

The Memory Hierarchy : Memory Hierarchy - Cache

Storage Trends

SRAM

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	19,200	2,900	320	256	100	75	60	320
access (ns)	300	150	35	15	3	2	1.5	200

DRAM → main mem

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	8,000	880	100	30	1	0.1	0.06	130,000
access (ns)	375	200	100	70	60	50	40	9
typical size (MB)	0.064	0.256	4	16	64	2,000	8,000	125,000

Disk

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	500	100	8	0.30	0.01	0.005	0.0003	1,600,000
access (ms)	87	75	28	10	8	4	3	29
typical size (MB)	1	10	160	1,000	20,000	160,000	1,500,000	1,500,000

CPU Clock Rates

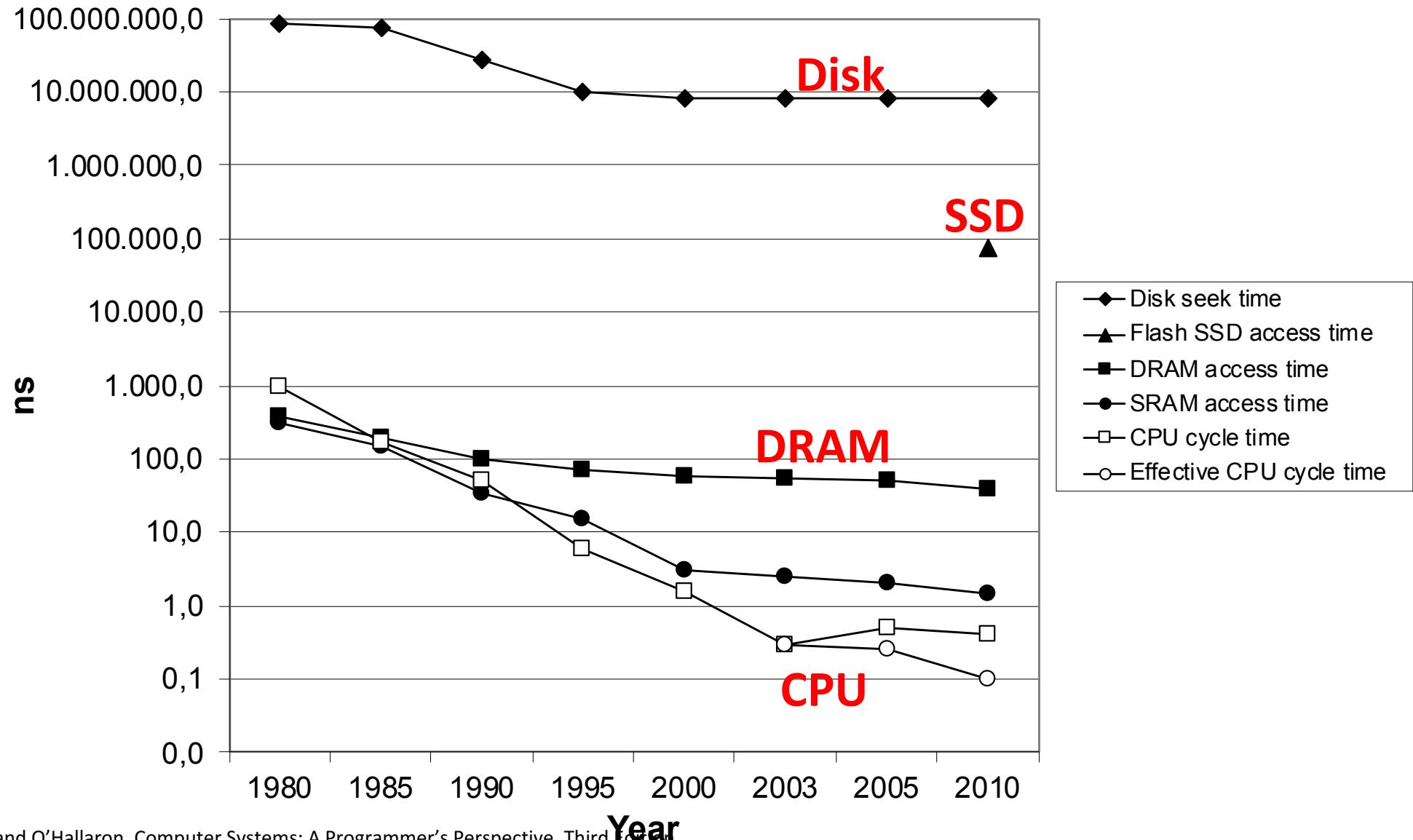
Inflection point in computer history
when designers hit the “Power Wall”



	1980	1990	1995	2000	2003	2005	2010	2010:1980
CPU	8080	386	Pentium	P-III	P-4	Core 2	Core i7	---
Clock rate (MHz)	1	20	150	600	3300	2000	2500	2500
Cycle time (ns)	1000	50	6	1.6	0.3	0.50	0.4	2500
Cores	1	1	1	1	1	2	4	4
Effective cycle time (ns)	1000	50	6	1.6	0.3	0.25	0.1	10,000

The CPU-Memory Gap

The gap widens between DRAM, disk, and CPU speeds.



Locality to the Rescue!

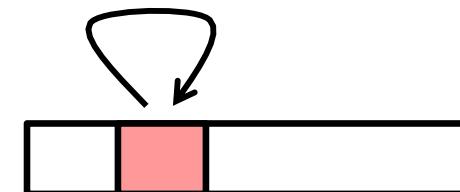
The key to bridging this CPU-Memory gap is a fundamental property of computer programs known as **locality**

Locality

- **Principle of Locality:** Programs tend to use data and instructions with addresses near or equal to those they have used recently

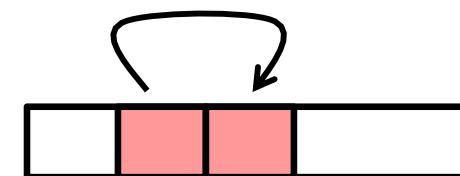
- **Temporal locality:**

- Recently referenced items are likely to be referenced again in the near future



- **Spatial locality:**

- Items with nearby addresses tend to be referenced close together in time



Locality Example

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

■ Data references

- Reference array elements in succession (stride-1 reference pattern).
- Reference variable `sum` each iteration.

Spatial locality

Temporal locality

■ Instruction references

- Reference instructions in sequence.
- Cycle through loop repeatedly.

Spatial locality

Temporal locality

Qualitative Estimates of Locality

- **Claim:** Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer.
- **Question:** Does this function have good locality with respect to array a?

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}
```



Locality Example

- **Question:** Does this function have good locality with respect to array a? *no → it doesn't have spatial locality*

```
int sum_array_cols(int a[M][N])
{
    int i, j, sum = 0;

    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}
```

Locality Example

- **Question:** Can you permute the loops so that the function scans the 3-d array **a** with a stride-1 reference pattern (and thus has good spatial locality)?

```
int sum_array_3d(int a[M] [N] [N])
{
    int i, j, k, sum = 0;

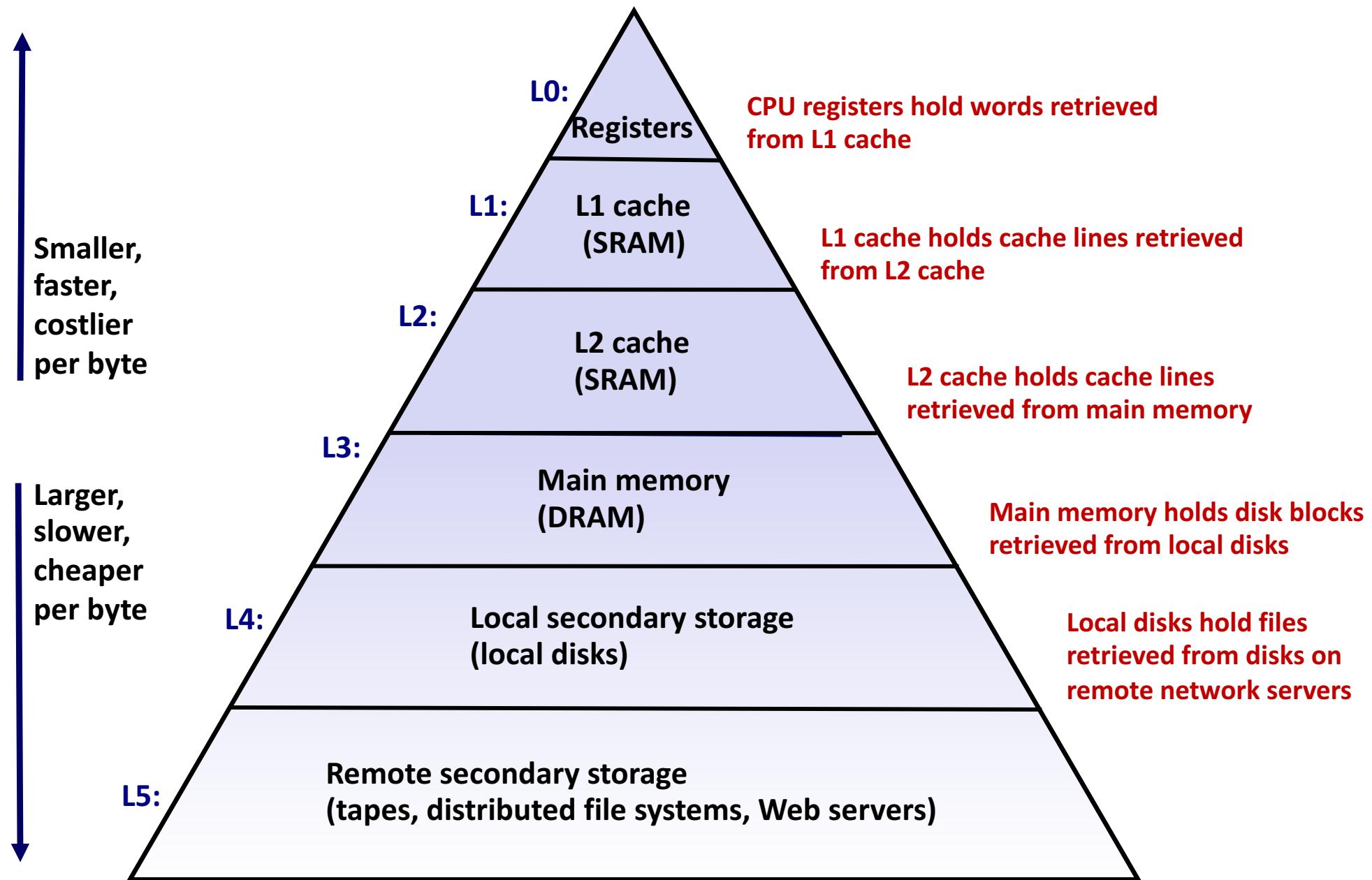
    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            for (k = 0; k < N; k++)
                sum += a[k] [i] [j];

