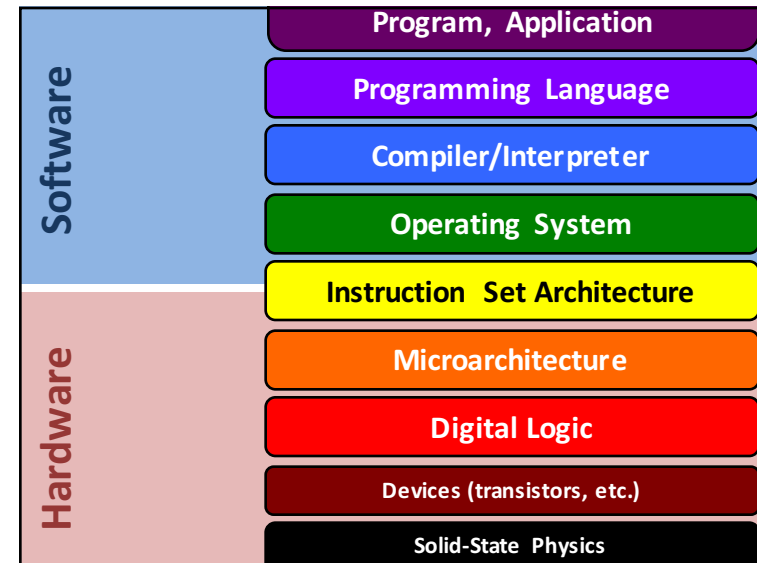


## Memory Hierarchy: Cache

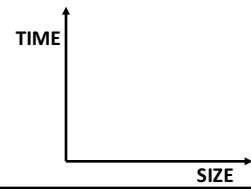
Memory hierarchy  
 Cache basics  
 Locality  
 Cache organization  
 Cache-aware programming



## How does execution time grow with SIZE?

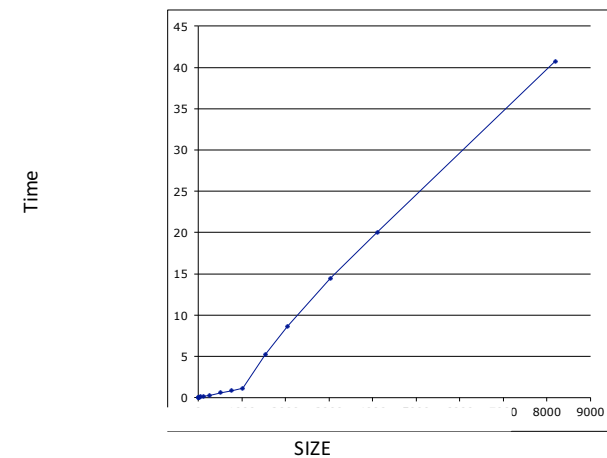
```
int[] array = new int[SIZE];
fillArrayRandomly(array);
int s = 0;

for (int i = 0; i < 200000; i++) {
  for (int j = 0; j < SIZE; j++) {
    s += array[j];
  }
}
```

