

## **Lecture 12**

### **Memory Performance and Dependable Memory Hierarchy**

# 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

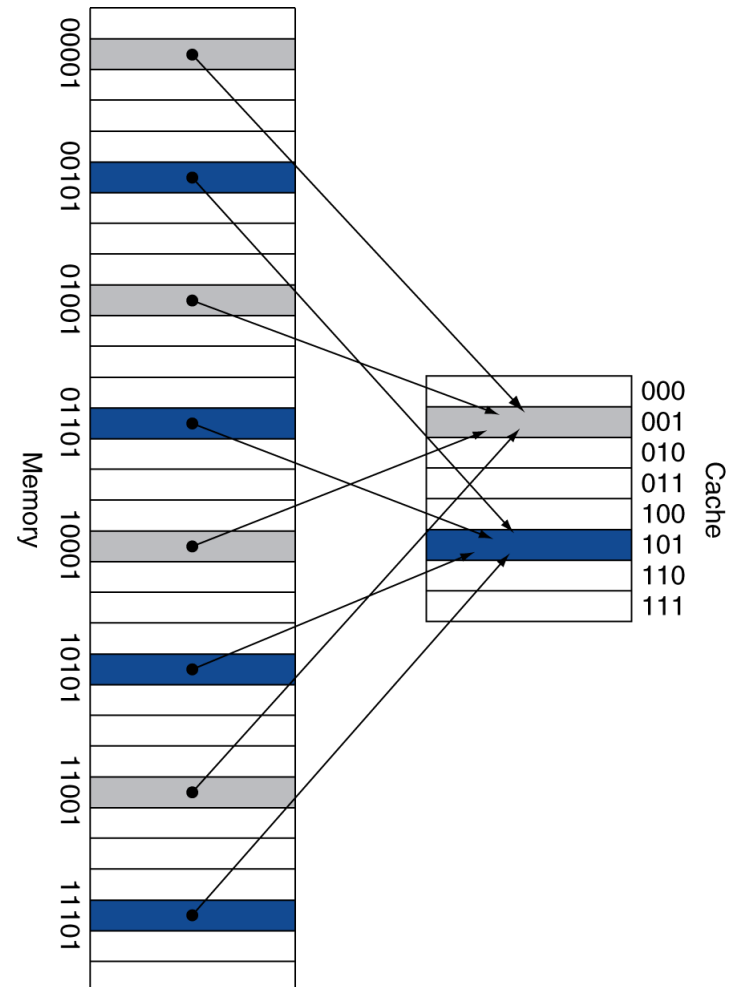
- Hit time is also important for performance
- Average memory access time (AMAT)
  - ◆  $\text{AMAT} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$
- Example
  - ◆ CPU with 1ns clock, hit time = 1 cycle, miss penalty = 20 cycles, l-cache miss rate = 5%
  - ◆  $\text{AMAT} = 1 + 0.05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

# Performance Summary

- When CPU performance increased
  - ◆ Miss penalty becomes more significant
  - ◆  $\text{CPI}=2$ ,  $\text{Miss}=3.44$ , % of memory stall:  $3.44/5.44=63\%$
  - ◆  $\text{CPI}=1$ ,  $\text{Miss}=3.44$ , % 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
<b>000</b>	<b>Y</b>	<b>10</b>	<b>Mem[10000]</b>
001	N		
010	Y	11	Mem[11010]
<b>011</b>	<b>Y</b>	<b>00</b>	<b>Mem[00011]</b>
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

# Cache Example

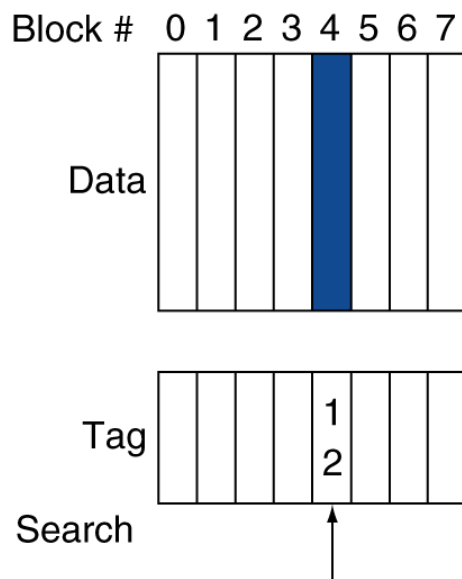
Word addr	Binary addr	Hit/miss	Cache block
18	10 010	Miss	010

