# Memory Hierarchy --- Caching

CSCI 2021: Machine Architecture and Organization

Antonia Zhai

Department Computer Science and Engineering
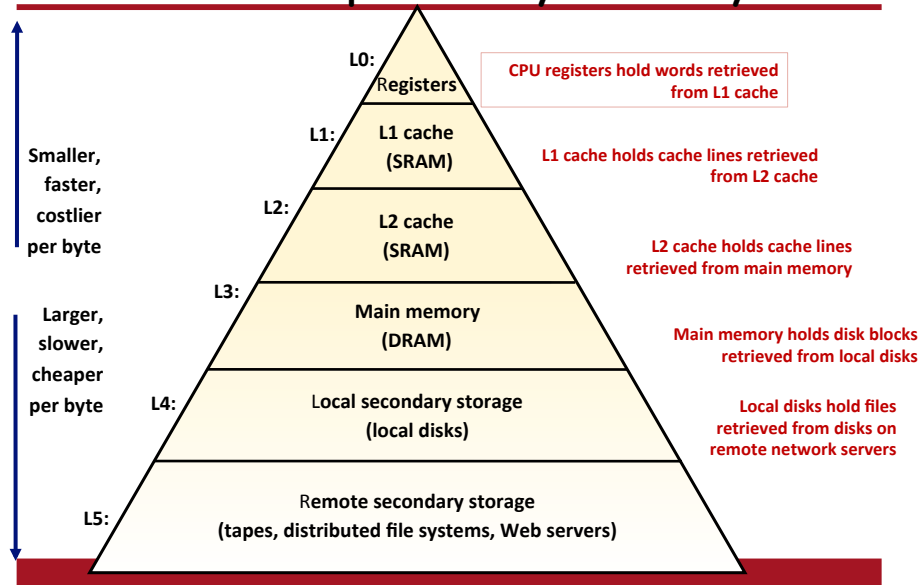
University of Minnesota

http://www.cs.umn.edu/~zhai

**With Slides from Bryant and O'Hallaron**

UNIVERSITY OF MINNESOTA

---

## An Example Memory Hierarchy

Smaller,
faster,
costlier
per byte

Larger,
slower,
cheaper
per byte

L0: Registers — CPU registers hold words retrieved from L1 cache

L1: L1 cache (SRAM) — L1 cache holds cache lines retrieved from L2 cache

L2: L2 cache (SRAM) — L2 cache holds cache lines retrieved from main memory

L3: Main memory (DRAM) — Main memory holds disk blocks retrieved from local disks

L4: Local secondary storage (local disks) — Local disks hold files retrieved from disks on remote network servers

L5: Remote secondary storage (tapes, distributed file systems, Web servers)
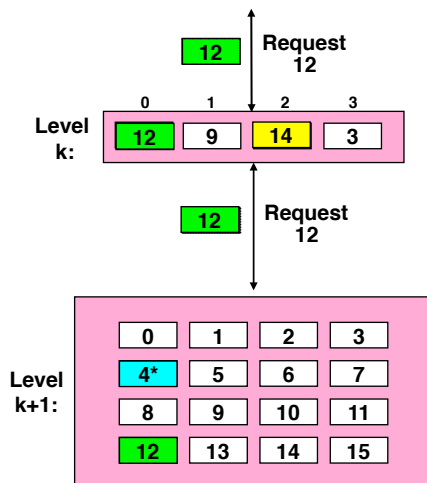
# Caches

- *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 Caching  Concepts



Program needs object d, which is stored in some block b.

Cache hit
- Program finds  b  in the cache at level k. E.g.,  block 14.

Cache miss
- b is not at level k, so level k cache  must fetch it from level k+1.   E.g.,  block 12.
- If level k cache is full, then some current block must be replaced (evicted). Which one is the "victim"?
  - Placement policy: where can the new block go? E.g., b mod 4
  - Replacement policy: which block should be evicted? E.g., LRU

2

# General Caching Concepts: Types of Cache Misses

- Cold (compulsory) miss
  - Cold misses occur because the cache is empty.
- Conflict miss
  - 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 mod 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.