    return sum;
}
```

Memory Hierarchies

- Some fundamental and enduring properties of hardware and software:
 - Fast storage technologies cost more per byte, have less capacity, and require more power (heat!).
 - The gap between CPU and main memory speed is widening.
 - Well-written programs tend to exhibit good locality.
- These fundamental properties complement each other beautifully.
- They suggest an approach for organizing memory and storage systems known as a **memory hierarchy**.

An Example Memory Hierarchy



Caches

- ***Cache:*** A smaller, faster storage device that acts as a staging area for a subset of the data in a larger, slower device.
- **Fundamental idea of a memory hierarchy:**
 - For each k , the faster, smaller device at level k serves as a cache for the larger, slower device at level $k+1$.
- **Why do memory hierarchies work?**
 - Because of locality, programs tend to access the data at level k more often than they access the data at level $k+1$.
 - Thus, the storage at level $k+1$ can be slower, and thus larger and cheaper per bit.
- ***Big Idea:*** The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.

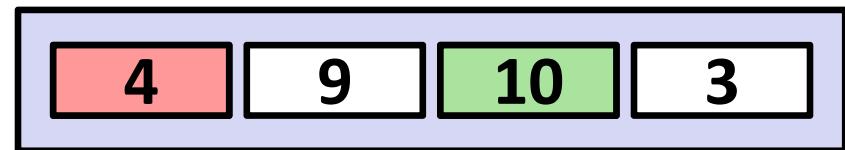
General Cache Concepts

Block: min size of bytes copied

cache hit : check for the content, if you find it in that block it is a hit

X
miss

Cache

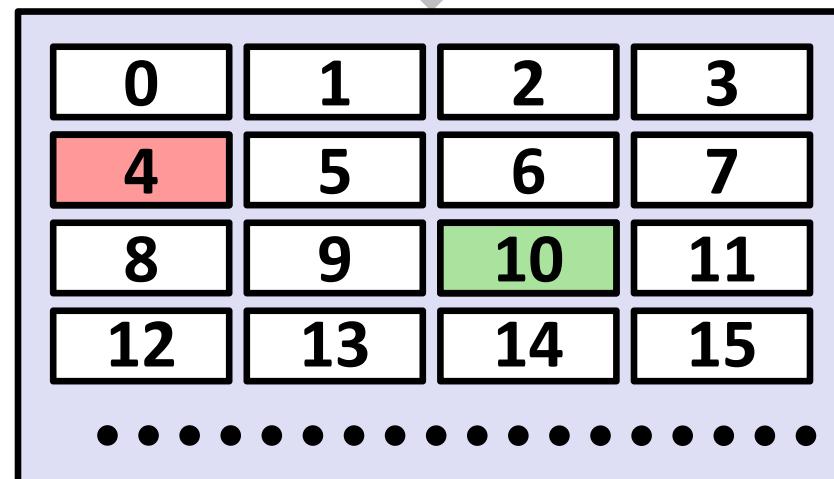


Smaller, faster, more expensive memory caches a subset of the blocks



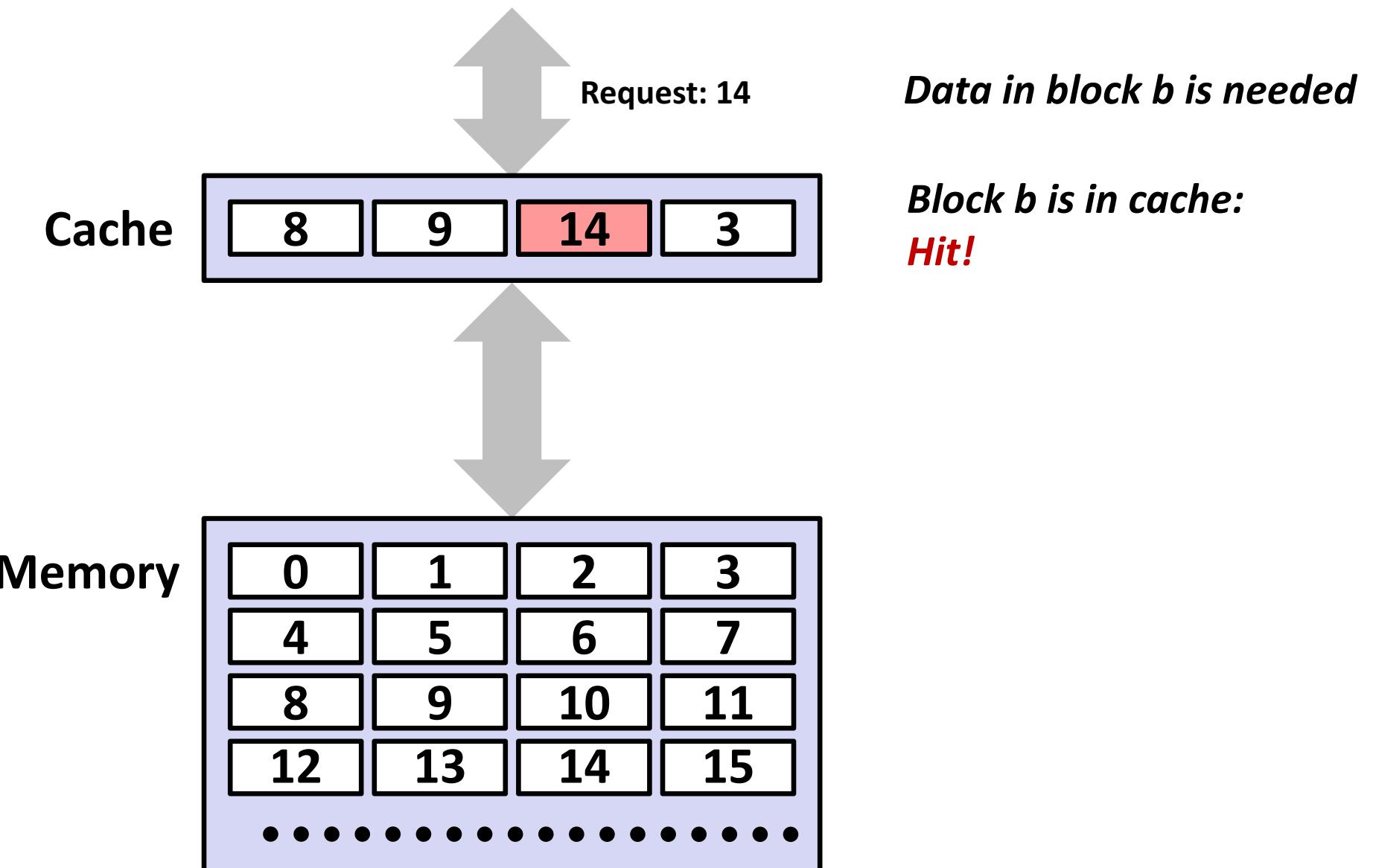
Data is copied in block-sized transfer units

Memory

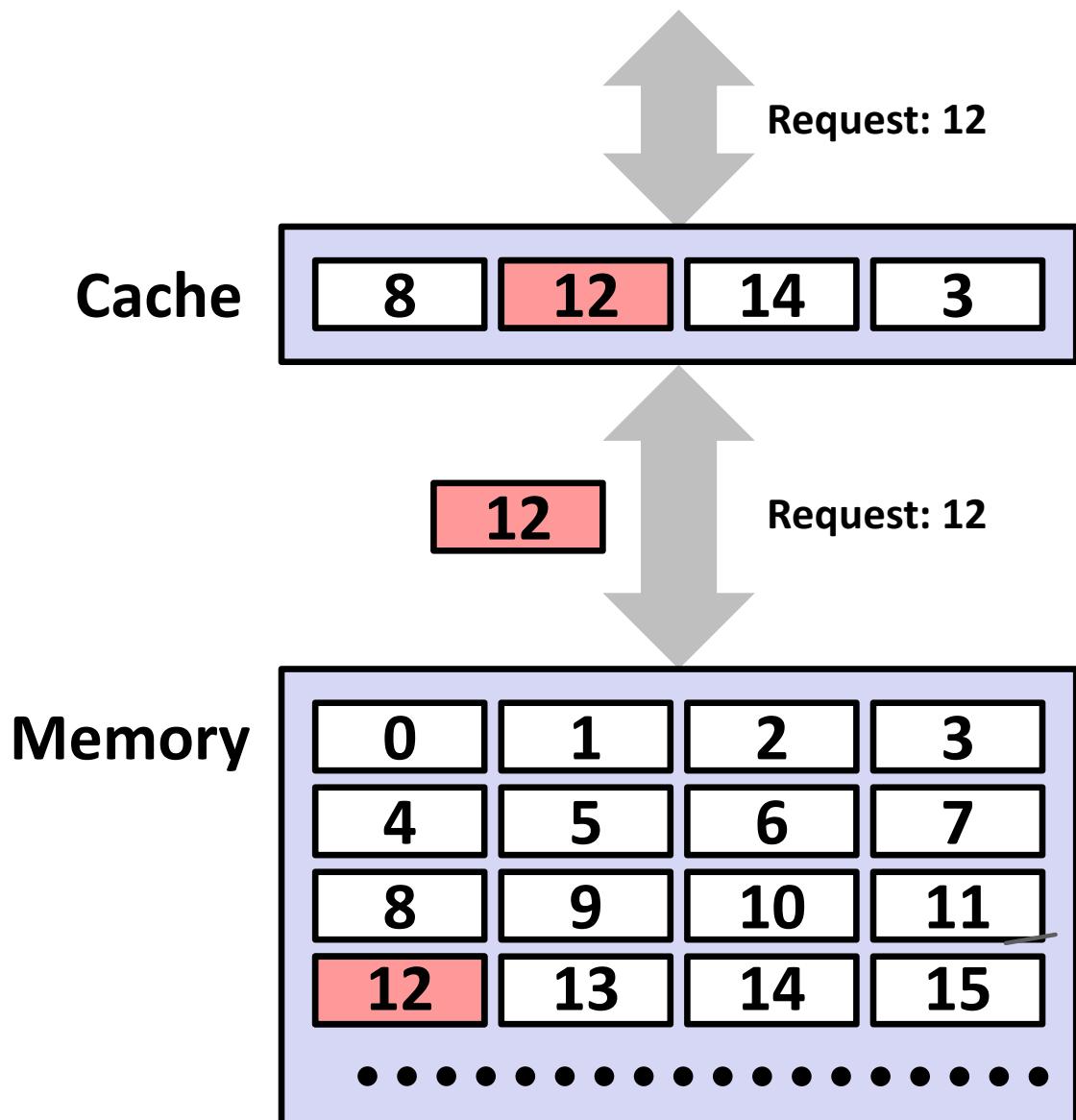


Larger, slower, cheaper memory viewed as partitioned into “blocks”

General Cache Concepts: Hit



General Cache Concepts: Miss



Data in block b is needed

*Block b is not in cache:
Miss!*

*Block b is fetched from
memory*

Block b is stored in cache

- **Placement policy:**
determines where b goes
- **Replacement policy:**
determines which block gets evicted (victim)

Examples of Caching in the Hierarchy

Cache Type	What is Cached?	Where is it Cached?	Latency (cycles)	Managed By
Registers	4-8 bytes words	CPU core	0	Compiler
TLB	Address translations	On-Chip TLB	0	Hardware
L1 cache	64-bytes block	On-Chip L1	1	Hardware
L2 cache	64-bytes block	On/Off-Chip L2	10	Hardware
Virtual Memory	4-KB page	Main memory	100	Hardware + OS
Buffer cache	Parts of files	Main memory	100	OS
Disk cache	Disk sectors	Disk controller	100,000	Disk firmware
Network buffer cache	Parts of files	Local disk	10,000,000	AFS/NFS client
Browser cache	Web pages	Local disk	10,000,000	Web browser
Web cache	Web pages	Remote server disks	1,000,000,000	Web proxy server

Summary

Principal of locality → Temporal locality

↳ if I look at A, in near time I
may look A again

Spatial locality

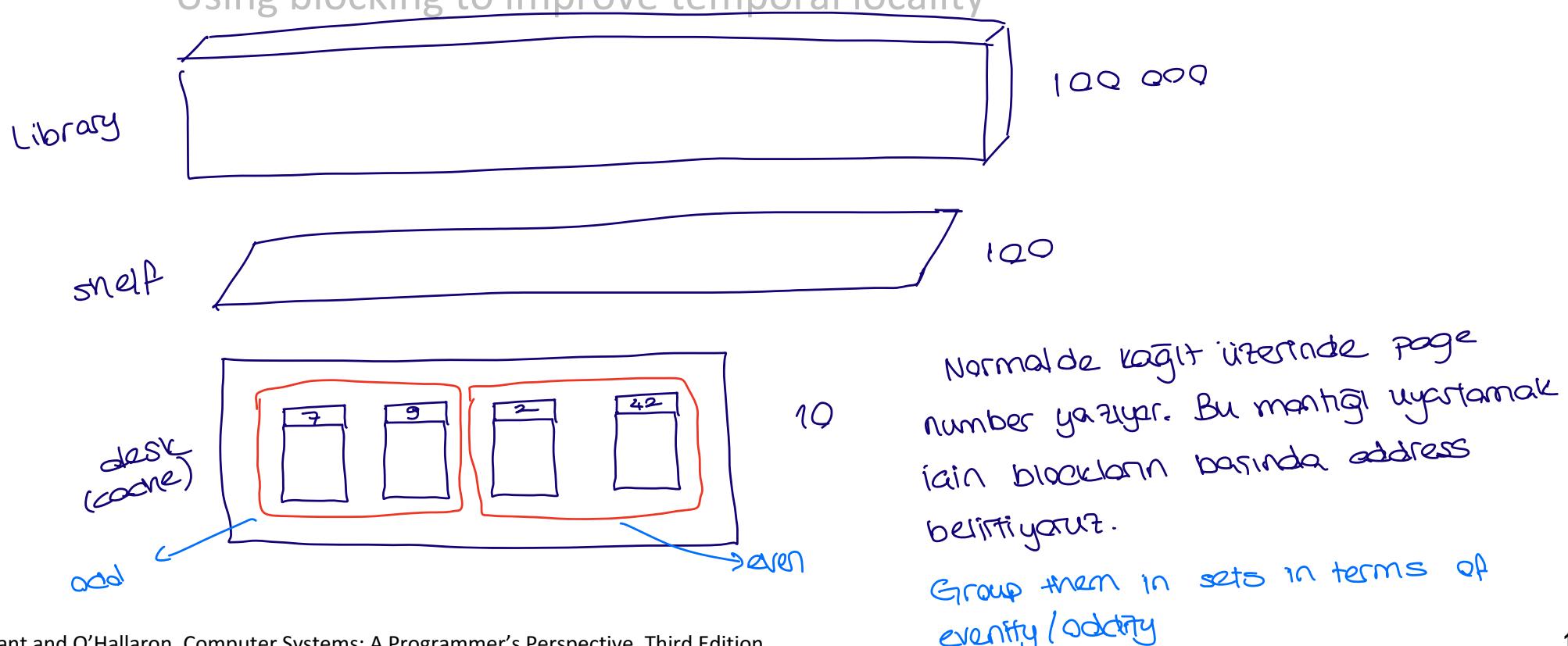
- The speed gap between CPU, memory and mass storage continues to widen.
- Well-written programs exhibit a property called locality.