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



# Associative Cache Example

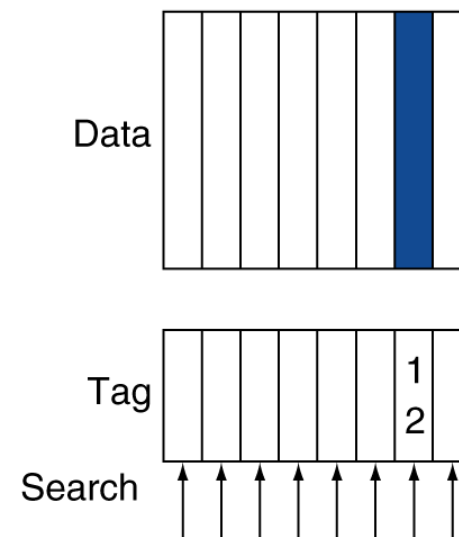
**Direct mapped**



**Set associative**



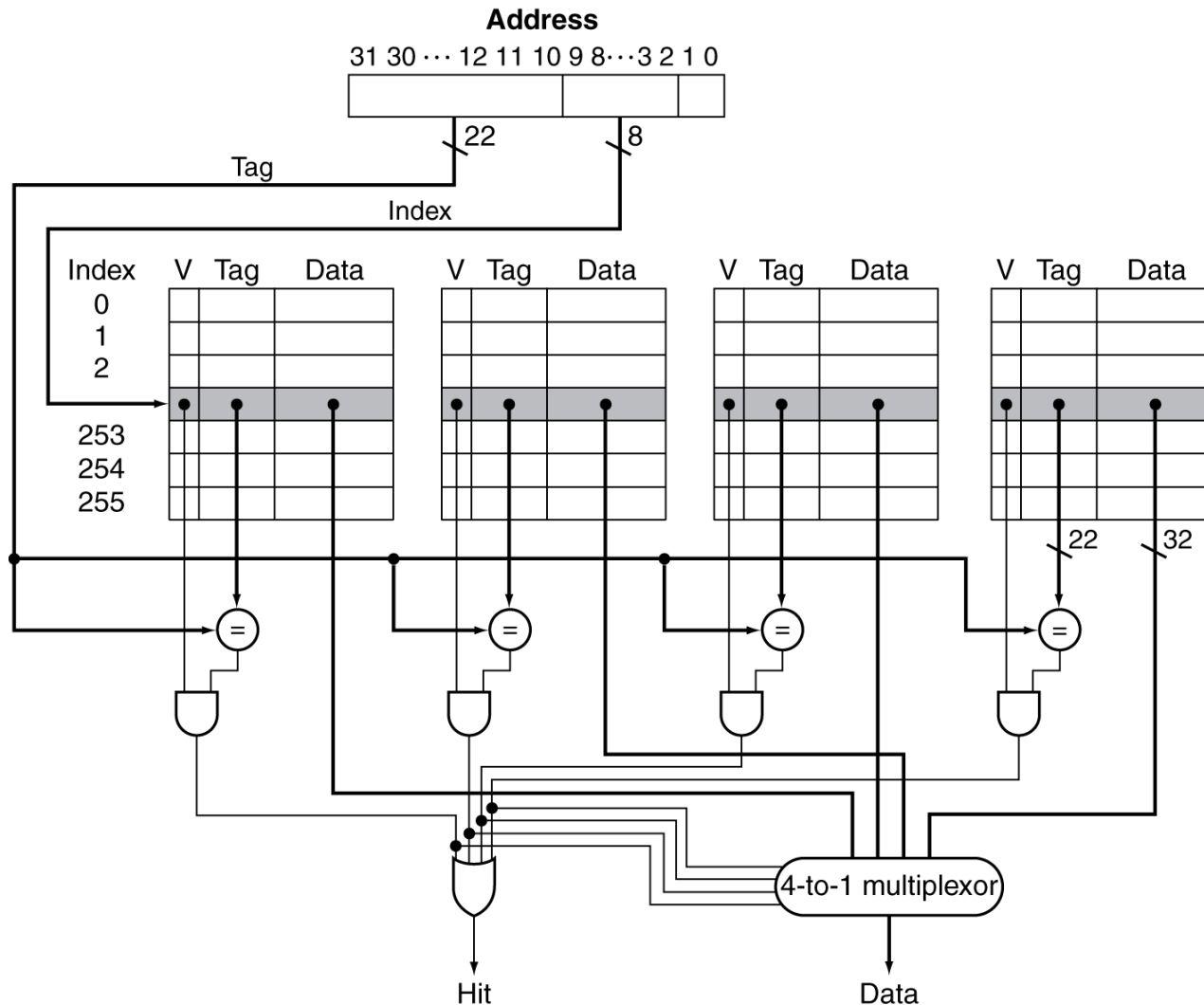
**Fully associative**



# 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
    - (Block number) modulo (#Sets in cache)
  - ◆ Search all entries in a given set at once
  - ◆  $n$  comparators (less expensive)

# Set Associative Cache Organization



- For a cache with 8 blocks

0		
1		
2		
3		
4		
5		
6		
7		

Set	Tag	Data	Tag	Data
0				
1				
2				
3				

Set	Tag	Data	Tag	Data	Tag	Data	Tag	Data
0								
1								

[illegible]