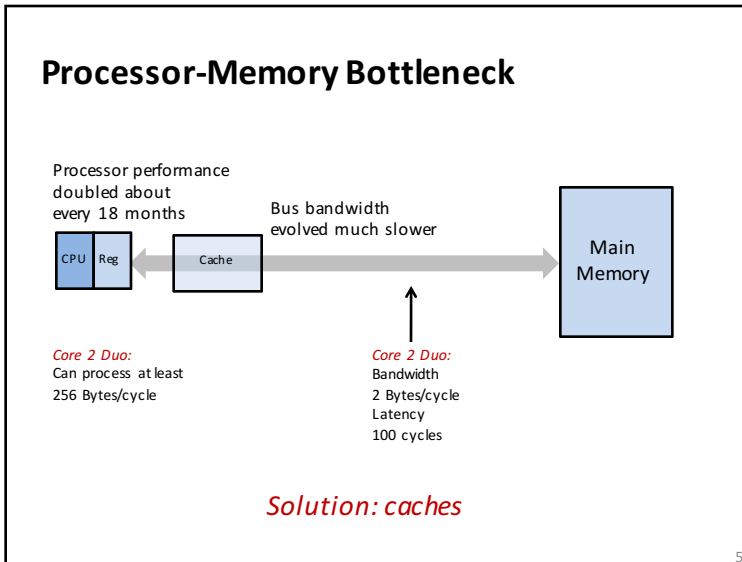


3

## reality beyond $O(\dots)$



4

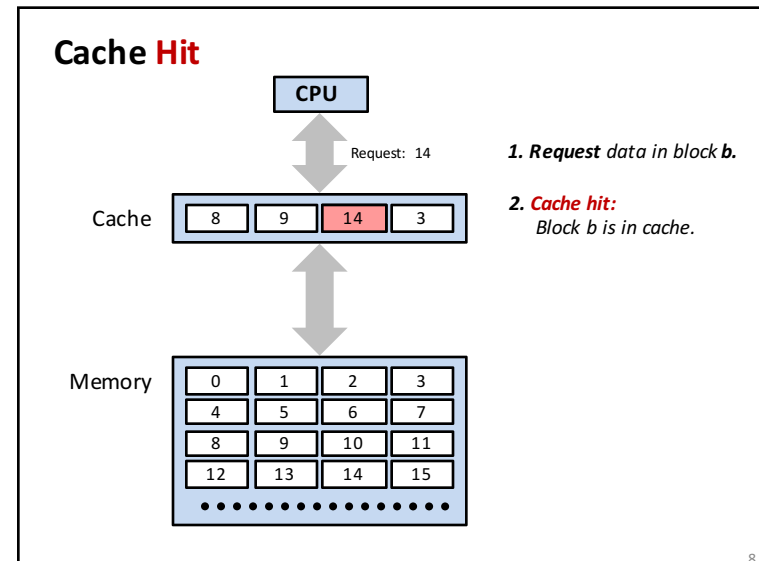
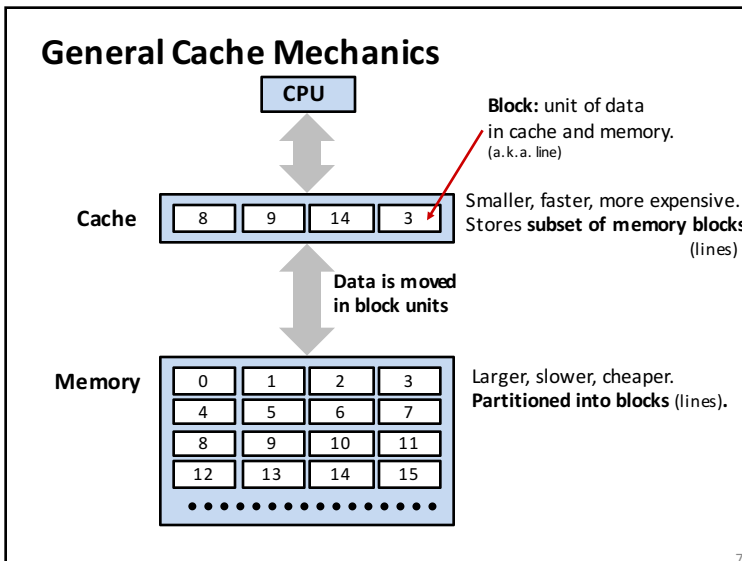


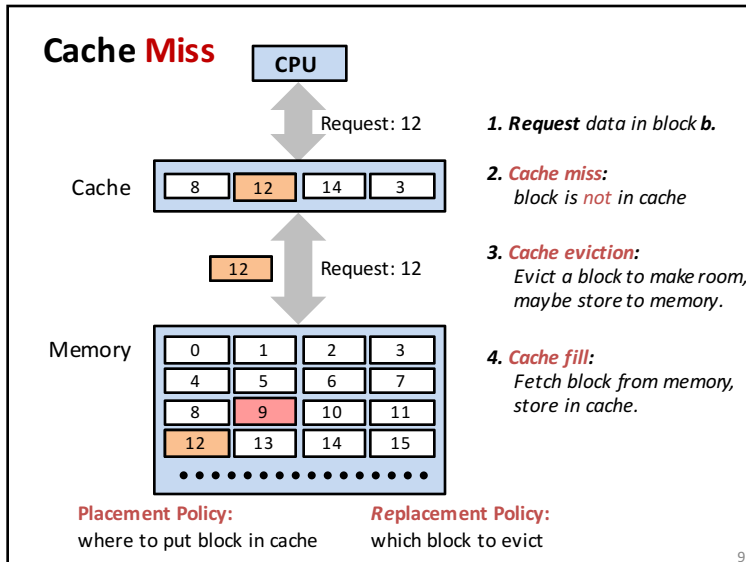
### Cache

**English:**  
*n.* a hidden storage space for provisions, weapons, or treasures  
*v.* to store away in hiding for future use

**Computer Science:**  
*n.* a computer memory with short access time used to store frequently or recently used instructions or data  
*v.* to store [data/instructions] temporarily for later quick retrieval

Also used more broadly in CS: software caches, file caches, etc.





### Locality: why caches work

Programs tend to use data and instructions at 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 are *likely* to be referenced close together in time.