- Memory hierarchies based on caching close the gap by exploiting locality.

Cache

- Cache memory organization and operation

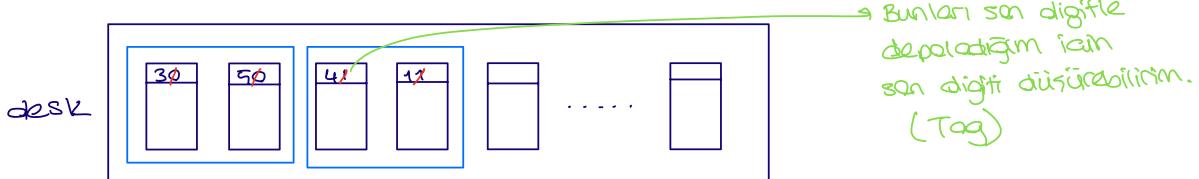
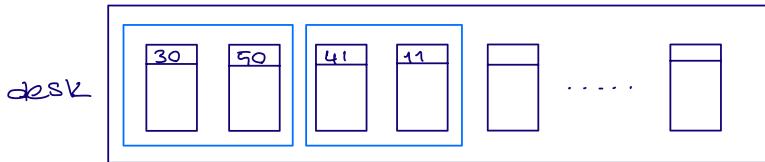
- Performance impact of caches

- The memory mountain
- Rearranging loops to improve spatial locality
- Using blocking to improve temporal locality

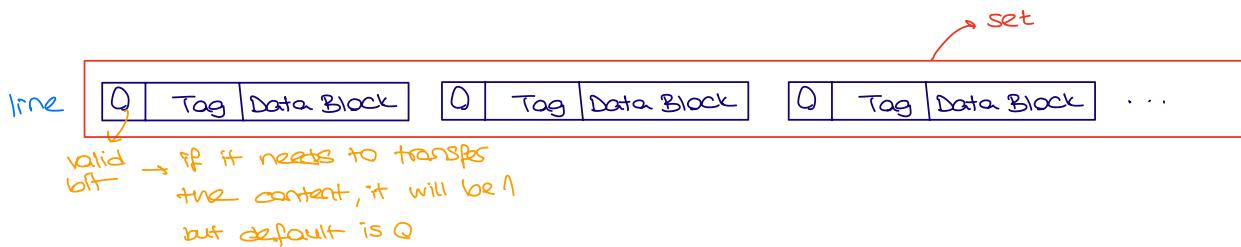


QZ

Group them in sets by their last digits



each set contains multiple lines



B: 2^b bytes per cache block

E: 2^e lines per set

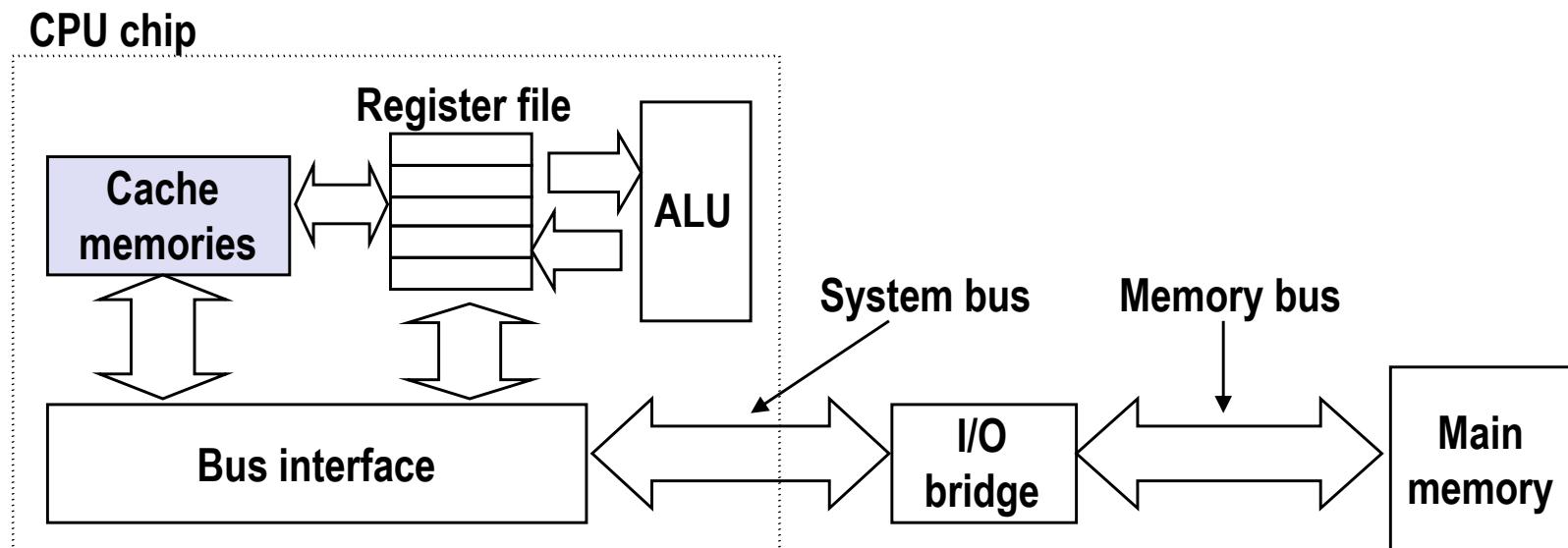
S: 2^s sets

$$B = 6 \text{ bytes} = 2^6$$

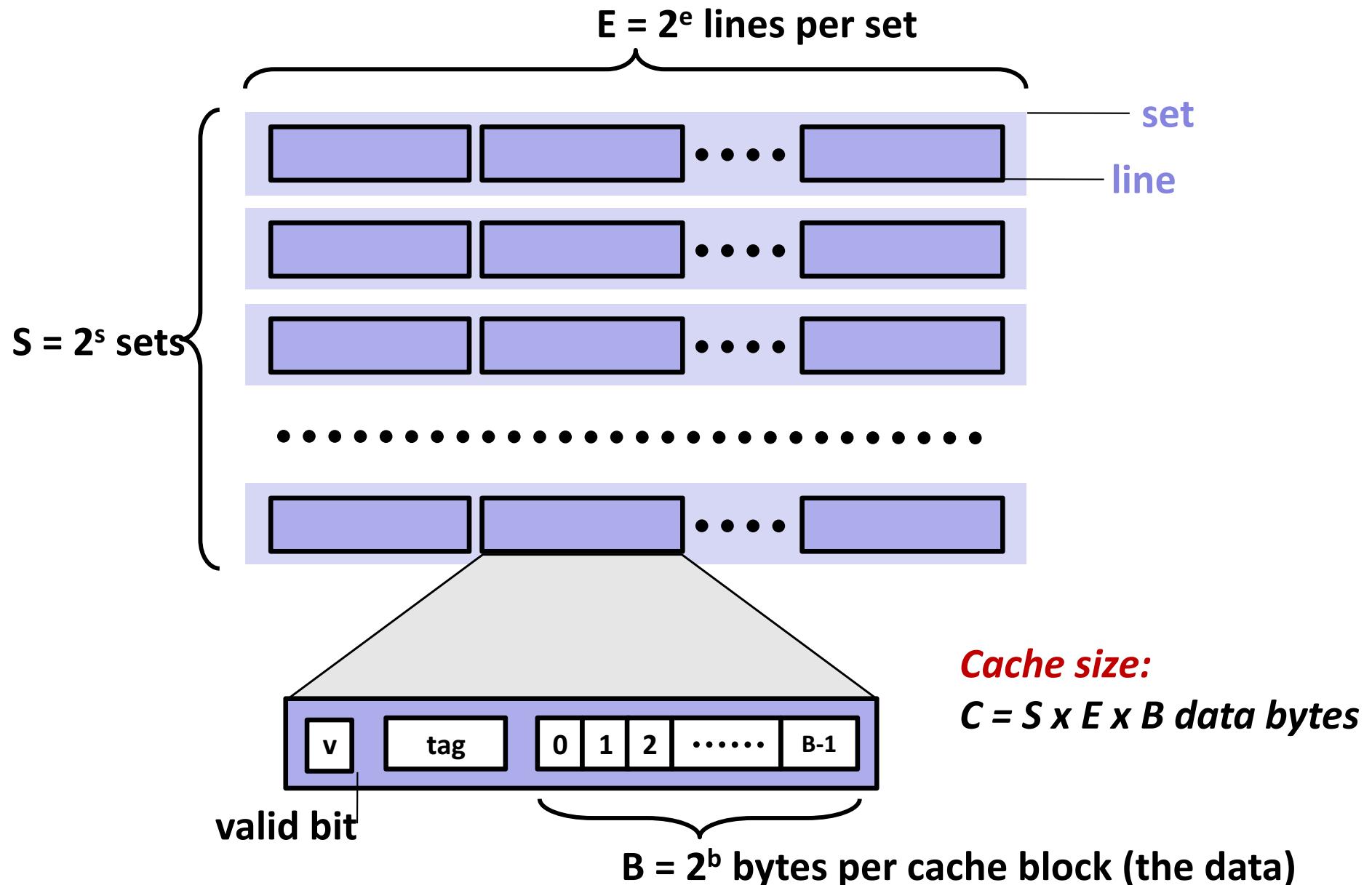


Cache Memories

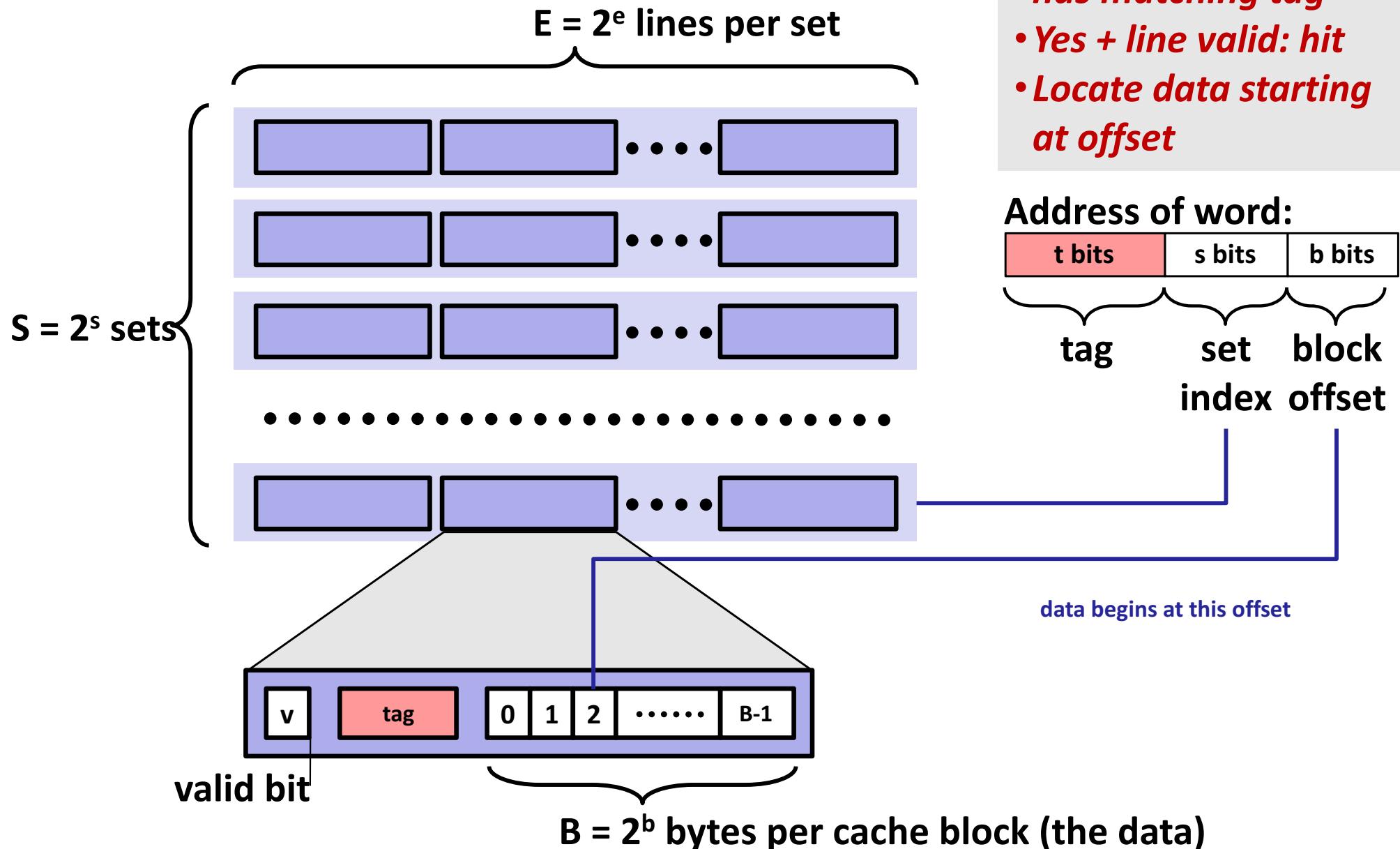
- Cache memories are small, fast SRAM-based memories managed automatically in hardware.
 - Hold frequently accessed blocks of main memory
- CPU looks first for data in caches (e.g., L1, L2, and L3), then in main memory.
- Typical system structure:



General Cache Organization (S, E, B)



Cache Read



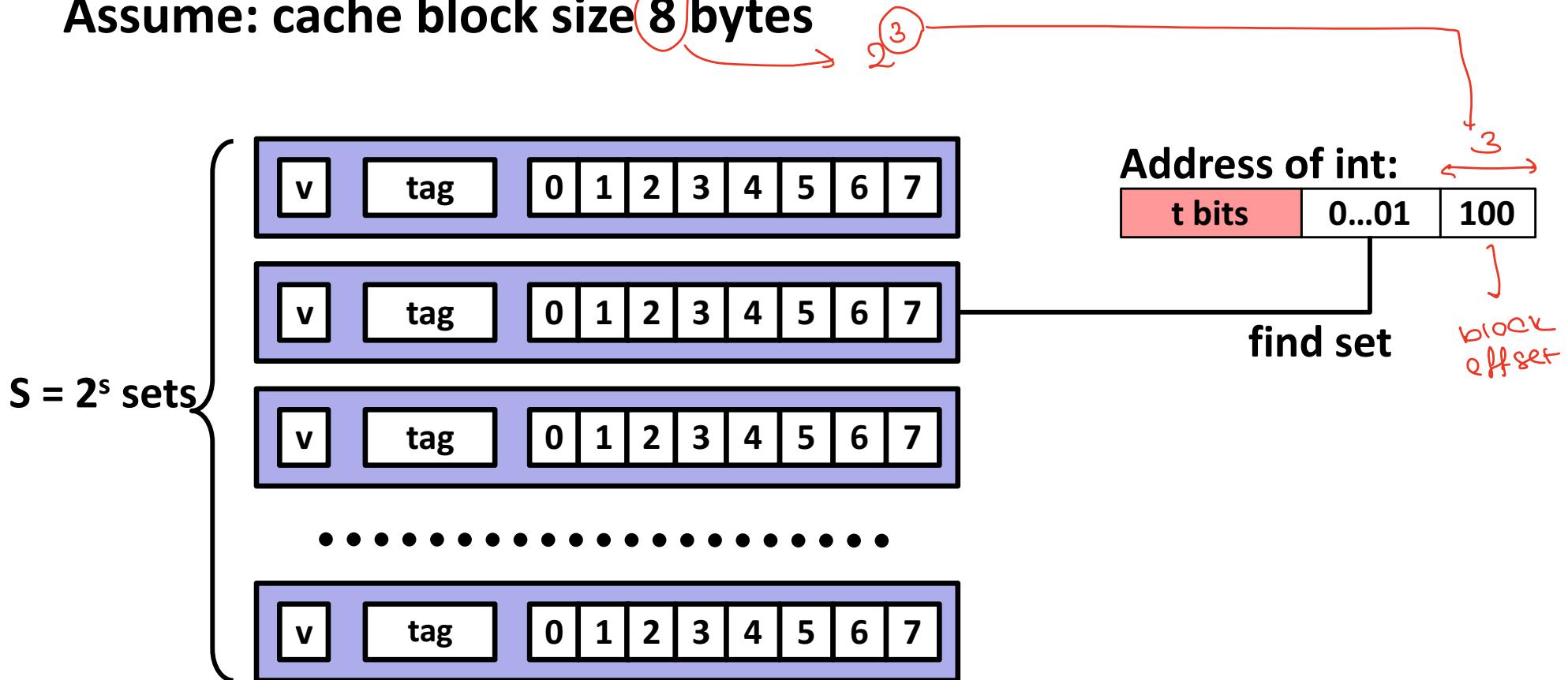
- Locate set
- Check if any line in set has matching tag
- Yes + line valid: hit
- Locate data starting at offset

Example: Direct Mapped Cache ($E = 1$)

Direct mapped: One line per set

Assume: cache block size $\textcircled{8}$ bytes

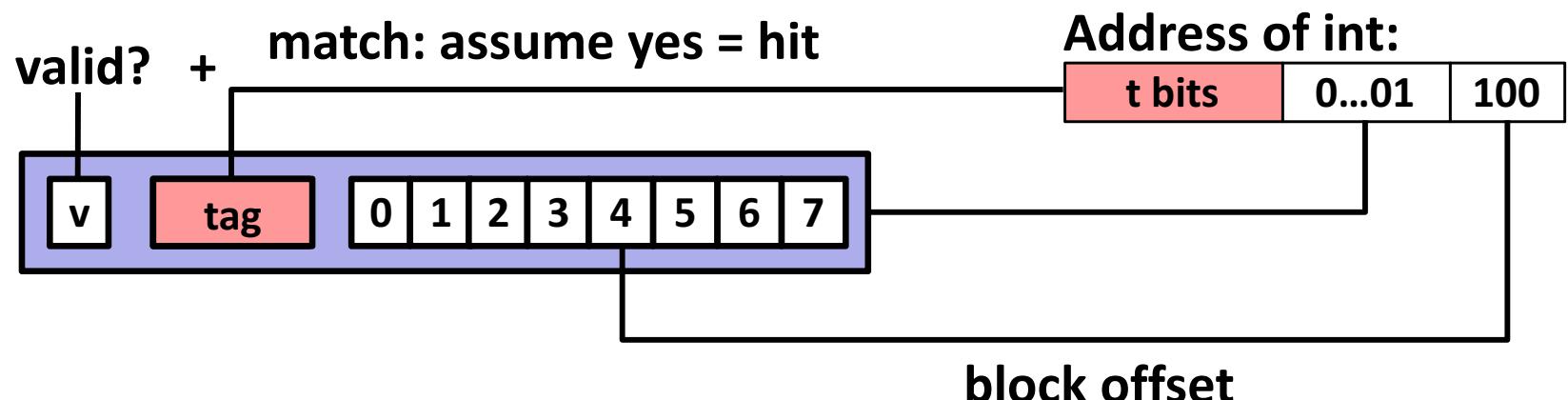
1 lines per set



Example: Direct Mapped Cache ($E = 1$)

Direct mapped: One line per set

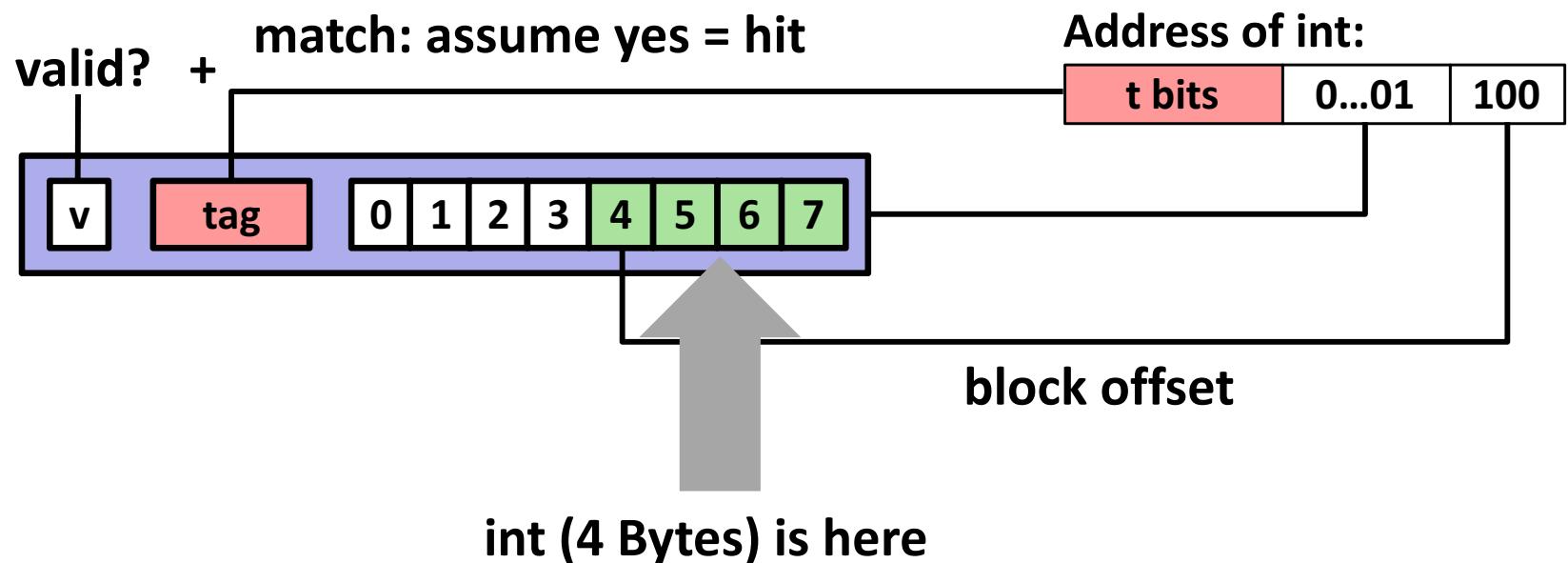
Assume: cache block size 8 bytes



Example: Direct Mapped Cache ($E = 1$)

Direct mapped: One line per set

Assume: cache block size 8 bytes



No match: old line is evicted and replaced

Direct-Mapped Cache Simulation

$t=1 \quad s=2 \quad b=1$

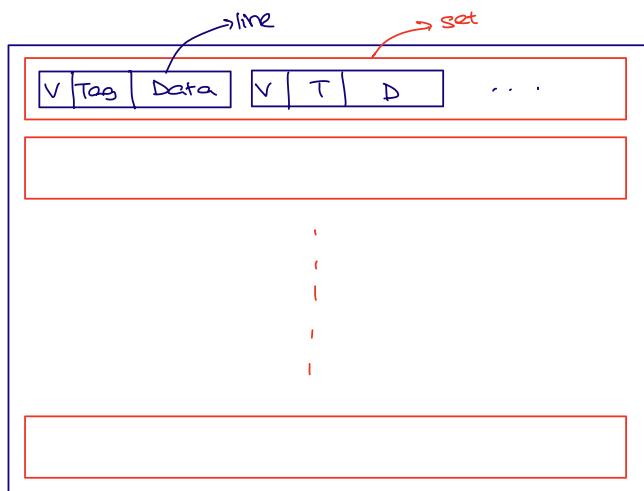
x	xx	x
---	----	---

M=16 byte addresses, B=2 bytes/block,
S=4 sets, E=1 Blocks/set

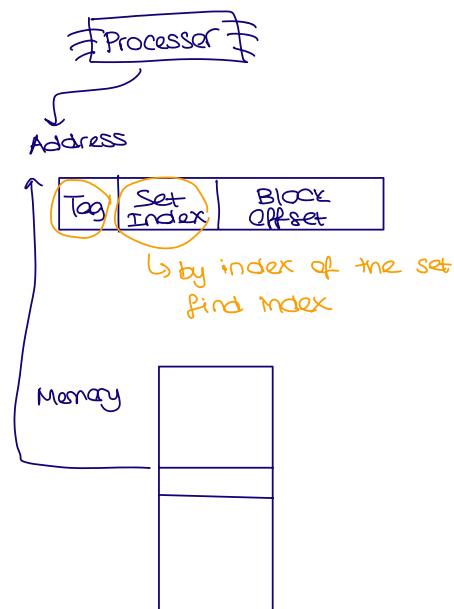
Address trace (reads, one byte per read):

0	[<u>0000</u> ₂],	miss
1	[<u>0001</u> ₂],	hit
7	[<u>0111</u> ₂],	miss
8	[<u>1000</u> ₂],	miss
0	[<u>0000</u> ₂]	miss

	v	Tag	Block
Set 0	1	0	M[0-1]
Set 1			
Set 2			
Set 3	1	0	M[6-7]



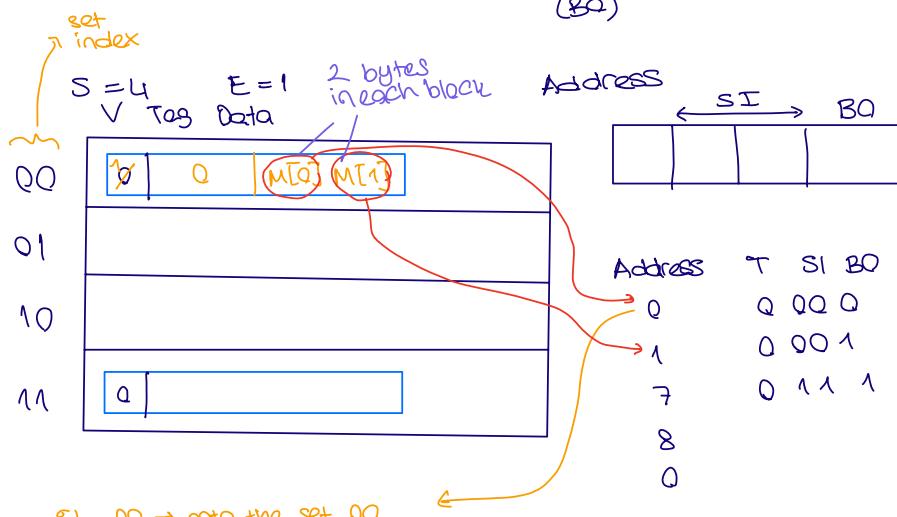
$$\text{coarse size} = C = S \times E \times B$$



|||||| ||||| |||||

$$\text{Memory} = 16 \text{ bytes} = 2^4 \Rightarrow |\text{Address}| = 4 \text{ bits}$$

$$B = 2 \text{ bytes} = 2^1 \Rightarrow \text{Byte offset} = 1 \text{ bit} \\ (\text{Eq})$$



$g1 = \{Q\} \rightarrow$ goto the set $\{Q\}$

check the valid bit \Rightarrow it is 0 \Rightarrow cache miss

Go to the memory

first consecutive 2 bit before or after our bit Q

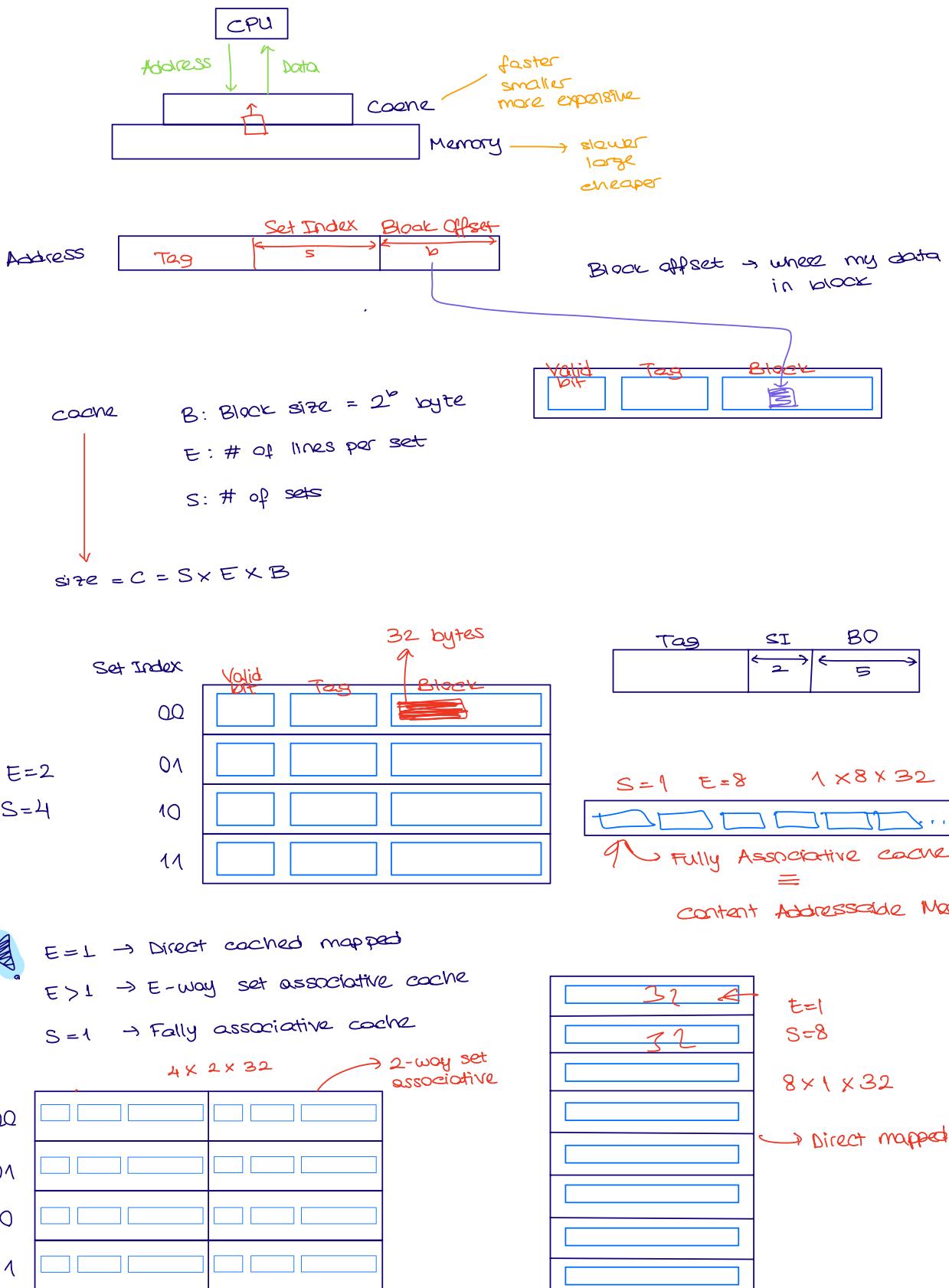
(α ve β 'i global icsin) (Bizim bitimiz 1 oluyor da Q ile 15 olurdu)

Put one into the other post as $M[0] M[1]$

Set valid last 1

Set to as Q

0000	M[0]
0001	M[1]
0010	M[2]
0011	M[3]
0100	.
0101	.
0110	.
0111	.
1000	.
1001	M[5]



A Higher Level Example