# 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]		

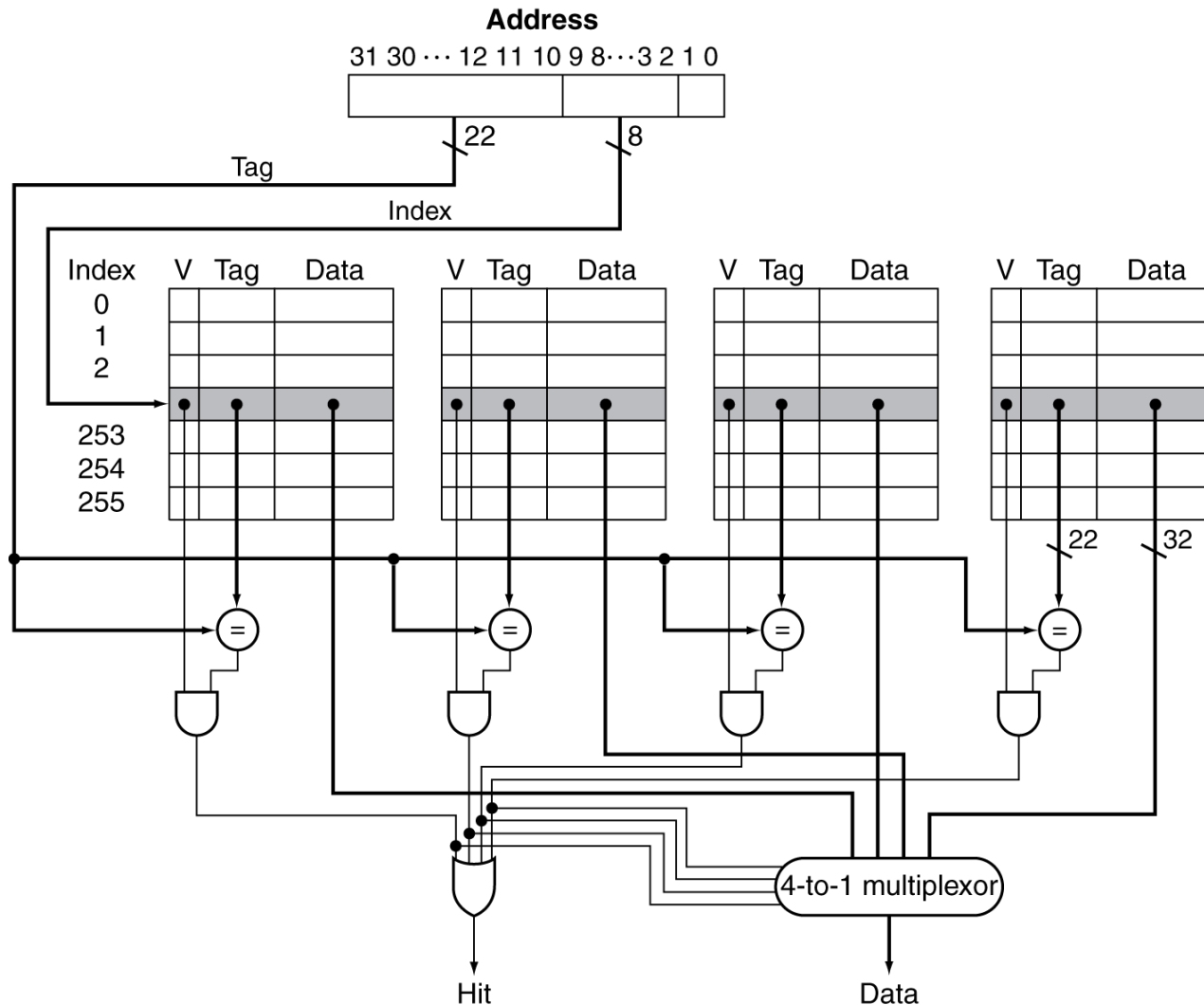
## ■ Fully associative

Block address		Hit/miss	Cache content after access			
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]	

# 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%

# Set Associative Cache Organization





# 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

- With just primary cache

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