**How do caches exploit temporal and spatial locality?**

### Example: Locality?

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

What is stored in memory?

**Data:**  
Temporal: **sum** referenced in each iteration  
Spatial: array **a** accessed in *stride-1* pattern

**Instructions:**  
Temporal: execute loop repeatedly  
Spatial: execute instructions in sequence

**Assessing locality in code is an important programming skill.**

### Locality Example #1

row-major M x N 2D array in C

```
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 #2

row-major M x N 2D array in C

```

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;
}
    
```

a[0][0]	a[0][1]	a[0][2]	a[0][3]
a[1][0]	a[1][1]	a[1][2]	a[1][3]
a[2][0]	a[2][1]	a[2][2]	a[2][3]

13

### Locality Example #3

```

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

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

What is "wrong" with this code?  
How can it be fixed?

14

### Cost of Cache Misses

**Huge difference between a hit and a miss**  
Could be 100x, if just L1 and main memory

**99% hits could be twice as good as 97%. How?**  
Cache hit time of 1 cycle, miss penalty of 100 cycles

Mean access time:

- 97% hits: 1 cycle + 0.03 \* 100 cycles = 4 cycles
- 99% hits: 1 cycle + 0.01 \* 100 cycles = 2 cycles

hit/miss rate

**This is why "miss rate" is used instead of "hit rate"**

15

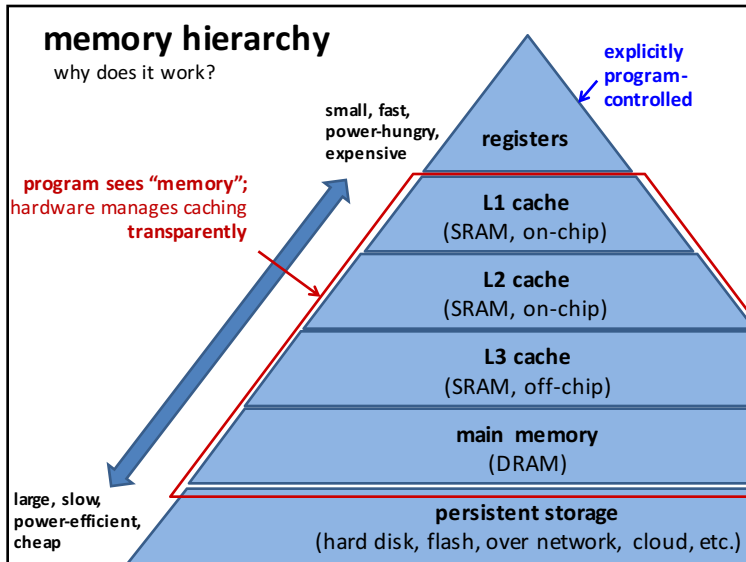
### Cache Performance Metrics

**Miss Rate**  
Fraction of memory accesses to data not in cache (misses / accesses)  
Typically: 3% - 10% for L1; maybe < 1% for L2, depending on size, etc.

**Hit Time**  
Time to find and deliver a block in the cache to the processor.  
Typically: 1 - 2 clock cycles for L1; 5 - 20 clock cycles for L2

**Miss Penalty**  
Additional time required on cache miss = main memory access time  
Typically 50 - 200 cycles for L2 (trend: increasing!)