```
int sum_array_rows(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```

```
int sum_array_cols(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (j = 0; j < 16; j++)
        for (i = 0; i < 16; i++)
            sum += a[i][j];
    return sum;
}
```



Ignore the variables sum, i, j

assume: cold (empty) cache,
a[0][0] goes here



32 B = 4 doubles

E tag comparisons per access

E-way Set Associative Cache (Here: E = 2)

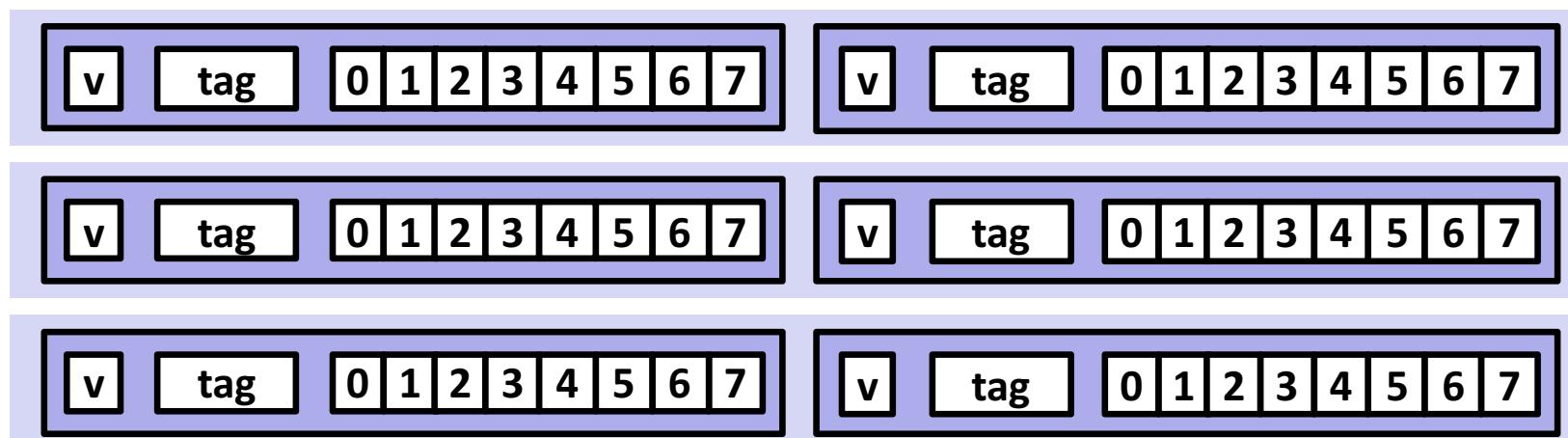
E = 2: Two lines per set

Assume: cache block size 8 bytes

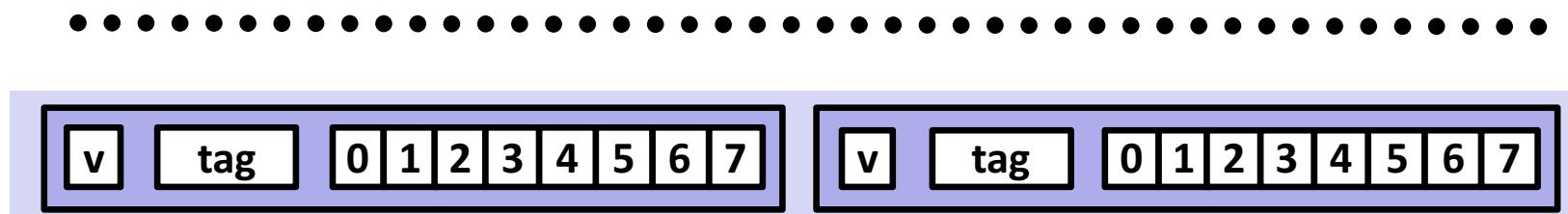


Address

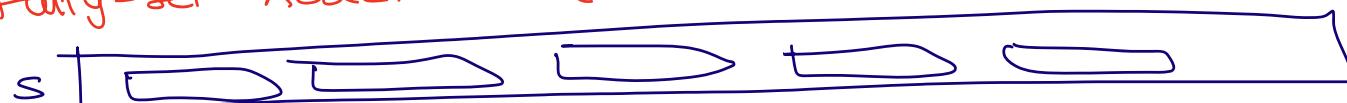
Address of short int:



find set



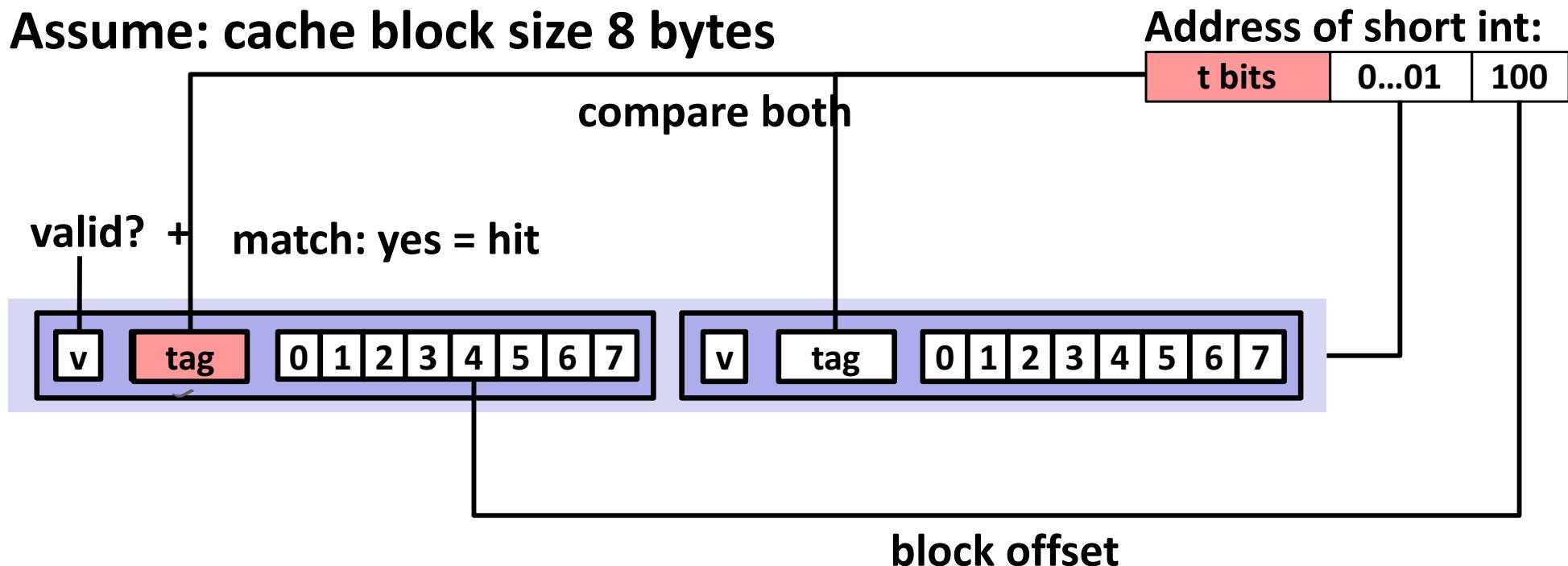
Fully-set Associative ($S=1$) \rightarrow content addressable memory



E-way Set Associative Cache (Here: E = 2)

E = 2: Two lines per set

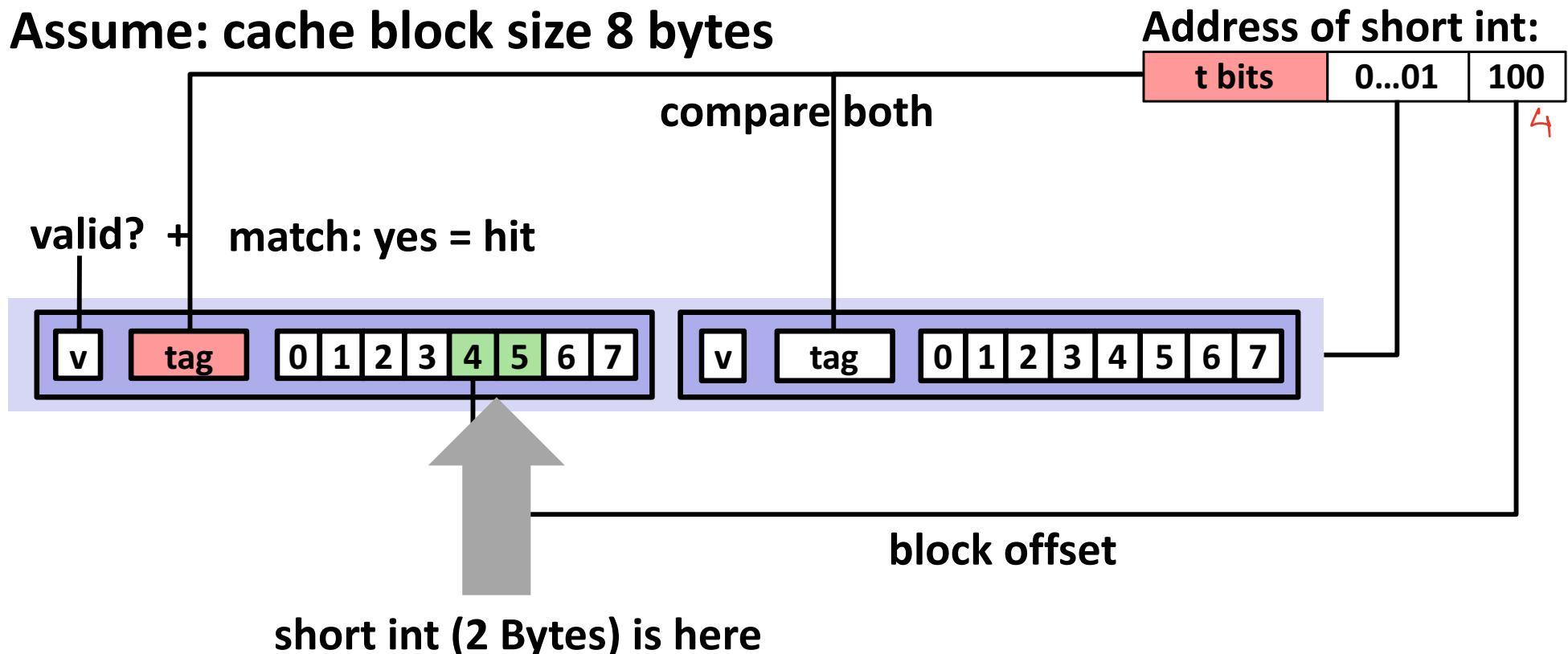
Assume: cache block size 8 bytes



E-way Set Associative Cache (Here: E = 2)

$E = 2$: Two lines per set

Assume: cache block size 8 bytes



No match:

- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

2-Way Set Associative Cache Simulation

$t=2$ $s=1$ $b=1$

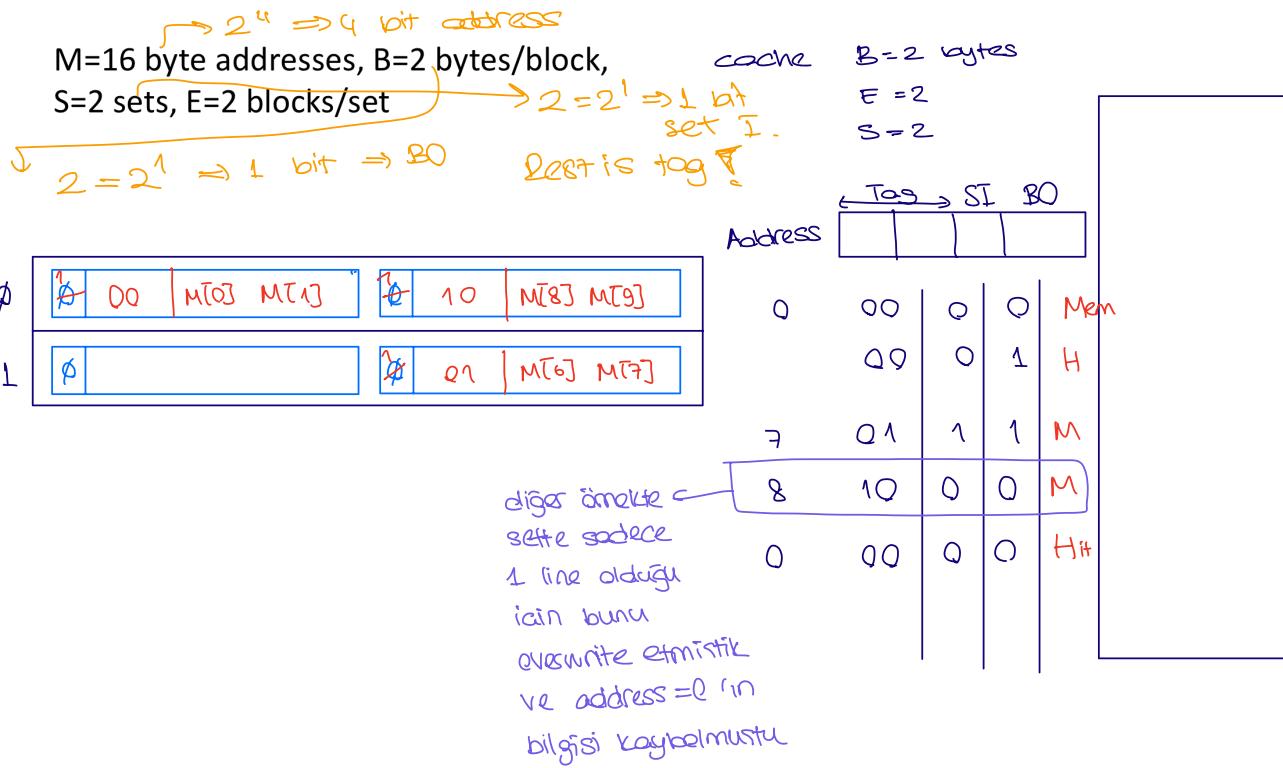
xx	x	x
----	---	---

$M=16$ byte addresses, $B=2$ bytes/block,
 $S=2$ sets, $E=2$ blocks/set

Address trace (reads, one byte per read):

0	$[0000_2]$,	miss
1	$[0001_2]$,	hit
7	$[0111_2]$,	miss
8	$[1000_2]$,	miss
0	$[0000_2]$	hit