## Example (cont.)

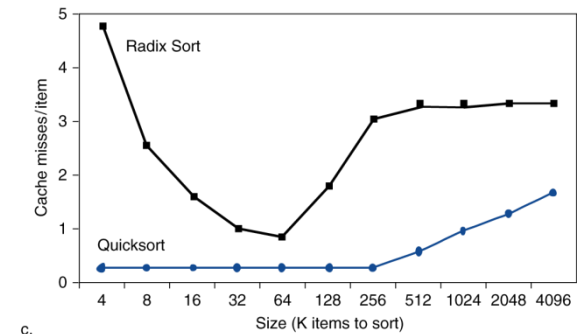
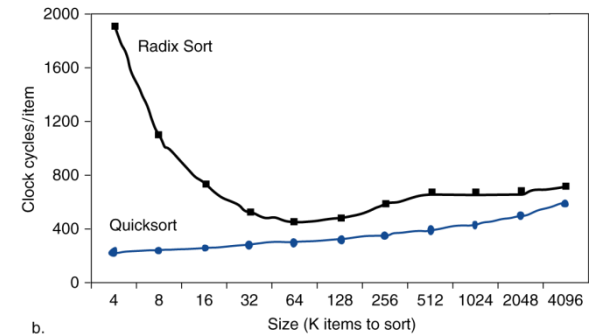
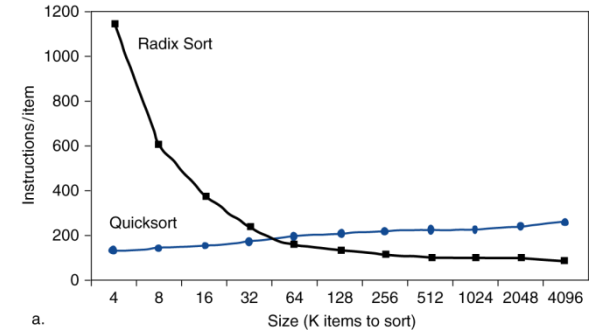
- Now add L-2 cache
  - ◆ Access time = 5ns
  - ◆ Global miss rate to main memory = 0.5%
- Primary miss with L-2 hit
  - ◆ Penalty =  $5\text{ns} / 0.25\text{ns} = 20$  cycles
- Primary miss with L-2 miss
  - ◆ Extra penalty = 400 cycles
- $\text{CPI} = 1 + 0.02 \times 20 + 0.005 \times 400 = 3.4$
-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