16



### Cache Organization: Key Points

**Block**  
Fixed-size **unit of data** in memory/cache

**Placement Policy**  
Where should a given block be stored in the cache?  

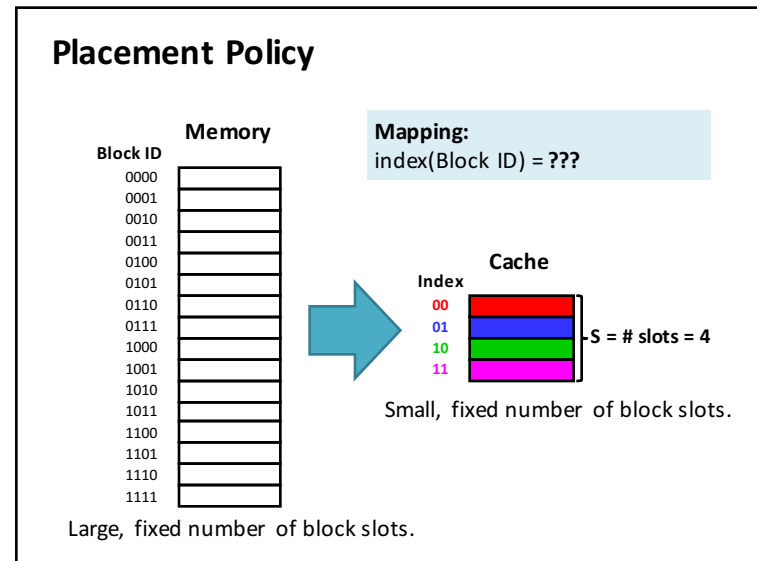
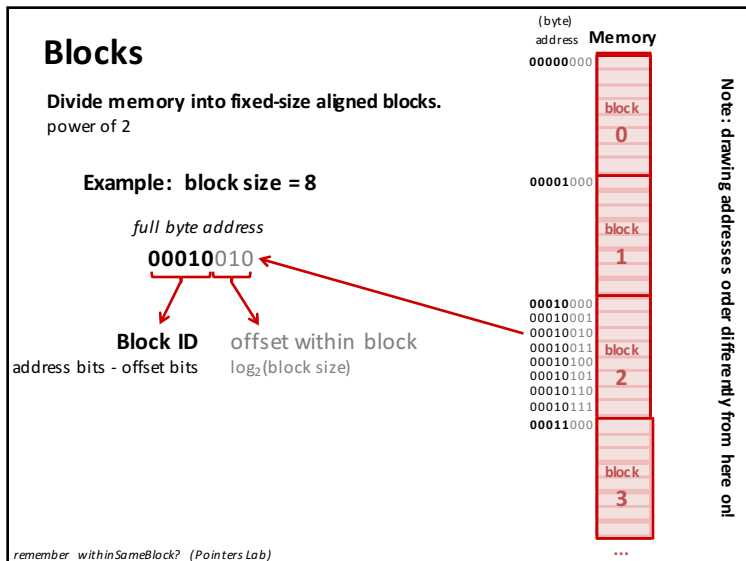
- direct-mapped, set associative