	v	Tag	Block
Set 0	1	00	$M[0-1]$
	1	10	$M[8-9]$
Set 1	1	01	$M[6-7]$
	0		



A Higher Level Example

```
int sum_array_rows(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```

(Handwritten annotations: A blue box highlights the inner loop iteration a[0][0] through a[0][3]. Above the box, 'M' is written above each row. To the right of the box, 'H' is written above each column, and 'Q03' is written below the box.)

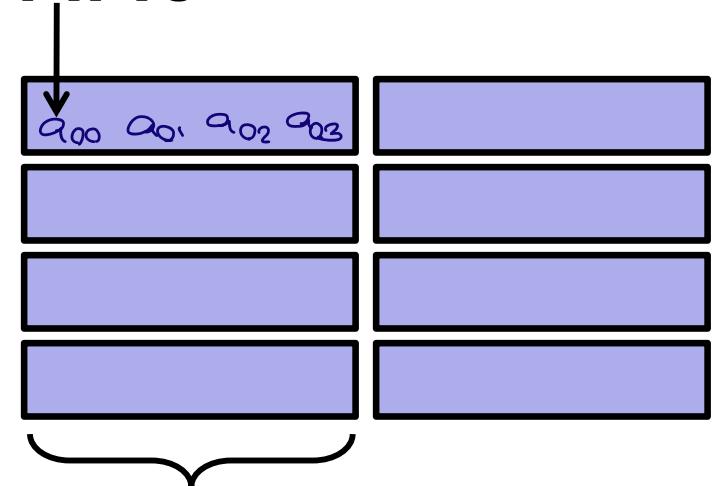
```
int sum_array_rows(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (j = 0; j < 16; j++)
        for (i = 0; i < 16; i++)
            sum += a[i][j];
    return sum;
}
```

(Handwritten annotations: A blue box highlights the inner loop iteration a[0][0] through a[3][0]. Above the box, 'M' is written above each row. To the right of the box, 'M' is written above each column, and 'Q30' is written below the box. Below the box, 'Q40 ...' is written.)

Ignore the variables sum, i, j

**assume: cold (empty) cache,
a[0][0] goes here**



32 B = 4 doubles

What about writes?

■ Multiple copies of data exist:

- L1, L2, Main Memory, Disk

■ What to do on a write-hit?

- Write-through (write immediately to memory)
- Write-back (defer write to memory until replacement of line)

along with the valid bit, hold a dirty bit

■ Need a dirty bit (line different from memory or not)

■ What to do on a write-miss?

- Write-allocate (load into cache, update line in cache)
 - if the contents size < block size
- Good if more writes to the location follow
- No-write-allocate (writes immediately to memory)

■ Typical

- Write-through + No-write-allocate

- Write-back + Write-allocate

Typically

The line you want to write is in the cache

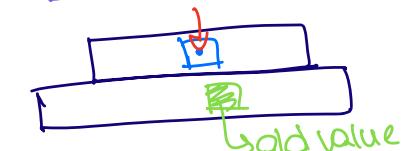
write both in cache and to mem.



→ write it to the lower level



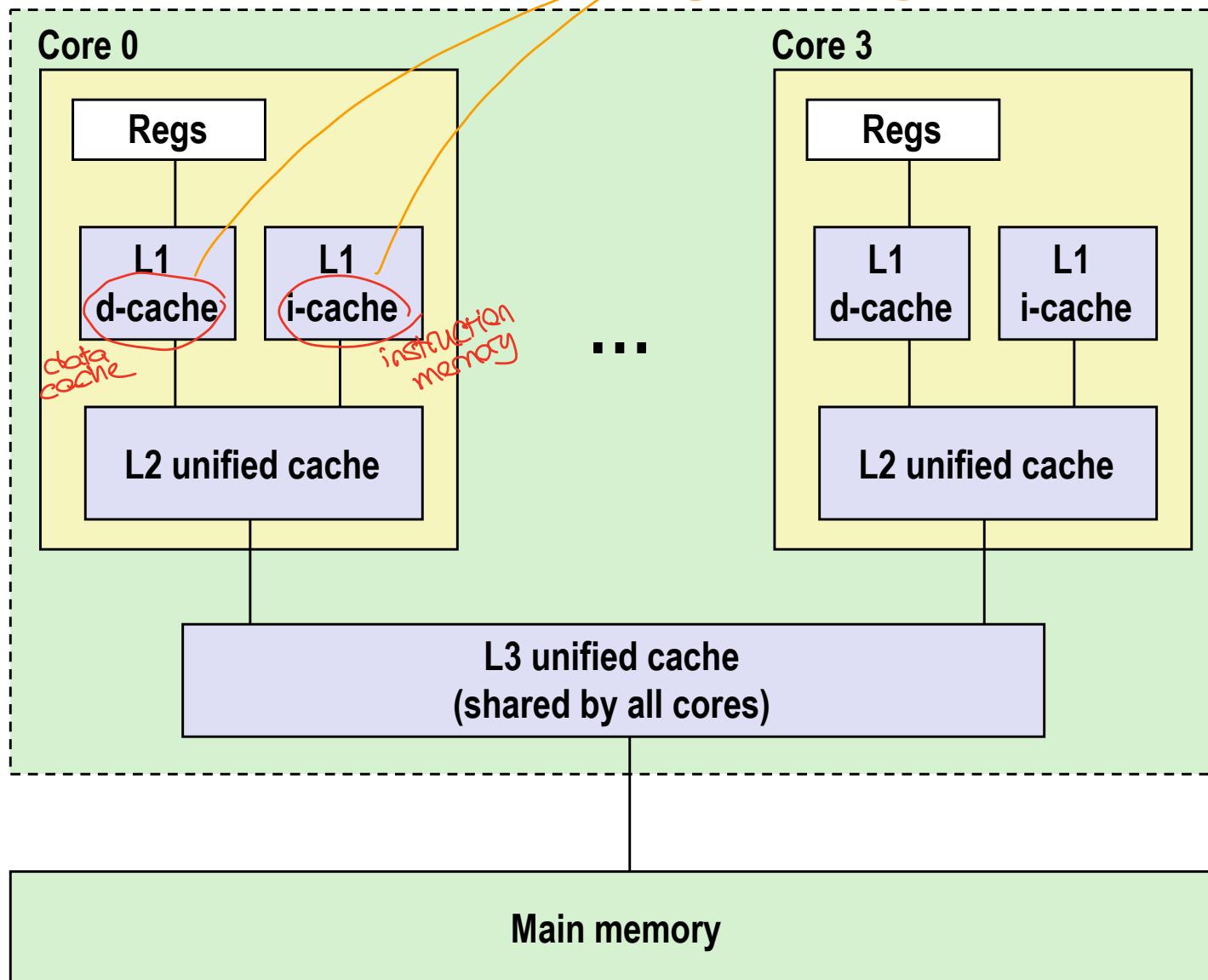
D=0 → mem has old value
D=1 → write-back



old value

Intel Core i7 Cache Hierarchy

Processor package



L1 i-cache and d-cache:

32 KB, 8-way,
Access: 4 cycles

L2 unified cache:

256 KB, 8-way,
Access: 11 cycles

L3 unified cache:

8 MB, 16-way,
Access: 30-40
cycles

Block size: 64 bytes
for all caches.

L1 i-cache and d-cache

32 KB, 8-way

$$C = 32 \text{ KB} = 2^{15} \text{ bytes}$$

$$E = 8 = 2^3$$

$$S = 2^6$$

$$B = 64 = 2^6$$

$$C = S \times E \times B = 2^6 \times 2^3 \times 2^6 = 2^{15}$$



L1 i-cache and d-cache:
32 KB, $E=8$, 8-way,
Access: 4 cycles

Block size: 64 bytes
for all caches.

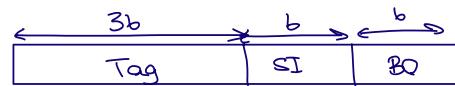
$$B = 64 \text{ bytes} \Rightarrow 2^6 \Rightarrow 6 \text{ bit BO}$$

$$E = 8 \Rightarrow 2^3$$

$$C = 2^{15} = S \times 2^3 \times 2^6 \Rightarrow S = 2^6$$

SI \Rightarrow 6 bits

Mem address \Rightarrow 48 bits

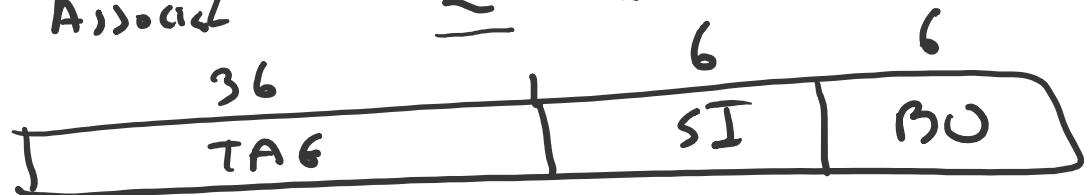


i7 L1 d-cache

Block: 64 bytes $\rightarrow 2^6$

Cache size = 32 KB $\rightarrow 2^{15} = C$
8-Way Set Assoc

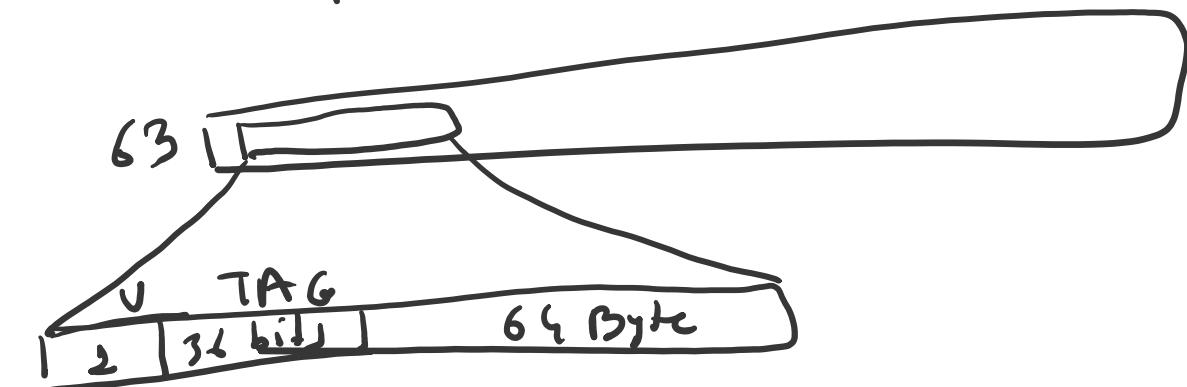
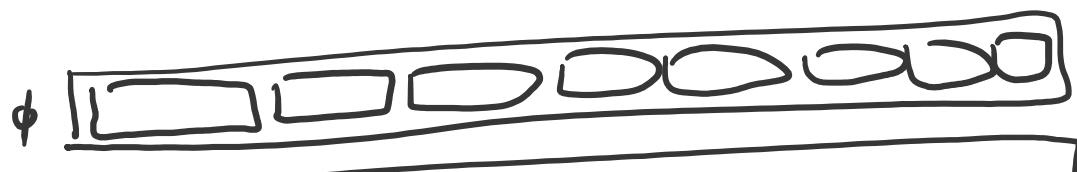
Address = 48 bits



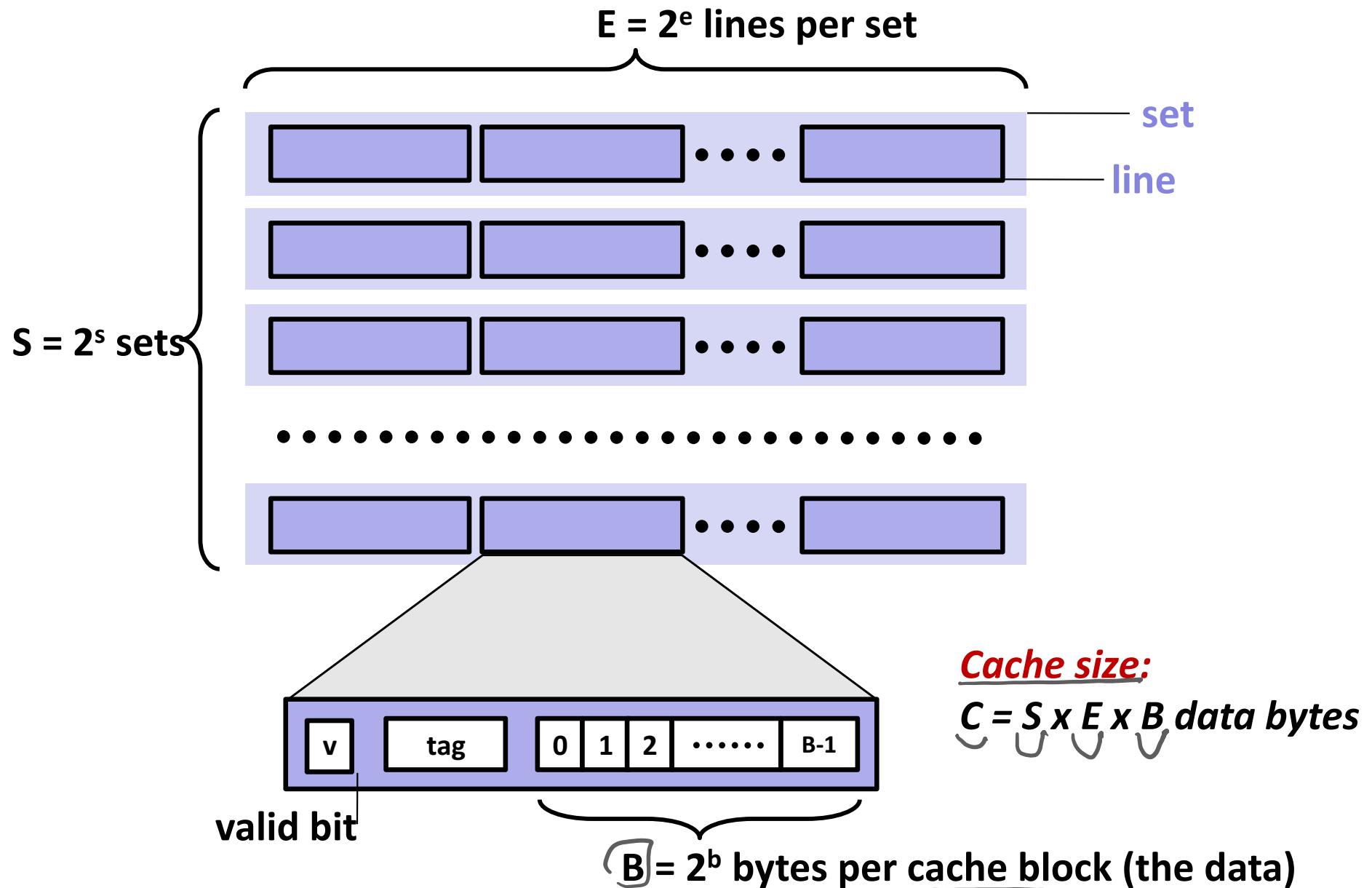
$$C = S \cdot E \cdot B$$

$$2^{15} = S \cdot 2^3 \cdot 2^6$$

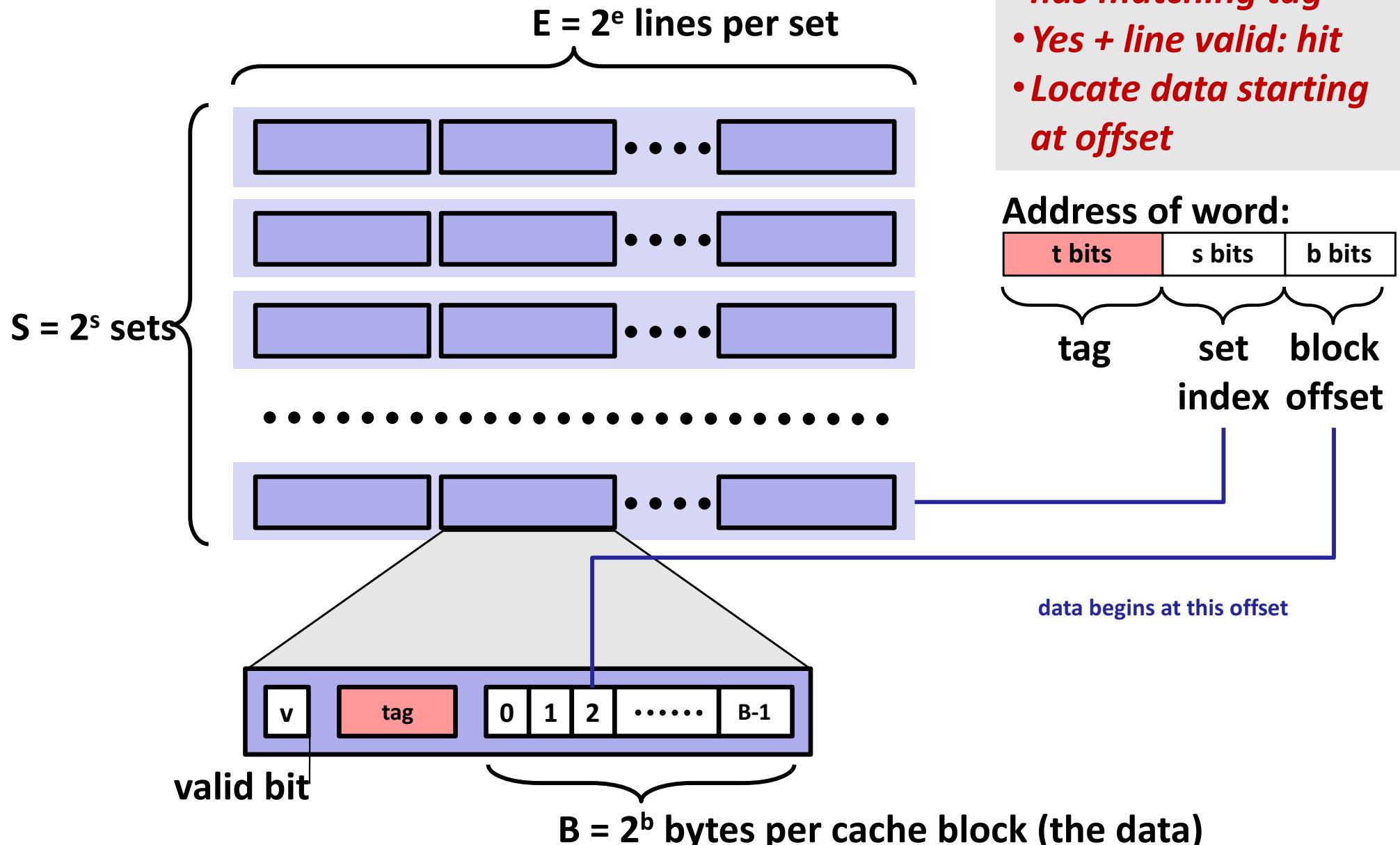
$$\Rightarrow S \Rightarrow 2^6$$



General Cache Organization (S , E , B)

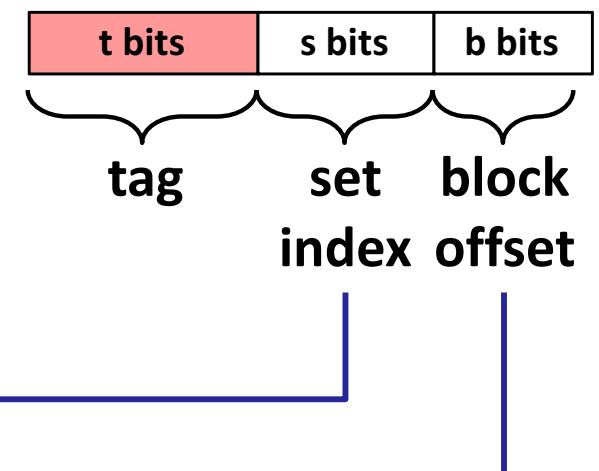


Cache Read

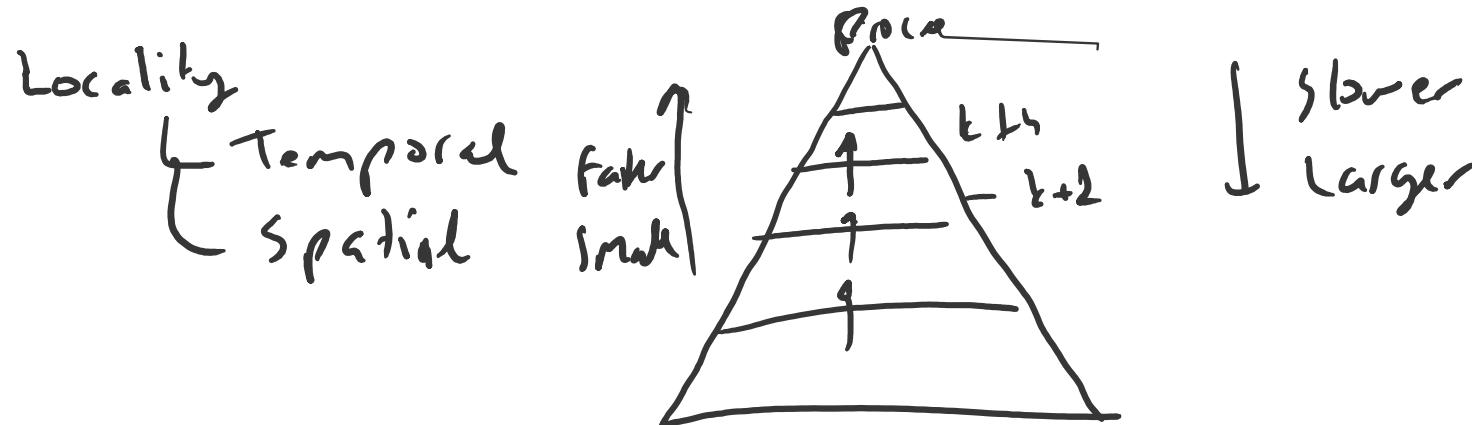


- **Locate set**
- **Check if any line in set has matching tag**
- **Yes + line valid: hit**
- **Locate data starting at offset**

Address of word:



data begins at this offset



$$C = S \times E \times B$$

↳ # bytes per block
 ↳ # of lines per set
 ↳ # of sets

Direct-mapped

$E=1$



One line per set!

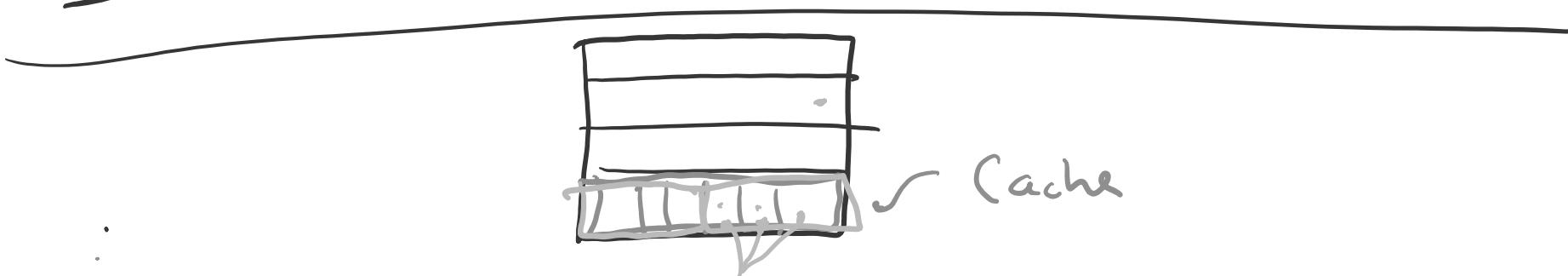
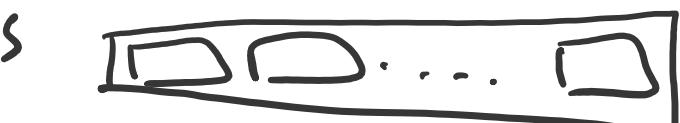
 K -way-set associative

$E = K$ K lines per set



Fully-set Associative

$S=1$



General Caching Concepts:

Types of Cache Misses

■ Cold (compulsory) miss

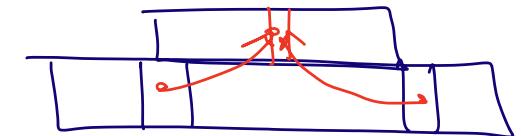
- Cold misses occur because the cache is empty.

■ Conflict miss *→ more likely to happen in direct-mapped access*

- Most caches limit blocks at level $k+1$ to a small subset (sometimes a singleton) of the block positions at level k .
 - E.g. Block i at level $k+1$ must be placed in block $(i \bmod 4)$ at level k .
- Conflict misses occur when the level k cache is large enough, but multiple data objects all map to the same level k block.
 - E.g. Referencing blocks 0, 8, 0, 8, 0, 8, ... would miss every time.

■ Capacity miss

- Occurs when the set of active cache blocks (**working set**) is larger than the cache.



Cache Performance Metrics

■ Miss Rate

- Fraction of memory references not found in cache (misses / accesses)
= 1 – hit rate
- Typical numbers (in percentages):
 - 3-10% for L1
 - can be quite small (e.g., < 1%) for L2, depending on size, etc.

■ Hit Time

- Time to deliver a line in the cache to the processor
 - includes time to determine whether the line is in the cache
- Typical numbers:
 - 1-2 clock cycle for L1
 - 5-20 clock cycles for L2

■ Miss Penalty

- Additional time required because of a miss
 - typically 50-200 cycles for main memory (Trend: increasing!)

Lets think about those numbers

■ Huge difference between a hit and a miss

- Could be 100x, if just L1 and main memory

■ Would you believe 99% hits is twice as good as 97%?

- Consider:

cache hit time of 1 cycle

miss penalty of 100 cycles

- Average access time:

$$97\% \text{ hits: } 1 \text{ cycle} + 0.03 * 100 \text{ cycles} = \mathbf{4 \text{ cycles}}$$

$$99\% \text{ hits: } 1 \text{ cycle} + 0.01 * 100 \text{ cycles} = \mathbf{2 \text{ cycles}}$$

■ This is why “miss rate” is used instead of “hit rate”

Writing Cache Friendly Code

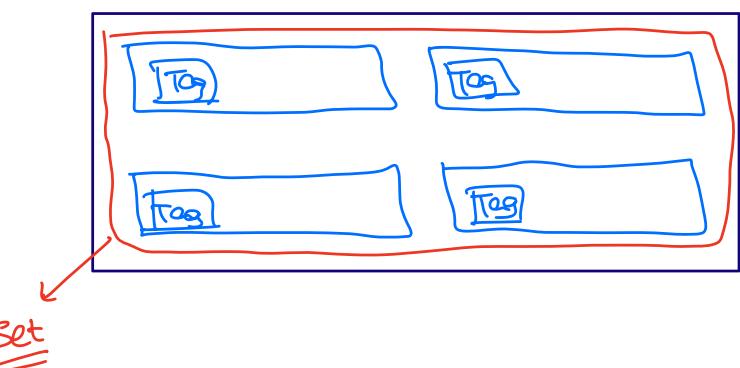
- **Make the common case go fast**
 - Focus on the inner loops of the core functions
- **Minimize the misses in the inner loops**
 - Repeated references to variables are good (**temporal locality**)
 - Stride-1 reference patterns are good (**spatial locality**)

Key idea: Our qualitative notion of locality is quantified through our understanding of cache memories.

Today

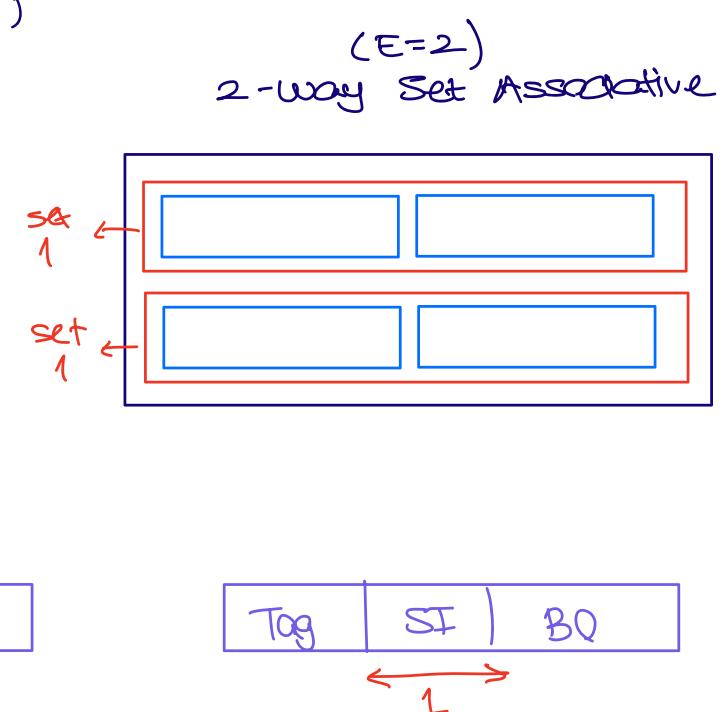
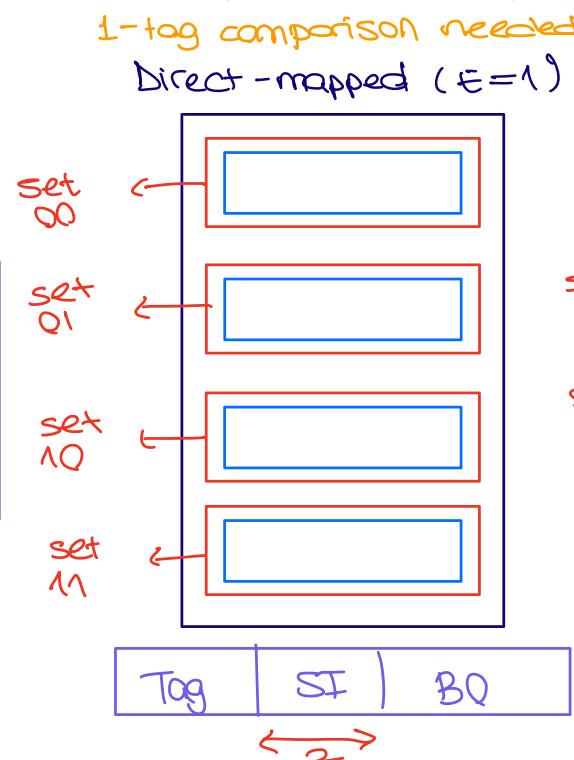
- Cache organization and operation
- Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

Given $C = S \times E \times B$
 Max number of Tag comparisons needed
 Fully Associative ($S=1$)



Address

Tag	BO
-----	----



The Memory Mountain

- **Read throughput (read bandwidth)**
 - Number of bytes read from memory per second (MB/s)
- **Memory mountain:** Measured read throughput as a function of spatial and temporal locality.
 - Compact way to characterize memory system performance.

Memory Mountain Test Function