# 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



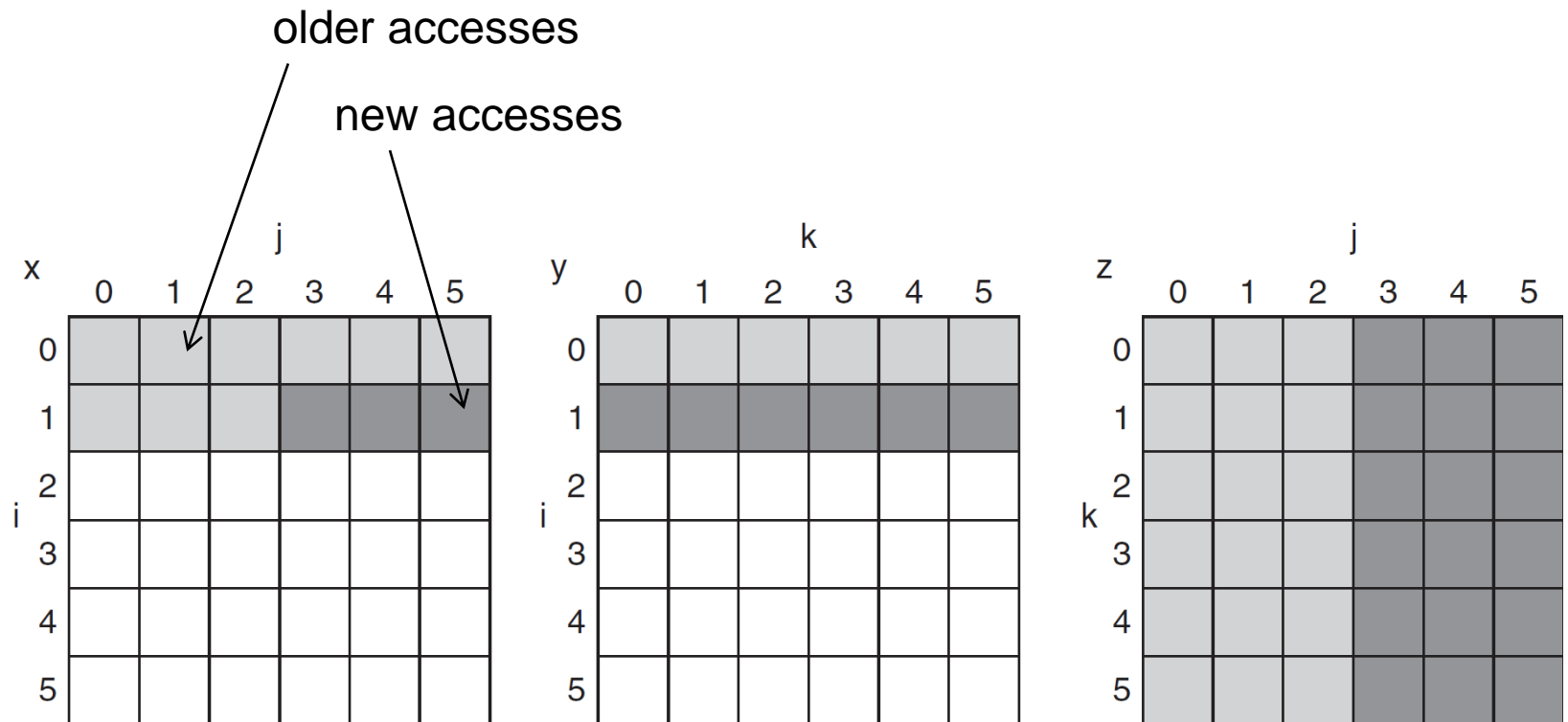
# Software Optimization via Blocking

- 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;
}
```

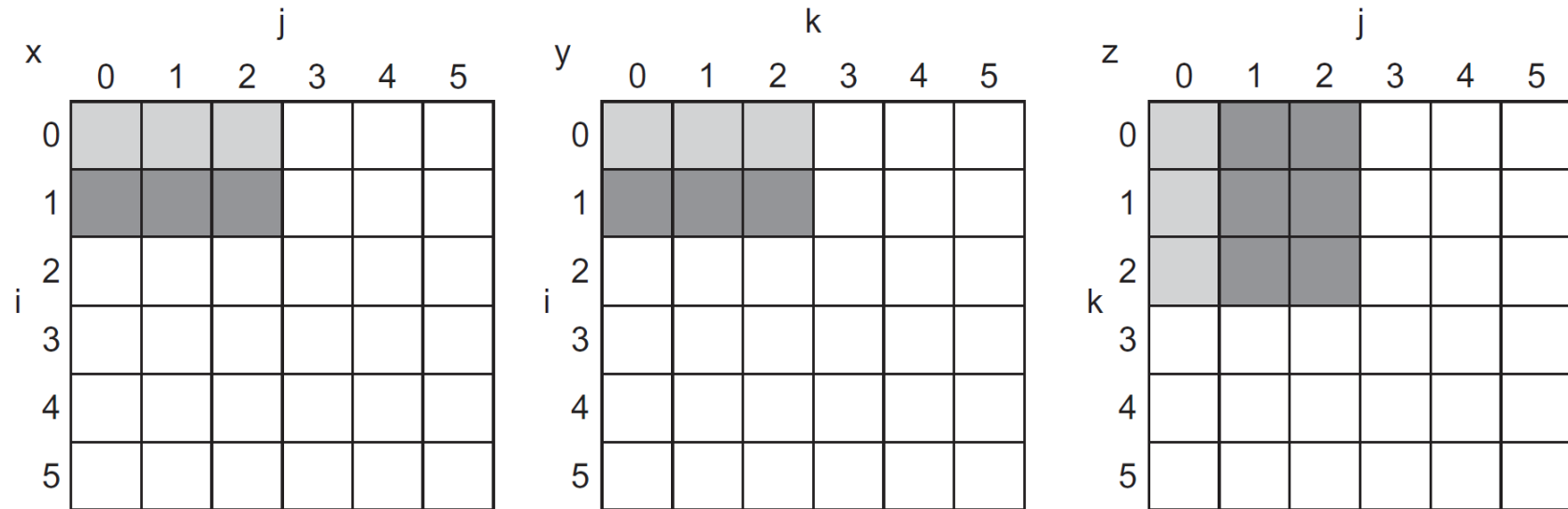
# DGEMM Access Pattern

- C, A, and B arrays

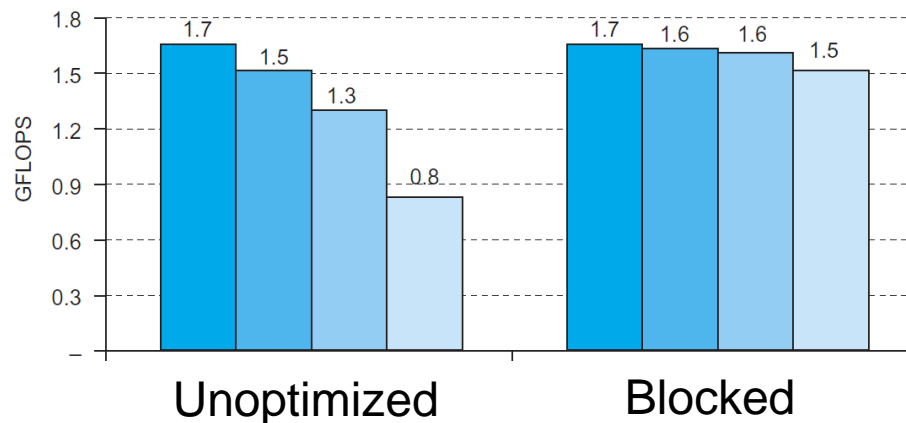




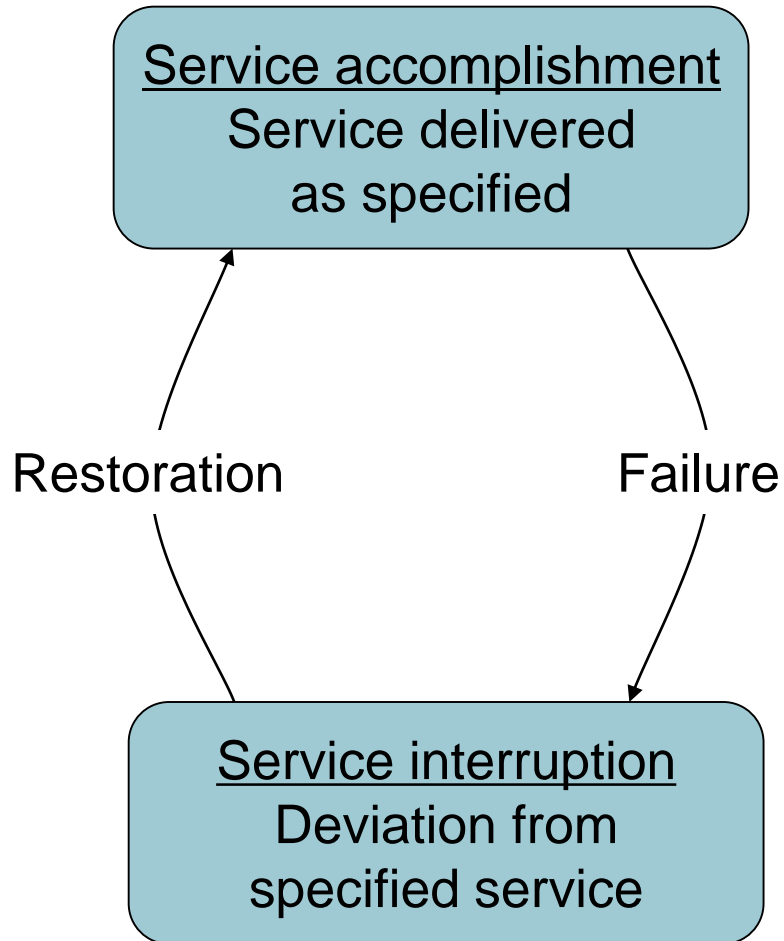
# Blocked DGEMM Access Pattern



■ 32x32 ■ 160x160 ■ 480x480 ■ 960x960



# Dependability



- Fault: failure of a component
  - ◆ May or may not lead to system failure

# Dependability Measures

- Reliability: mean time to failure (MTTF)
- Service interruption: mean time to repair (MTTR)
- Mean time between failures
  - ◆  $MTBF = MTTF + MTTR$
- Availability =  $MTTF / (MTTF + MTTR)$
- Improving Availability
  - ◆ Increase MTTF: fault avoidance, fault tolerance, fault forecasting