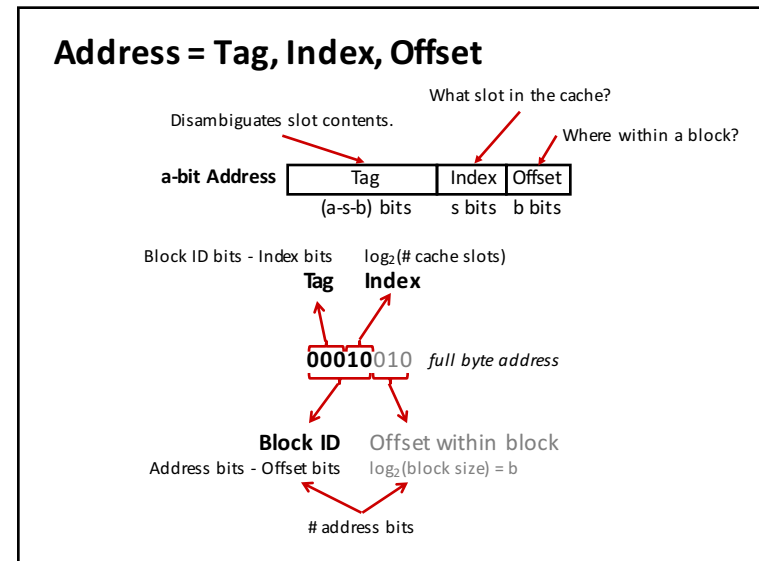
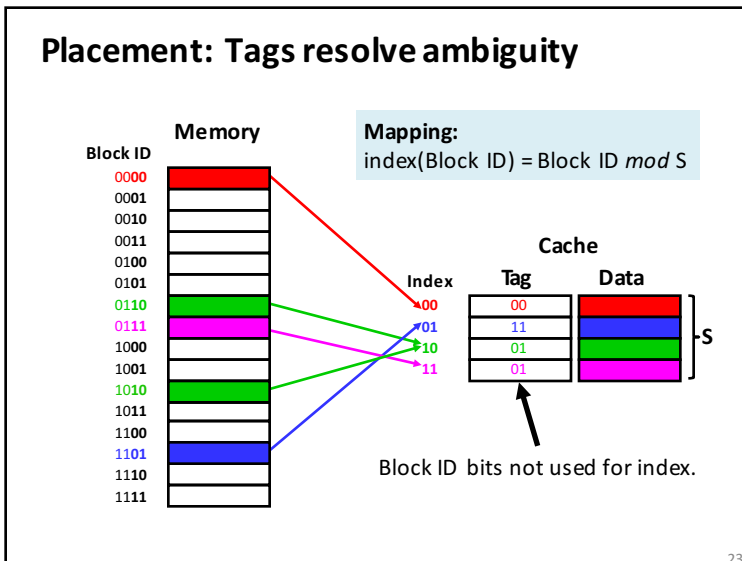
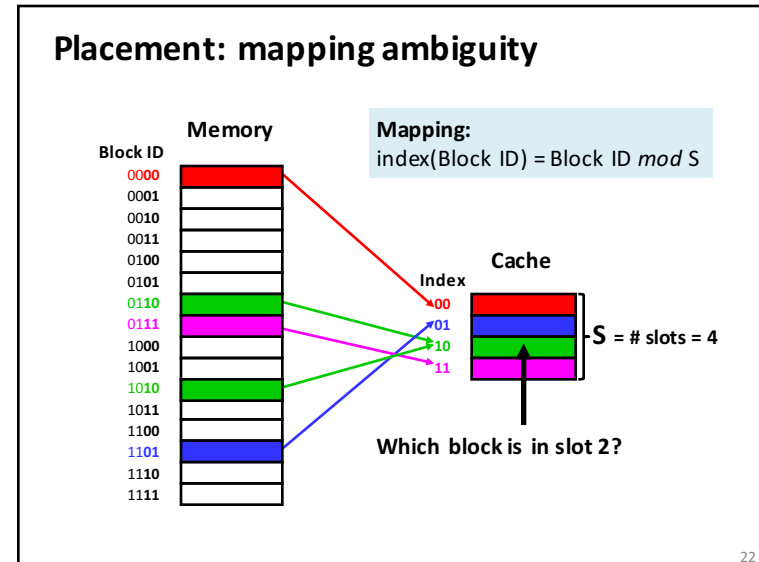
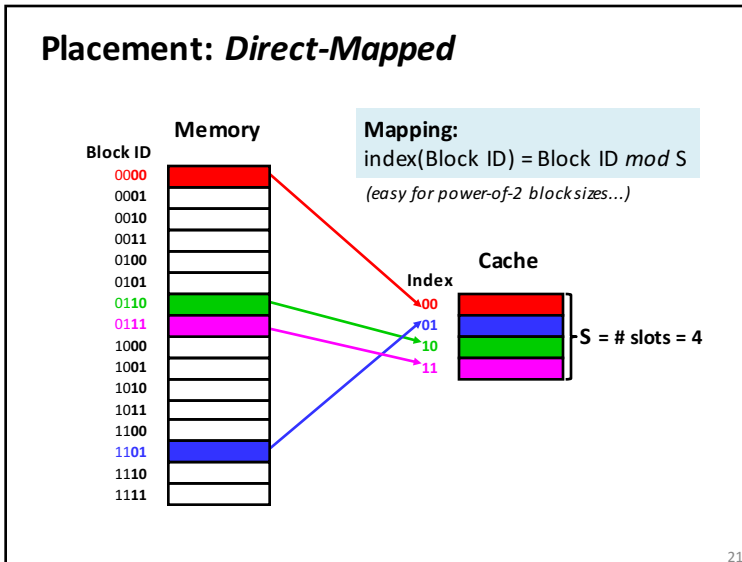
**Replacement Policy**  
What if there is no room in the cache for requested data?  

- least recently used, most recently used

**Write Policy**  
When should writes update lower levels of memory hierarchy?  

- write back, write through, write allocate, no write allocate





### Placement: ~~Direct-Mapped~~

**Memory**

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

**Why not this mapping?**  
 $\text{index}(\text{Block ID}) = \text{Block ID} / S$   
*(still easy for power-of-2 block sizes...)*

**Cache**

Index	
00	
01	
10	
11	

25

### A puzzle.

Cache starts *empty*.  
 Access (address, hit/miss) stream:

(10, miss), (11, hit), (12, miss)

What could the block size be?

26

### Placement: direct mapping conflicts

Block ID

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

**What happens when accessing in repeated pattern:**  
**0010, 0110, 0010, 0110, 0010...?**

**cache conflict**  
 Every access suffers a miss, evicts cache line needed by next access.

Index

00	
01	
10	
11	

27

### Placement: Set Associative

**sets**  
 $S = \# \text{ Slots in cache}$

Index per **set** of block slots.  
 Store block in **any** slot within set.

**Mapping:**  
 $\text{index}(\text{Block ID}) = \text{Block ID} \bmod S$

**1-way**  
8 sets,  
1 block each

Set	
0	
1	
2	
3	
4	
5	
6	
7	

**direct mapped**

**2-way**  
4 sets,  
2 blocks each

Set	
0	.....
1	.....
2	.....
3	.....

**4-way**  
2 sets,  
4 blocks each

Set	
0	.....
1	.....

**8-way**  
1 set,  
8 blocks

Set	
0	.....