```

/* The test function */
void test(int elems, int stride) {
    int i, result = 0;
    volatile int sink;

    for (i = 0; i < elems; i += stride)
        result += data[i];
    sink = result; /* So compiler doesn't optimize away the loop */
}

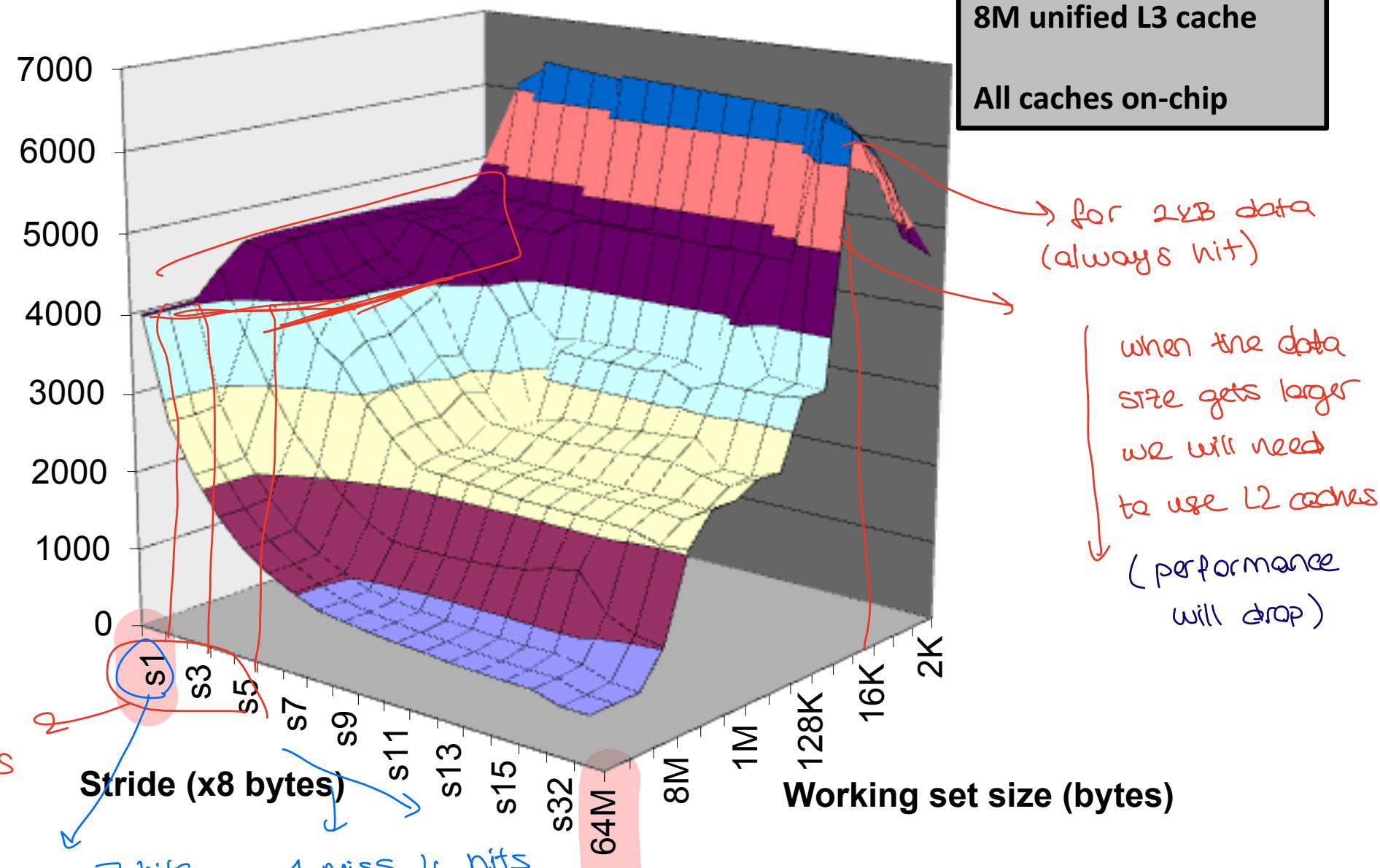
/* Run test(elems, stride) and return read throughput (MB/s) */
double run(int size, int stride, double Mhz)
{
    double cycles;
    int elems = size / sizeof(int);

    test(elems, stride); warm the cache boyelice bazı elementler  
cache'de olacak
    cycles = fcyc2(test, elems, stride, 0); /* call test(elems,stride) */
    return (size / stride) / (cycles / Mhz); /* convert cycles to MB/s */
}

```

The Memory Mountain

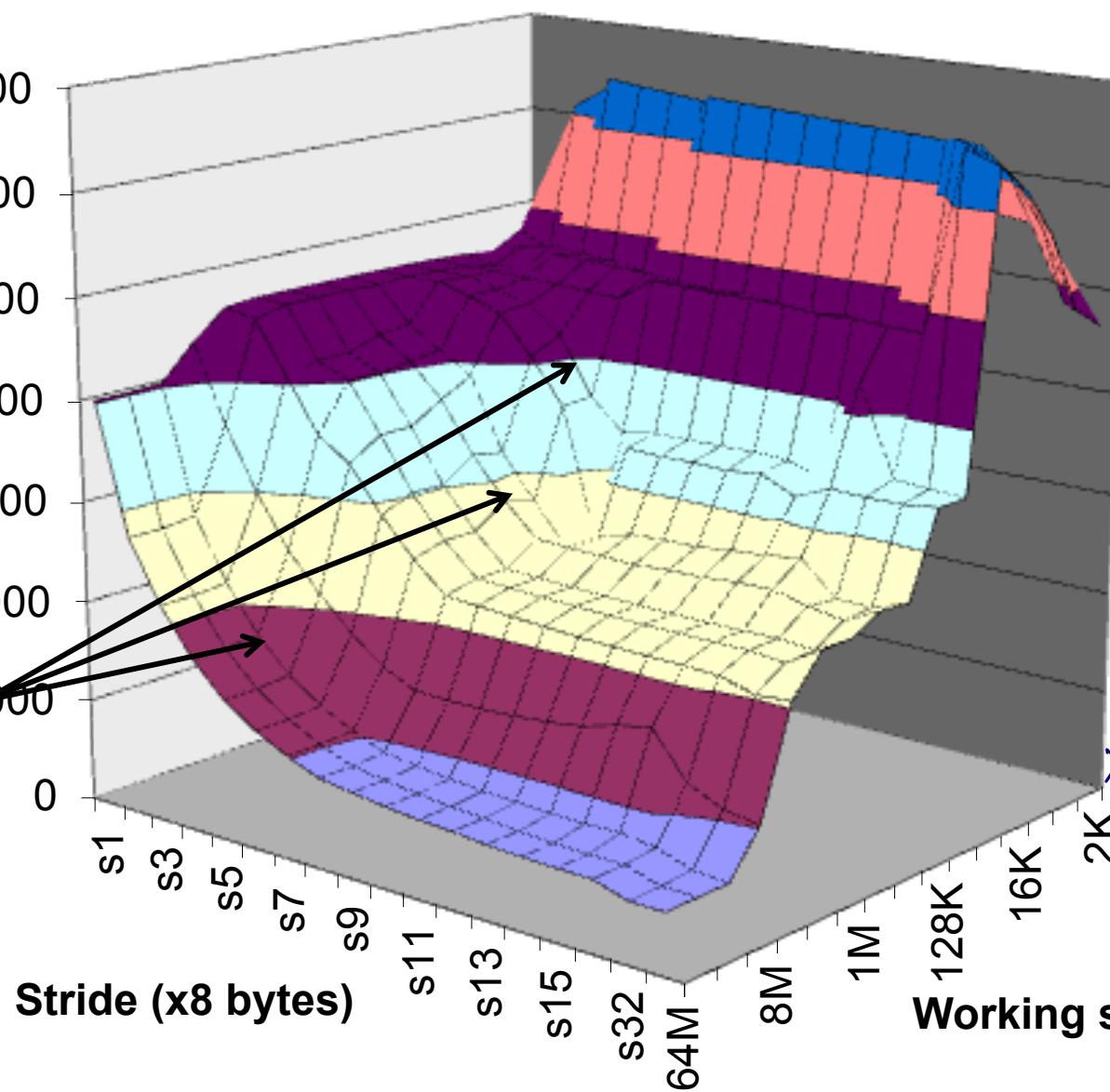
Read throughput (MB/s)



The Memory Mountain

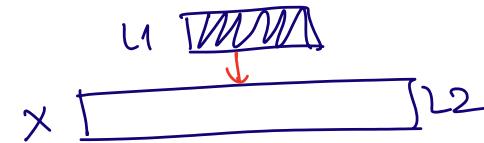
Read throughput (MB/s)

Slopes of spatial locality

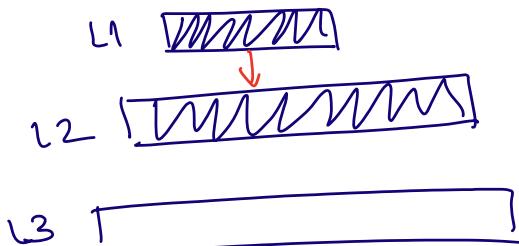


Intel Core i7
 32 KB L1 i-cache
 32 KB L1 d-cache
 256 KB unified L2 cache
 8M unified L3 cache
All caches on-chip

working set (ws) < 32 KB



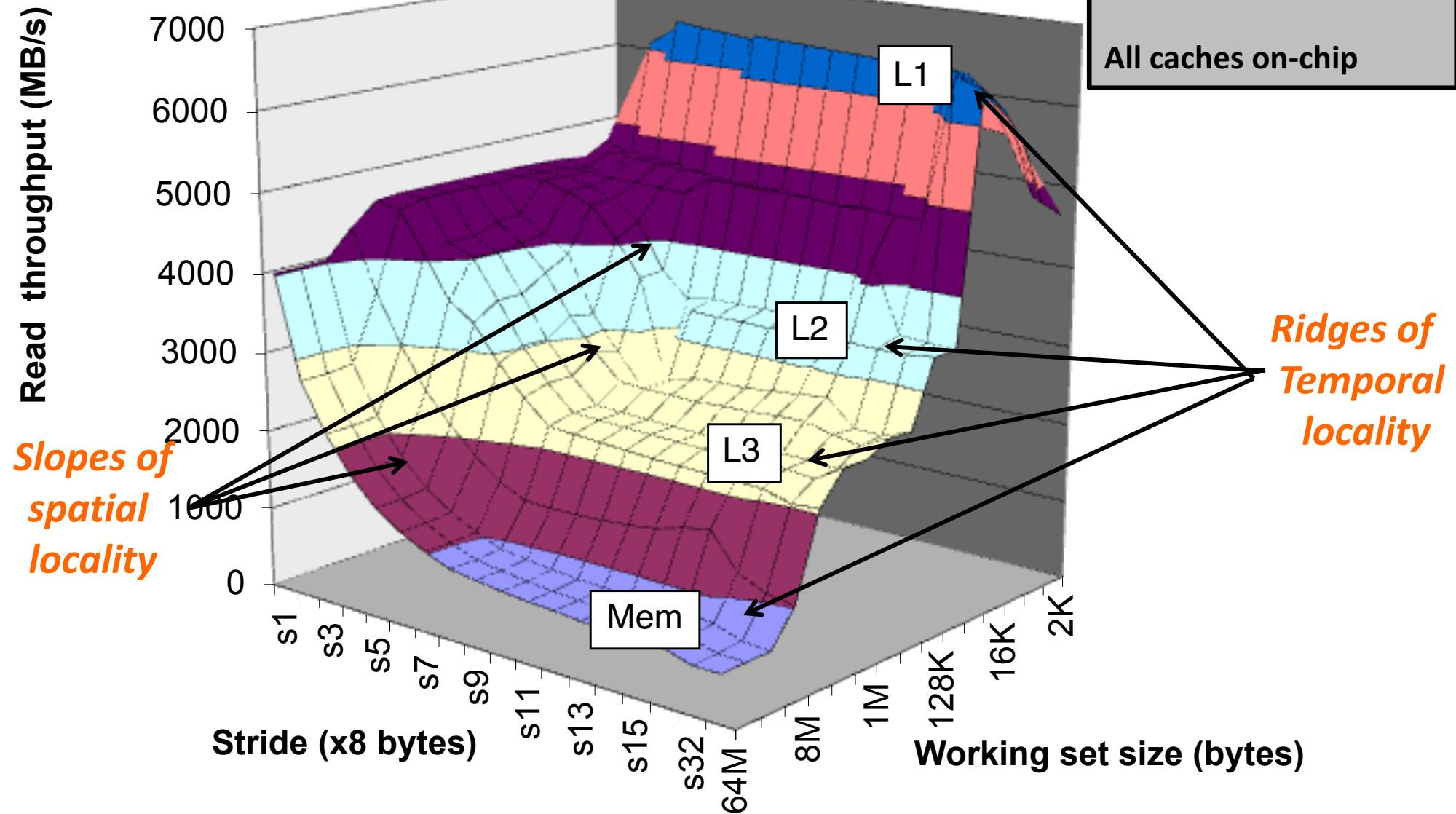
32 KB < ws < 256 KB



256 KB < ws



The Memory Mountain



Today

- Cache organization and operation
- Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

Miss Rate Analysis for Matrix Multiply

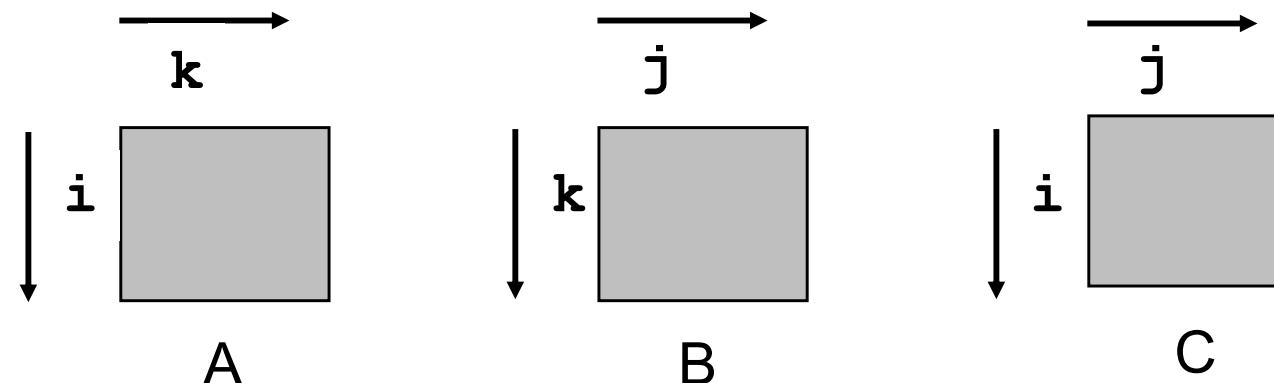
■ Assume:

- Block size = $32B$ (big enough for four 64-bit words)
- Matrix dimension (N) is very large
 - Approximate $1/N$ as 0.0
 - Cache is not even big enough to hold multiple rows

■ Analysis Method:

- Look at access pattern of inner loop

$$A \times B = C$$



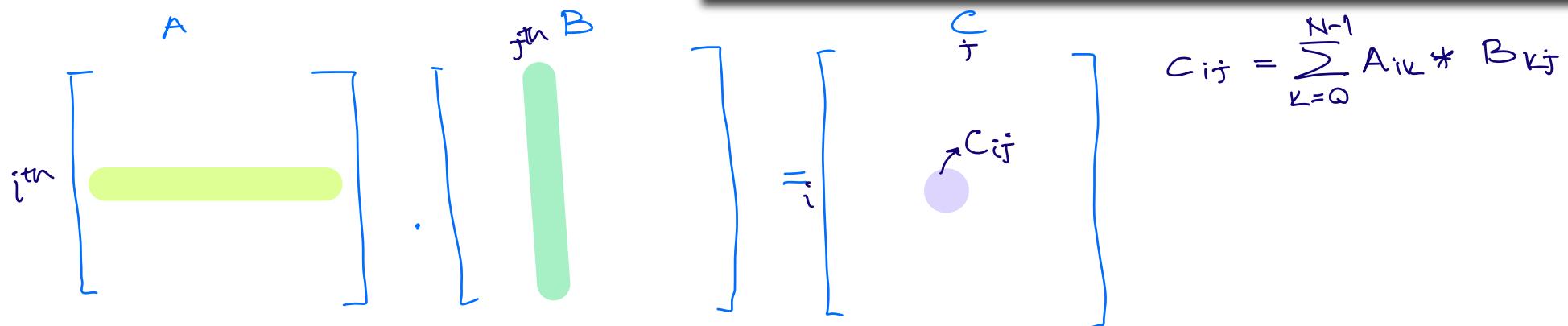
Matrix Multiplication Example

Variable sum held in register

■ Description:

- Multiply $N \times N$ matrices
- $O(N^3)$ total operations
- N reads per source element
- N values summed per destination
 - but may be able to hold in register