  - ◆ Reduce MTTR: fault detection, fault diagnosis and fault repair

# The Hamming SEC Code

- Hamming distance
  - ◆ Number of bits that are different between two bit patterns
- Minimum distance = 2 provides single bit error detection
  - ◆ E.g. parity code
- Minimum distance = 3 provides single error correction, 2 bit error detection

# Encoding SEC

- To calculate Hamming code:
  - ◆ Number bits from 1 on the left
  - ◆ All bit positions that are a power 2 are parity bits
  - ◆ Each parity bit checks certain data bits:

Bit position		1	2	3	4	5	6	7	8	9	10	11	12
Encoded data bits		p1	p2	d1	p4	d2	d3	d4	p8	d5	d6	d7	d8
Parity bit coverage	p1	X		X		X		X		X		X	
	p2		X	X			X	X			X	X	
	p4				X	X	X	X					X
	p8								X	X	X	X	X

# Decoding SEC

- Value of parity bits indicates which bits are in error
  - ◆ Use numbering from encoding procedure
  - ◆ E.g.
    - Parity bits = 0000 indicates no error
    - Parity bits = 1010 indicates bit 10 was flipped

# SEC/DED Code

- Add an additional parity bit for the whole word ( $p_n$ )
- Make Hamming distance = 4
- Decoding:
  - ◆ Let  $H$  = SEC parity bits
    - $H$  even,  $p_n$  even, no error
    - $H$  odd,  $p_n$  odd, correctable single bit error
    - $H$  even,  $p_n$  odd, error in  $p_n$  bit
    - $H$  odd,  $p_n$  even, double error occurred
- Note: ECC DRAM uses SEC/DED with 8 bits protecting each 64 bits

# Summary

- Cache Performance
  - ◆ Mainly depends on miss rate and miss penalty
- To improve cache performance:
  - ◆ Fully associative cache
  - ◆ Set-associative cache
  - ◆ Replacement policy
  - ◆ Multilevel cache
- Dependability
  - ◆ MTTF, MTTR, reliability, availability
  - ◆ Hamming code: SEC/DED code



# Homework

---

- Exercise 5.6, 5.9.