**fully associative**

28

### Example: Tag, Index, Offset?

4-bit Address 

Tag	Index	Offset
-----	-------	--------

Direct-mapped  
4 slots  
2-byte blocks

tag bits \_\_\_\_\_  
set index bits \_\_\_\_\_  
block offset bits \_\_\_\_\_

index(1101) = \_\_\_\_\_

### Example: Tag, Index, Offset?

E-way set-associative  
S slots  
16-byte blocks

16-bit Address 

Tag	Index	Offset
-----	-------	--------

E = 1-way  
S = 8 sets

E = 2-way  
S = 4 sets

E = 4-way  
S = 2 sets

tag bits \_\_\_\_\_  
set index bits \_\_\_\_\_  
block offset bits \_\_\_\_\_  
index(0x1833) \_\_\_\_\_

### Replacement Policy

If set is full, what block should be replaced?  
Common: **least recently used (LRU)**  
(hardware usually implements "not most recently used")

1-way associativity  
8 sets, 1 block each

direct mapped

2-way associativity  
4 sets, 2 blocks each

4-way associativity  
2 sets, 4 blocks each

31

### Another puzzle.

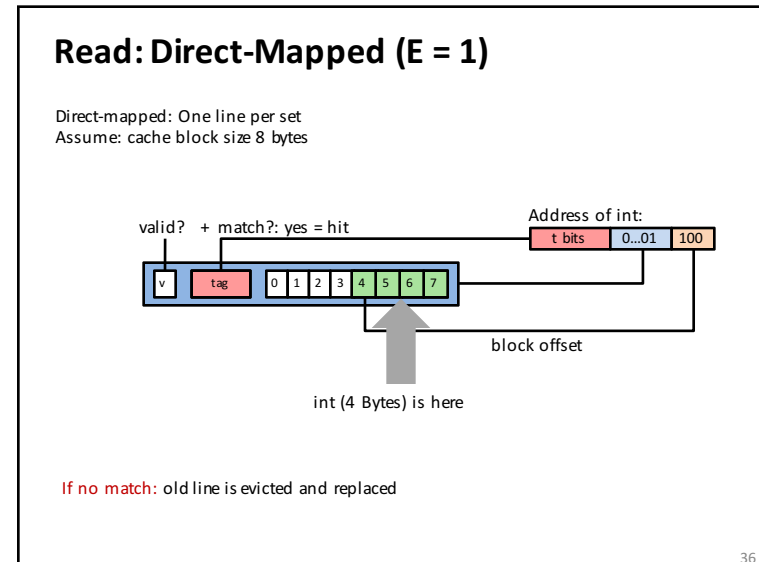
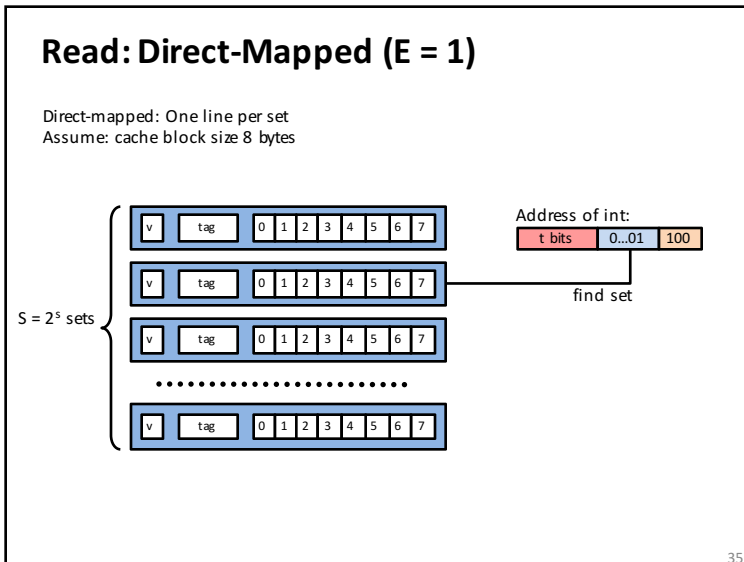
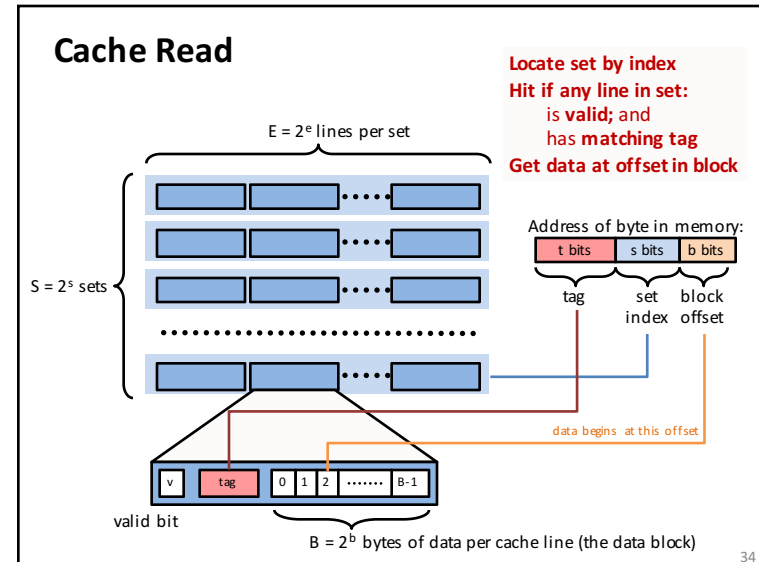
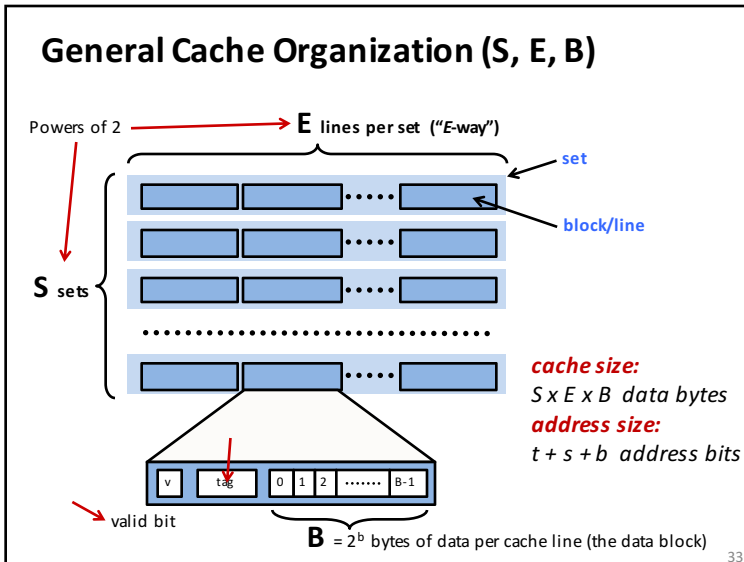
Cache starts *empty*, uses LRU.  
Access (address, hit/miss) stream