```
/* ijk */
for (i=0; i<n; i++) {
    for (j=0; j<n; j++) {
        sum = 0.0; ←
        for (k=0; k<n; k++)
            sum += a[i][k] * b[k][j];
        c[i][j] = sum;
    }
}
```



Layout of C Arrays in Memory (review)

C arrays allocated in row-major order

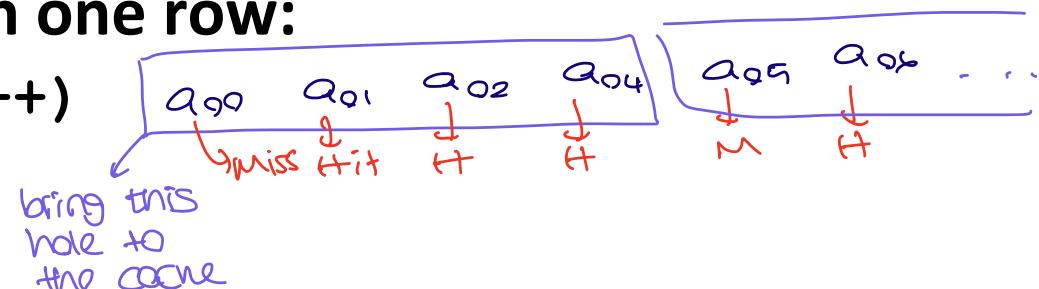
- each row in contiguous memory locations



Stepping through columns in one row:

$B = 32 \text{ bytes}$
 $|a_{00}| = 8 \text{ bytes}$

```
for (i = 0; i < N; i++)
    sum += a[0][i];
```



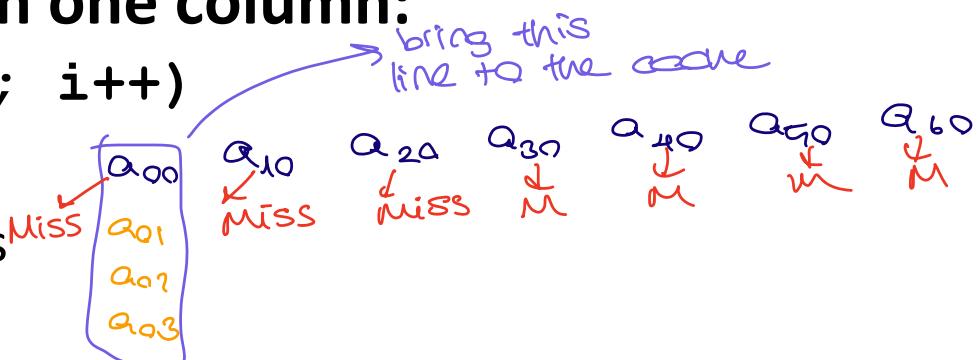
- accesses successive elements
- if block size (B) > 4 bytes, exploit spatial locality

compulsory miss rate = $\frac{8}{4}$ bytes / $B = \frac{8}{32} = \frac{1}{4}$ miss rate

Stepping through rows in one column:

```
for (i = 0; i < n; i++)
    sum += a[i][0];
```

- accesses distant elements
- no spatial locality!

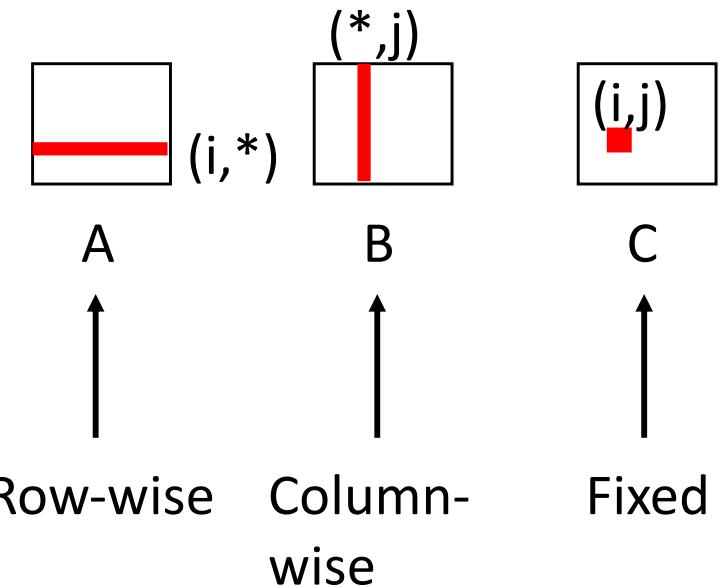


compulsory miss rate = 1 (i.e. 100%)

Matrix Multiplication (ijk)

```
/* ijk */
for (i=0; i<n; i++) {
    for (j=0; j<n; j++) {
        sum = 0.0;
        for (k=0; k<n; k++)
            sum += a[i][k] * b[k][j];
        c[i][j] = sum;
    }
}
```

Inner loop:



Misses per inner loop iteration:

A
0.25

B
1.0

C
0.0

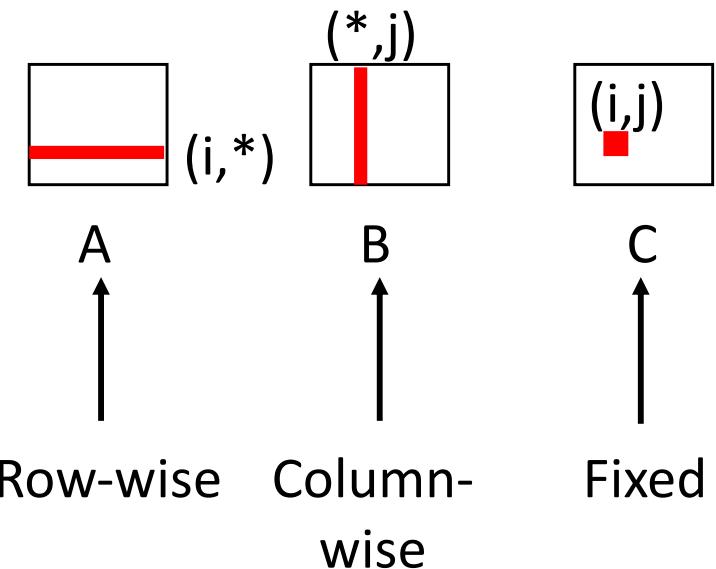
→ I'm accessing this
1 time
(negligible)

$|a_{ij}| = 8$ bytes
Block size = 32 bytes

Matrix Multiplication (jik)

```
/* jik */
for (j=0; j<n; j++) {
    for (i=0; i<n; i++) {
        sum = 0.0;
        for (k=0; k<n; k++)
            sum += a[i][k] * b[k][j];
        c[i][j] = sum
    }
}
```

Inner loop:



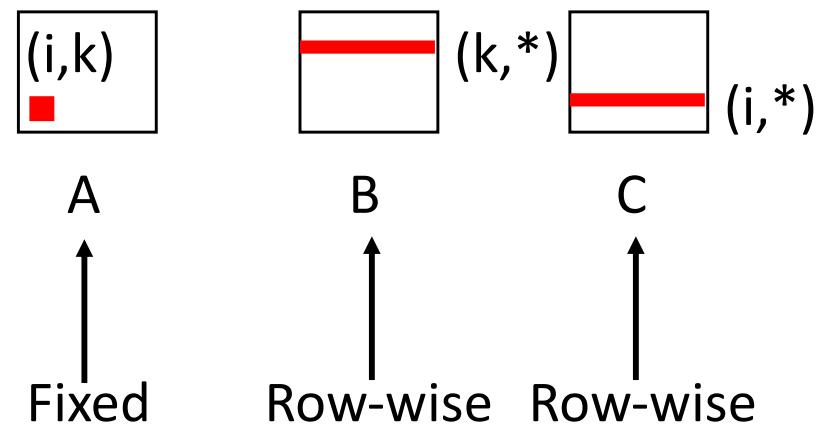
Misses per inner loop iteration:

<u>A</u>	<u>B</u>	<u>C</u>
0.25	1.0	0.0

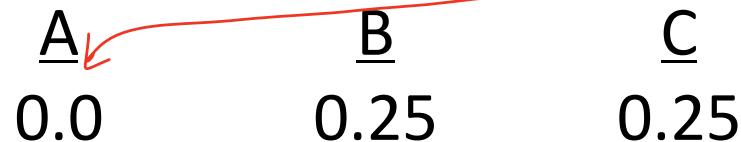
Matrix Multiplication (kij)

```
/* kij */
for (k=0; k<n; k++) {
    for (i=0; i<n; i++) {
        r = a[i][k];
        for (j=0; j<n; j++)
            c[i][j] += r * b[k][j];
    }
}
```

Inner loop:



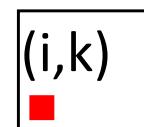
Misses per inner loop iteration:



Matrix Multiplication (ikj)

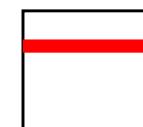
```
/* ikj */
for (i=0; i<n; i++) {
    for (k=0; k<n; k++) {
        r = a[i][k];
        for (j=0; j<n; j++)
            c[i][j] += r * b[k][j];
    }
}
```

Inner loop:



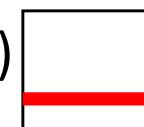
A

Fixed



B

Row-wise



C

Row-wise

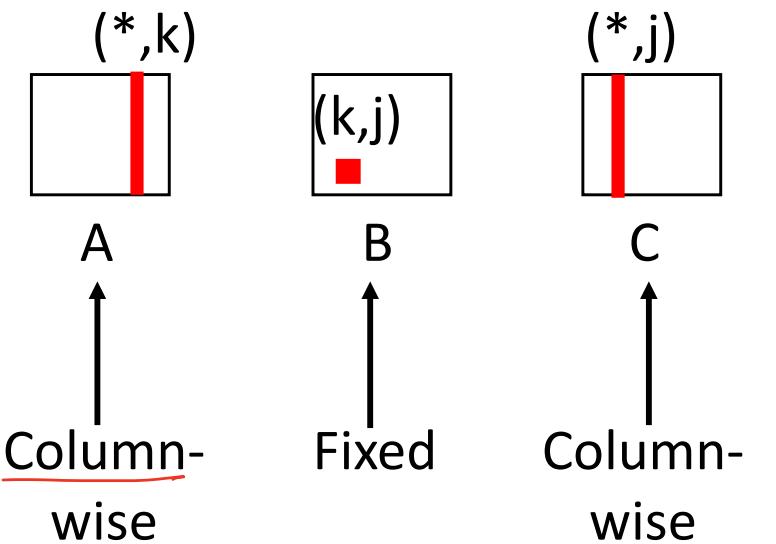
Misses per inner loop iteration:

<u>A</u>	<u>B</u>	<u>C</u>
0.0	0.25	0.25

Matrix Multiplication (jki)

```
/* jki */
for (j=0; j<n; j++) {
    for (k=0; k<n; k++) {
        r = b[k][j];
        for (i=0; i<n; i++)
            c[i][j] += a[i][k] * r;
    }
}
```

Inner loop:



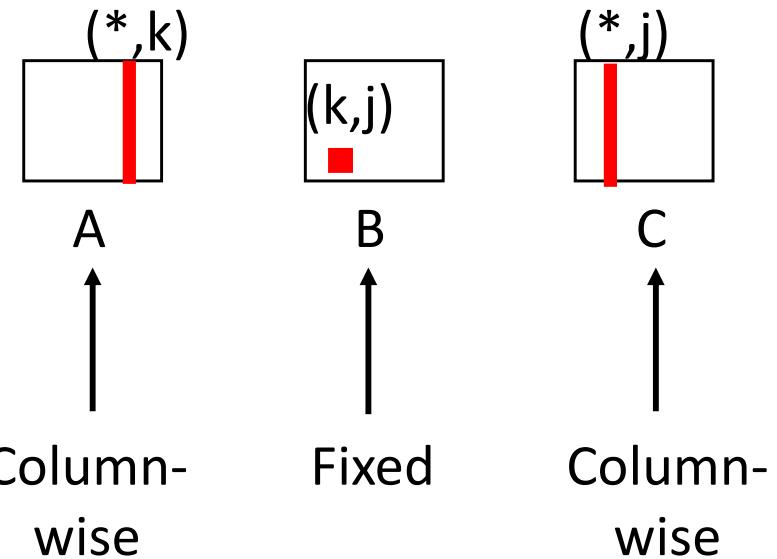
Misses per inner loop iteration:

<u>A</u>	<u>B</u>	<u>C</u>
1.0	0.0	1.0

Matrix Multiplication (kji)

```
/* kji */
for (k=0; k<n; k++) {
    for (j=0; j<n; j++) {
        r = b[k][j];
        for (i=0; i<n; i++)
            c[i][j] += a[i][k] * r;
    }
}
```

Inner loop:



Misses per inner loop iteration:

<u>A</u>	<u>B</u>	<u>C</u>
1.0	0.0	1.0

Summary of Matrix Multiplication

```

for (i=0; i<n; i++) {
    for (j=0; j<n; j++) {
        sum = 0.0;
        for (k=0; k<n; k++)
            sum += a[i][k] * b[k][j];
        c[i][j] = sum;
    }
}

```

$\Theta(n^3)$

```

for (k=0; k<n; k++) {
    for (i=0; i<n; i++) {
        r = a[i][k];
        for (j=0; j<n; j++)
            c[i][j] += r * b[k][j];
    }
}

```

```

for (j=0; j<n; j++) {
    for (k=0; k<n; k++) {
        r = b[k][j];
        for (i=0; i<n; i++)
            c[i][j] += a[i][k] * r;
    }
}

```

ijk (& jik):

- 2 loads, 0 stores
- misses/iter = **1.25**

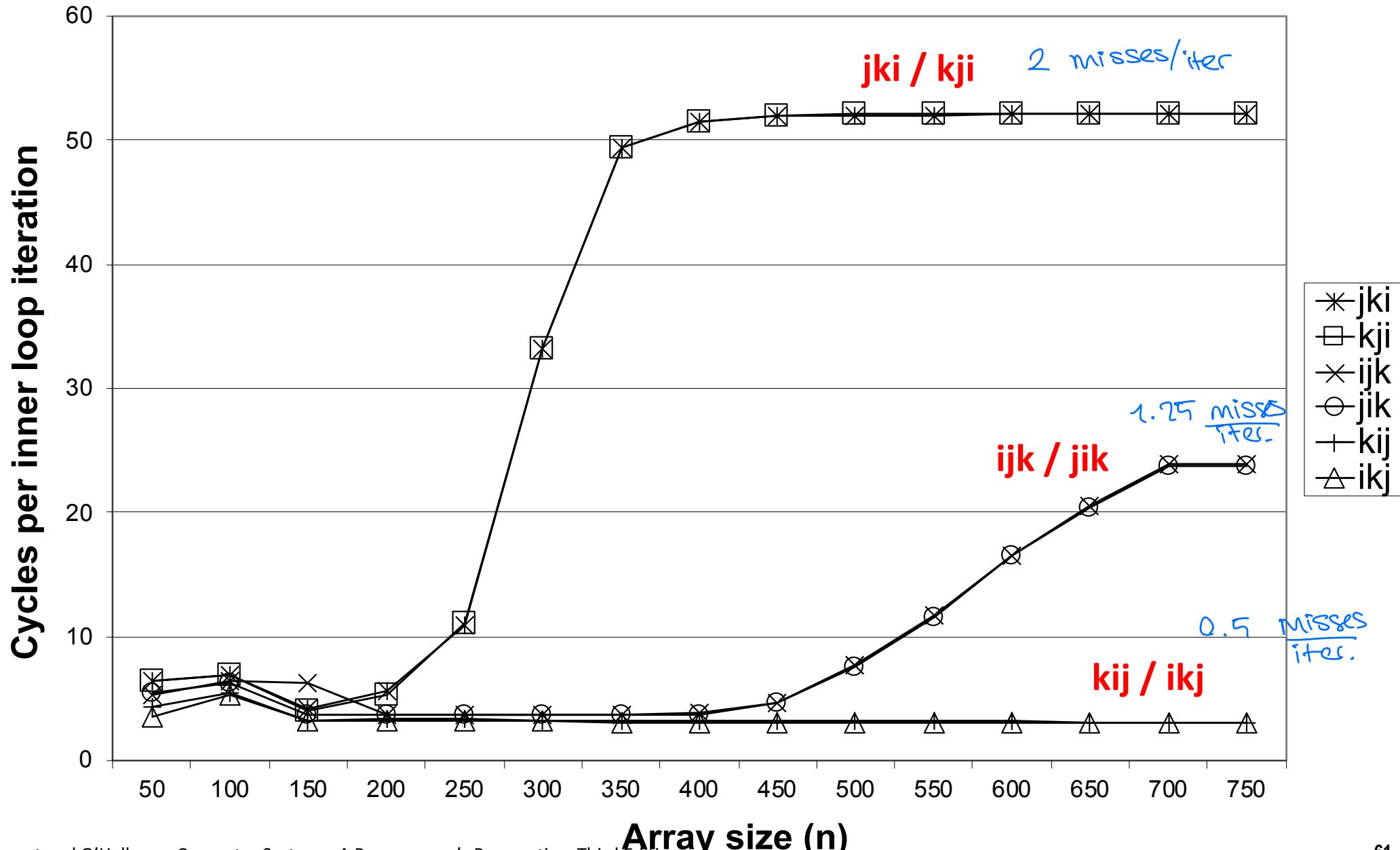
kij (& ikj):

- 2 loads, 1 store
- misses/iter = **0.5**

jki (& kji):

- 2 loads, 1 store
- misses/iter = **2.0**

Core i7 Matrix Multiply Performance

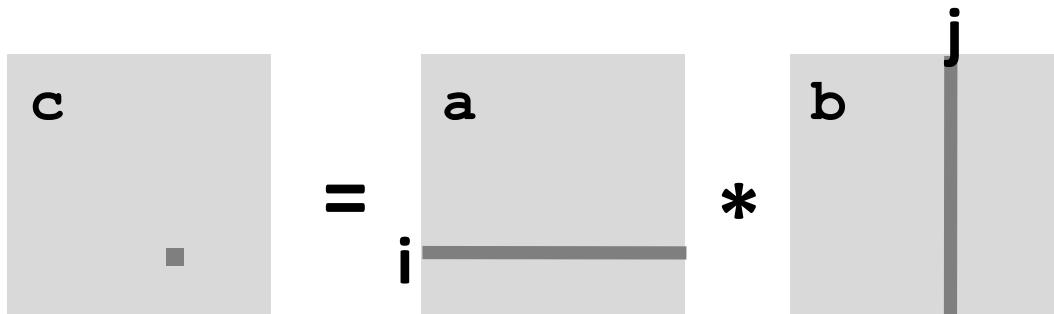


Today

- Cache organization and operation
- Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

Example: Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);  
  
/* Multiply n x n matrices a and b */  
void mmm(double *a, double *b, double *c, int n) {  
    int i, j, k;  
    for (i = 0; i < n; i++)  
        for (j = 0; j < n; j++)  
            for (k = 0; k < n; k++)  
                c[i*n+j] += a[i*n + k]*b[k*n + j];  
}
```



