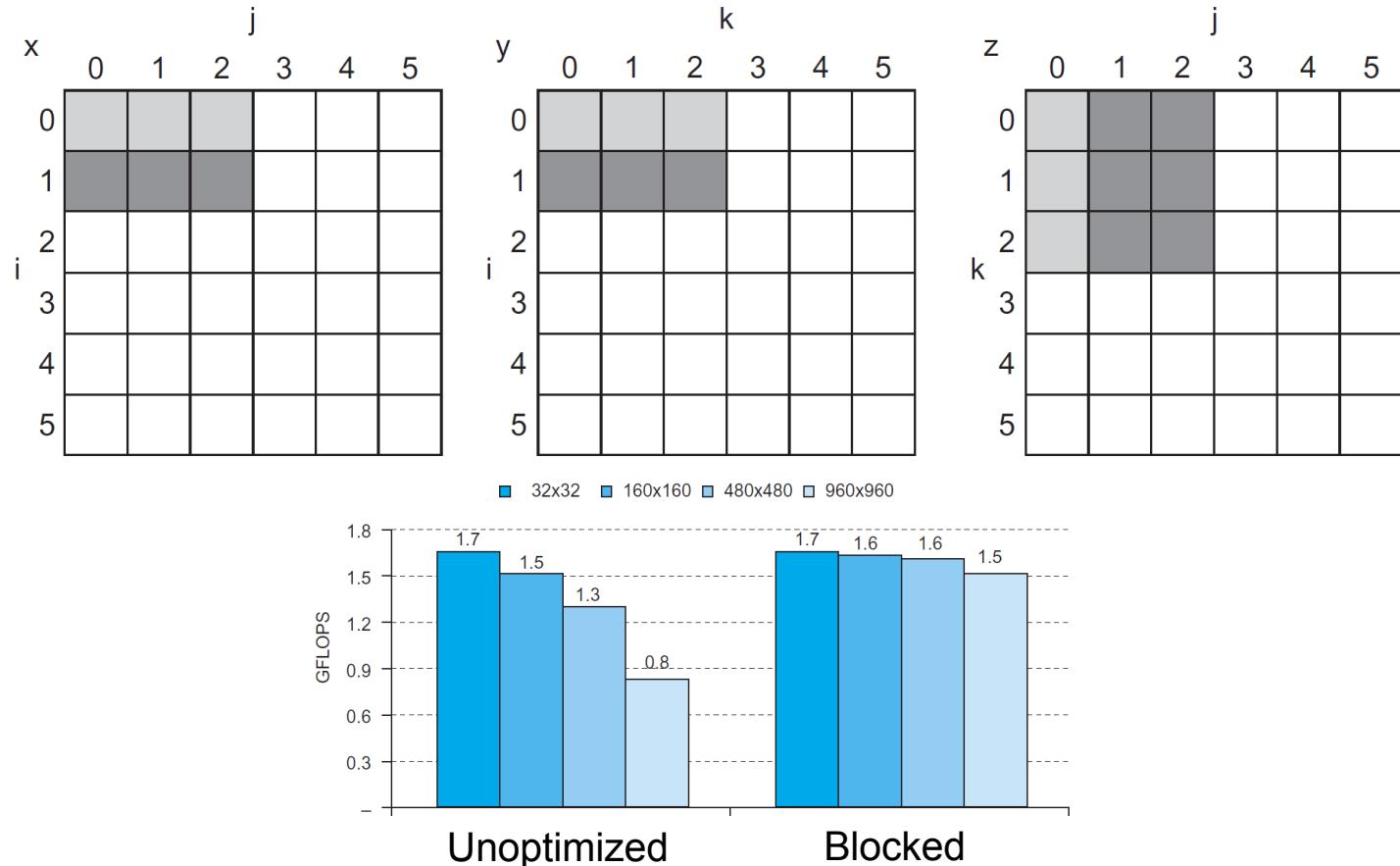


DGEMM Access Pattern

- C, A, and B arrays



Blocked DGEMM Access Pattern



Homework

- Exercise 5.3 5.6