(10, miss); (12, miss); (10, miss)

associativity of cache?

32





### Example (E = 1)

Assume *sum, i, j* in registers  
Address of an aligned element of *a*: aa...arrrrcccc000

```

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;
}
    
```

Assume: cold (empty) cache  
3-bit set index, 5-bit offset  
aa...arr rcc cc000  
0,0: aa...a000 000 00000

0,0	0,1	0,2	0,3
0,4	0,5	0,6	0,7
0,8	0,9	0,a	0,b
0,c	0,d	0,e	0,f
1,0	1,1	1,2	1,3
1,4	1,5	1,6	1,7
1,8	1,9	1,a	1,b
1,c	1,d	1,e	1,f

0,0	0,1	0,2	0,3

32 bytes = 4 doubles      32 bytes = 4 doubles  
4 misses per row of array      every access a miss  
4\*16 = 64 misses      16\*16 = 256 misses

37

### Example (E = 1)

block = 16 bytes; 8 sets in cache  
How many block offset bits?  
How many set index bits?

```

int dotprod(int x[8], int y[8]) {
  int sum = 0;
  int i;

  for (i = 0; i < 8; i++)
    sum += x[i]*y[i];
  return sum;
}
    
```

Address bits:  
B =  
S =

0:  
128:  
160:

if *x* and *y* have aligned starting addresses, e.g., &x[0] = 0, &y[0] = 128

x[0]	x[1]	x[2]	x[3]

if *x* and *y* have unaligned starting addresses, e.g., &x[0] = 0, &y[0] = 160

x[0]	x[1]	x[2]	x[3]
x[4]	x[5]	x[6]	x[7]
y[0]	y[1]	y[2]	y[3]
y[4]	y[5]	y[6]	y[7]

38

### Read: Set-Associative (Example: E = 2)

E = 2: Two lines per set  
Assume: cache block size 8 bytes

Address of int: t bits 0...01 100

v	tag	0	1	2	3	4	5	6	7
v	tag	0	1	2	3	4	5	6	7
v	tag	0	1	2	3	4	5	6	7
v	tag	0	1	2	3	4	5	6	7
v	tag	0	1	2	3	4	5	6	7
v	tag	0	1	2	3	4	5	6	7

find set

39

### Read: Set-Associative (Example: E = 2)

E = 2: Two lines per set  
Assume: cache block size 8 bytes

Address of int: t bits 0...01 100

compare both

valid? + match: yes = hit

v	tag	0	1	2	3	4	5	6	7
v	tag	0	1	2	3	4	5	6	7

block offset

int (4 Bytes) is here

If no match:

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

40

### Example (E = 2)

```
float dotprod(float x[8], float y[8])
{
    float sum = 0;
    int i;

    for (i = 0; i < 8; i++)
        sum += x[i]*y[i];
    return sum;
}
```

If x and y aligned,  
e.g. &x[0] = 0, &y[0] = 128,  
can still fit both because each set  
has space for two blocks/lines



41

### Types of Cache Misses

#### Cold (compulsory) miss

first access to a block

#### Conflict miss

cache has space for all needed blocks, but multiple blocks map to same slot  
e.g., referencing blocks 0, 8, 0, 8, ... would miss every time  
increasing associativity can reduce conflict misses

#### Capacity miss

*working set* of active cache blocks is larger than the cache

42

### What about writes?

#### Multiple copies of data exist:

L1, L2, possibly L3, main memory

#### Write-hit policy

**Write-through:** write immediately to memory, all caches in between.

**Write-back:** defer write to memory until line is evicted (replaced)

Need a *dirty bit* to indicate if line is different from memory or not

#### Write-miss policy

**Write-allocate:** load into cache, update line in cache.

Good if more nearby writes or reads follow

**No-write-allocate:** just write immediately to memory.

#### Typical caches:

Write-back + Write-allocate, usually

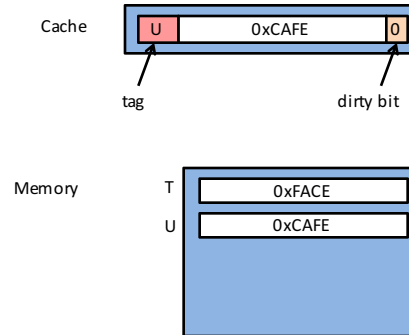
Write-through + No-write-allocate, occasionally

43

### Write-back, write-allocate example

```
eax =
ecx = T
edx = U
```

1. mov \$T, %ecx
2. mov \$U, %edx
3. mov \$0xFEEED, (%ecx)
  - a. Miss on T.



44

### Write-back, write-allocate example

eax =  
ecx = T  
edx = U

1. `mov $T, %ecx`
2. `mov $U, %edx`
3. `mov $0xFEEED, (%ecx)`
  - a. Miss on T.
  - b. Evict U (clean: discard).
  - c. Fill T (write-allocate).
  - d. Write T in cache (dirty).
4. `mov (%edx), %eax`
  - a. Miss on U.

Cache

T	0xFEEED	1
---	---------	---

tag dirty bit

Memory

T	0xFACE
U	0xCAFE

45

### Write-back, write-allocate example

eax = 0xCAFE  
ecx = T  
edx = U

1. `mov $T, %ecx`
2. `mov $U, %edx`
3. `mov $0xFEEED, (%ecx)`
  - a. Miss on T.
  - b. Evict U (clean: discard).
  - c. Fill T (write-allocate).
  - d. Write T in cache (dirty).
4. `mov (%edx), %eax`
  - a. Miss on U.
  - b. Evict T (dirty: write back).
  - c. Fill U.
  - d. Set %eax.
5. **DONE.**

Cache

U	0xCAFE	0
---	--------	---

tag dirty bit

Memory

T	0xFEEED
U	0xCAFE

46

### Example Memory Hierarchy

Processor package

**Intel Core i7 circa 2011**

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.

slower, but more likely to hit

47

### Aside: software caches

**Examples**

- File system buffer caches, web browser caches, database caches, network CDN caches, etc.

**Some design differences**

- Almost always fully-associative
  - so, no placement restrictions
  - index structures like hash tables are common (for placement)
- Often use complex replacement policies
  - misses are very expensive when disk or network involved
  - worth thousands of cycles to avoid them
- Not necessarily constrained to single "block" transfers
  - may fetch or write-back in larger units, opportunistically

48

## Cache-Friendly Code

**Locality, locality, locality.** ← locality

**Programmer can optimize for cache performance**

- Data structure organization

- Data access patterns

  - Nested loops

  - Blocking (see CSAPP 6.5)

**All systems favor “cache-friendly code”**

- Performance is hardware-specific

  - Cache size, line size, associativity, etc.

- Generic rules still capture most of advantages

  - Keep working set small (temporal locality)

  - Use small strides (spatial locality)

  - Focus on inner loop code