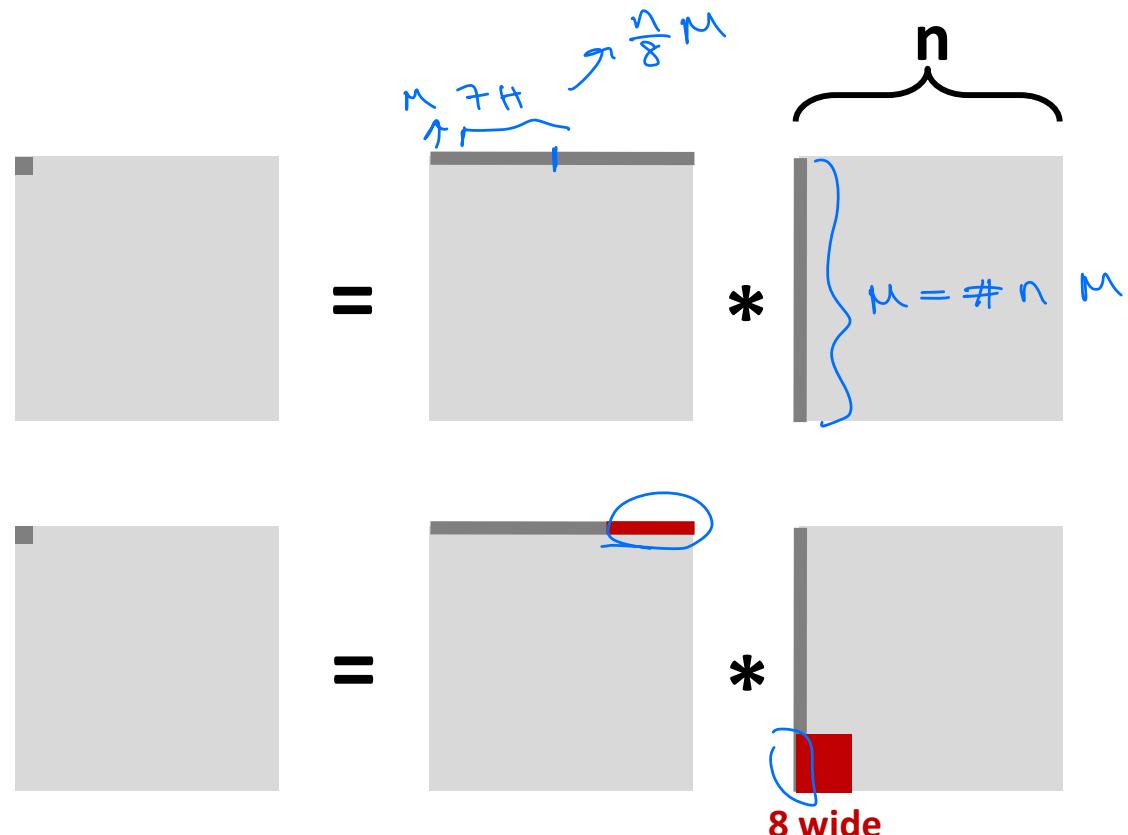
Cache Miss Analysis

■ Assume:

- Matrix elements are doubles
- Cache block = 8 doubles *b4 bytes*
- Cache size C << n (much smaller than n)

■ First iteration:

- $\frac{n}{8} + n = 9n/8$ misses
 A B



- Afterwards in cache:
 (schematic)

Cache Miss Analysis

2^{48}

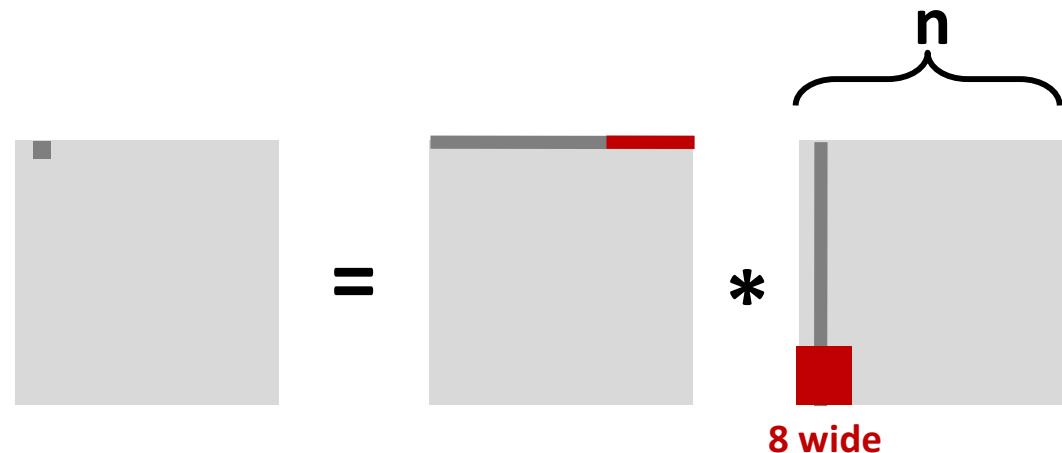
■ Assume:

- Matrix elements are doubles
- Cache block = 8 doubles
- Cache size $C \ll n$ (much smaller than n)

■ Second iteration:

- Again:
 $n/8 + n = 9n/8$ misses

for each $\frac{n}{8}$ size row,
 there will be 1 miss



8 wide

■ Total misses:

- $9n/8 * n^2 = (9/8) * n^3$ misses

of elements in C

$$\begin{array}{c}
 \text{C} \\
 \left[\begin{array}{c|c}
 C_{11} & C_{12} \\ \hline
 \hline
 C_{21} & C_{22} \end{array} \right] = \left[\begin{array}{c|c}
 A_{11} & A_{12} \\ \hline
 \hline
 A_{21} & A_{22} \end{array} \right] \times \left[\begin{array}{c|c}
 B_{11} & B_{12} \\ \hline
 \hline
 B_{21} & B_{22} \end{array} \right]
 \end{array}$$

A_{ij}, B_{ij}, C_{ij} are also matrices

$$C_{11} = A_{11} * B_{11} + A_{12} * B_{21}$$

$$C_{12} = A_{11} * B_{12} + A_{12} * B_{22}$$

$$C_{21} = A_{21} * B_{11} + A_{22} * B_{21}$$

$$C_{22} = A_{21} * B_{12} + A_{22} * B_{22}$$

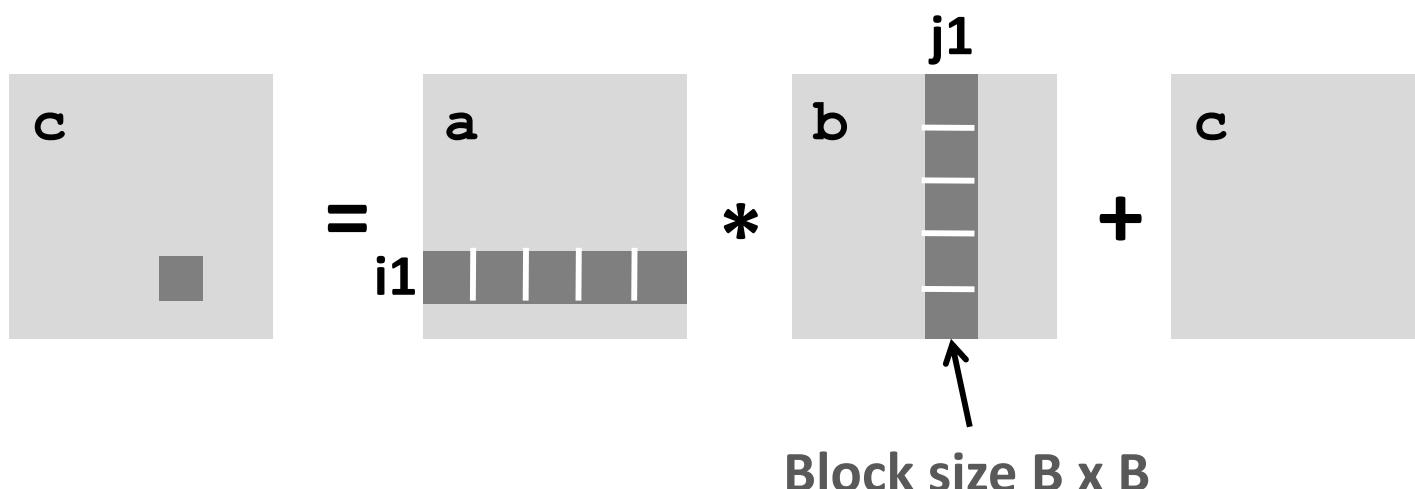
Loop unrolling
 16×16

pointer kullen
 \hookrightarrow pointer offif

Blocked Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);

/* Multiply n x n matrices a and b */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i+=B)
        for (j = 0; j < n; j+=B)
            for (k = 0; k < n; k+=B)
                /* B x B mini matrix multiplications */
                for (i1 = i; i1 < i+B; i++)
                    for (j1 = j; j1 < j+B; j++)
                        for (k1 = k; k1 < k+B; k++)
                            c[i1*n+j1] += a[i1*n + k1]*b[k1*n + j1];
}
```



Cache Miss Analysis

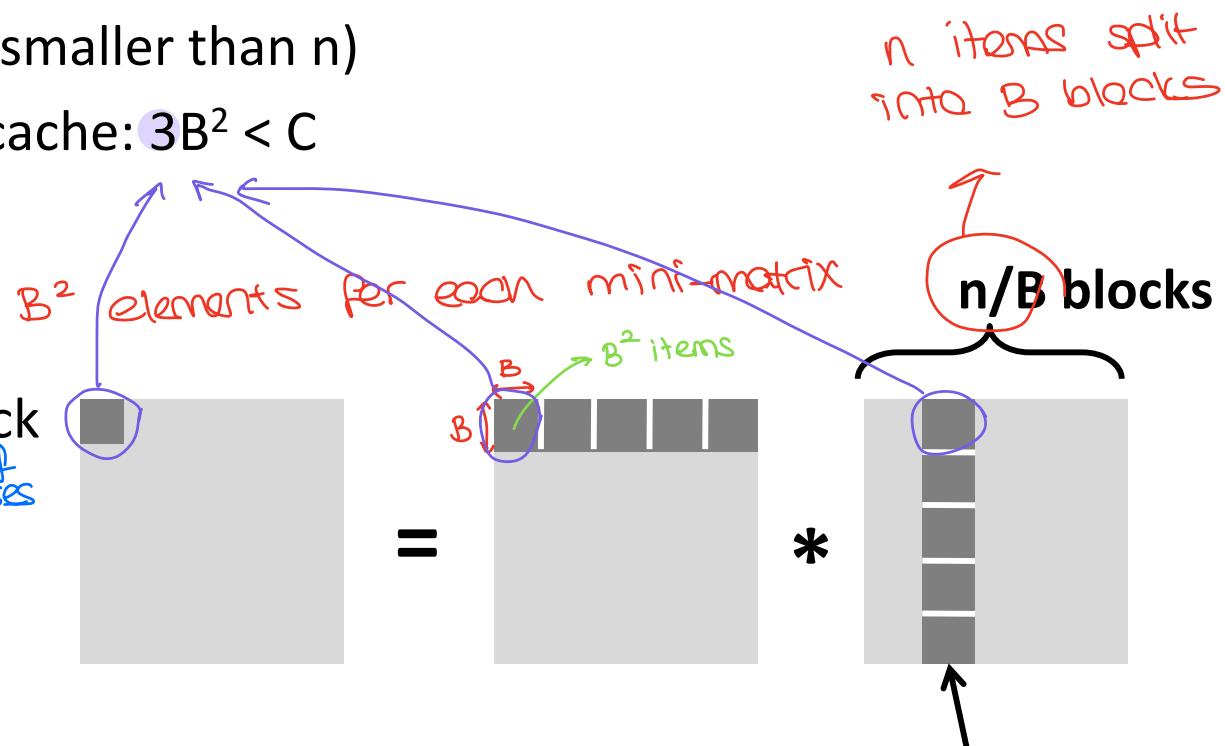
■ Assume:

- Cache block = 8 doubles
- Cache size $C \ll n$ (much smaller than n)
- Three blocks ■ fit into cache: $3B^2 < C$

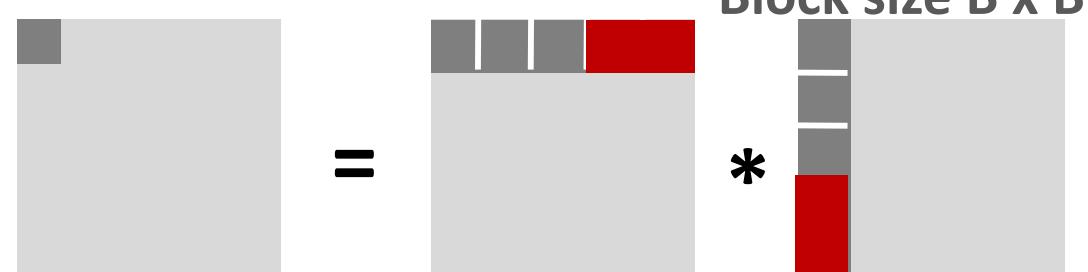
■ First (block) iteration:

- $B^2/8$ misses for each block
- $(2n/B * B^2/8) = nB/4$ $\xrightarrow{\text{# of misses}}$
(omitting matrix c)

we have 2 $\frac{n}{B}$ size block



- Afterwards in cache (schematic)



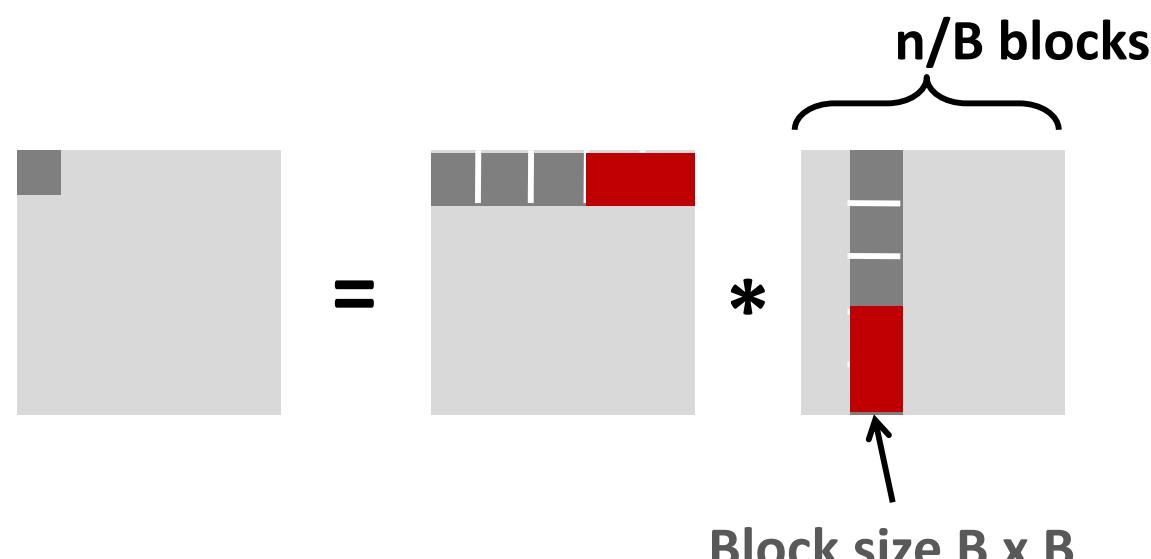
Cache Miss Analysis

■ Assume:

- Cache block = 8 doubles
- Cache size $C \ll n$ (much smaller than n)
- Three blocks ■ fit into cache: $3B^2 < C$

■ Second (block) iteration:

- Same as first iteration
- $2n/B * B^2/8 = nB/4$



■ Total misses:

$$\text{# of misses} = \frac{nB}{4} \times (\frac{n}{B})^2 = \frac{n^3}{(4B)}$$

Annotations explain the terms:

- $nB/4$: # of misses
- $(n/B)^2$: # of blocks in matrix C
- $n^3/(4B)$: mini matrices

Summary

■ No blocking: $(9/8) * n^3 \rightarrow \frac{9}{8} n^3$

■ Blocking: $1/(4B) * n^3 \rightarrow \frac{1}{256} n^3$

$\frac{1}{4*64} n^3$ $B = b4$

$3B^2 < 2^{17}$

$B^2 < \frac{2^{15}}{3}$

$B^2 < 2^{13}$

$B = 2^6 = 64$

■ Suggest largest possible block size B, but limit $3B^2 < C$!

■ Reason for dramatic difference:

- Matrix multiplication has inherent temporal locality:
 - Input data: $3n^2$, computation $2n^3$
 - Every array elements used $O(n)$ times!
- But program has to be written properly

Concluding Observations

- **Programmer can optimize for cache performance**
 - How data structures are organized
 - How data are accessed
 - Nested loop structure
 - Blocking is a general technique
- **All systems favor “cache friendly code”**
 - Getting absolute optimum performance is very platform specific
 - Cache sizes, line sizes, associativities, etc.
 - Can get most of the advantage with generic code
 - Keep working set reasonably small (temporal locality)
 - Use small strides (spatial locality)