---

# Locality

Principle of Locality:

- Temporal locality: Recently referenced items are likely to be referenced in the near future.

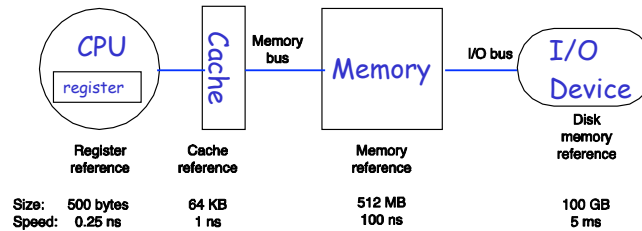- Spatial locality:  Items with nearby addresses tend to be referenced close together in time.

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

## Locality Example:
- **Data**
  - Reference array elements in succession  **Spatial locality**
  - Reference sum each iteration: **Temporal locality**
- **Instructions**
  - Reference instructions in sequence: **Spatial locality**
  - Cycle through loop repeatedly: **Temporal locality**

# Storage Units Organization

- Rule-of-Thumb: The larger the structure, the slower the access time
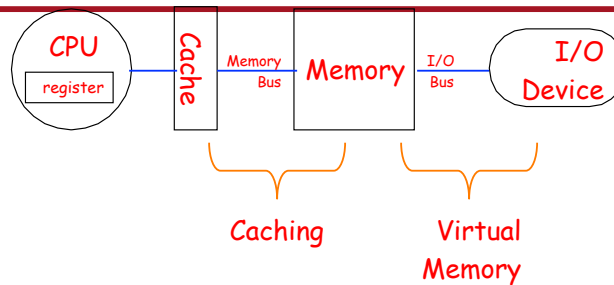- Solution: Hierarchical design



| | | |
|---|---|---|
| | Register reference | Cache reference | Memory reference | Disk memory reference |

| | | Register reference | Cache reference | Memory reference | Disk memory reference |
|---|---|---|---|---|---|
| Size: | 500 bytes | 64 KB | 512 MB | 100 GB |
| Speed: | 0.25 ns | 1 ns | 100 ns | 5 ms |

- There is locality in data access
- Put things that you are likely to use in the near future close to you

---

# Caching the Memory



Caching    Virtual Memory

- Cache uses *SRAM*: Static Random Access Memory
  - No refresh
- Main Memory is *DRAM*: Dynamic Random Access Memory
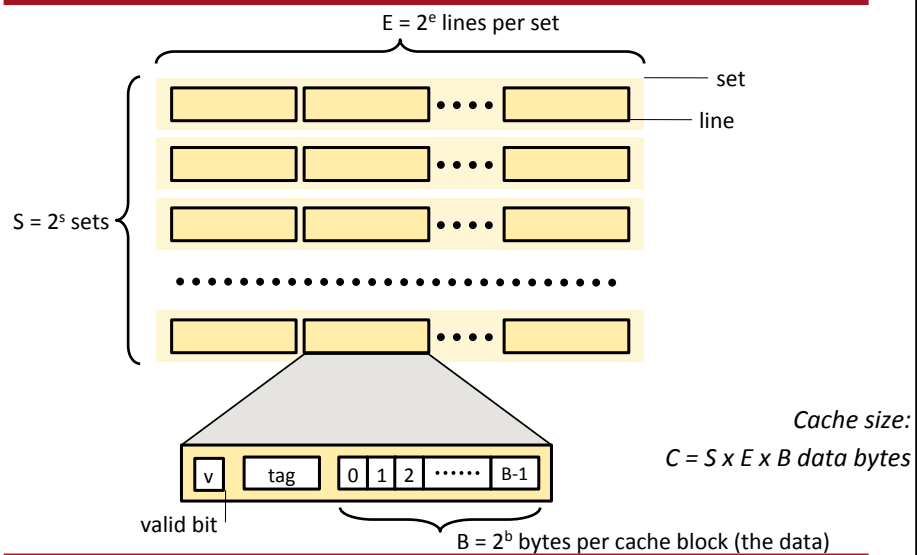  - Dynamic since needs to be refreshed periodically
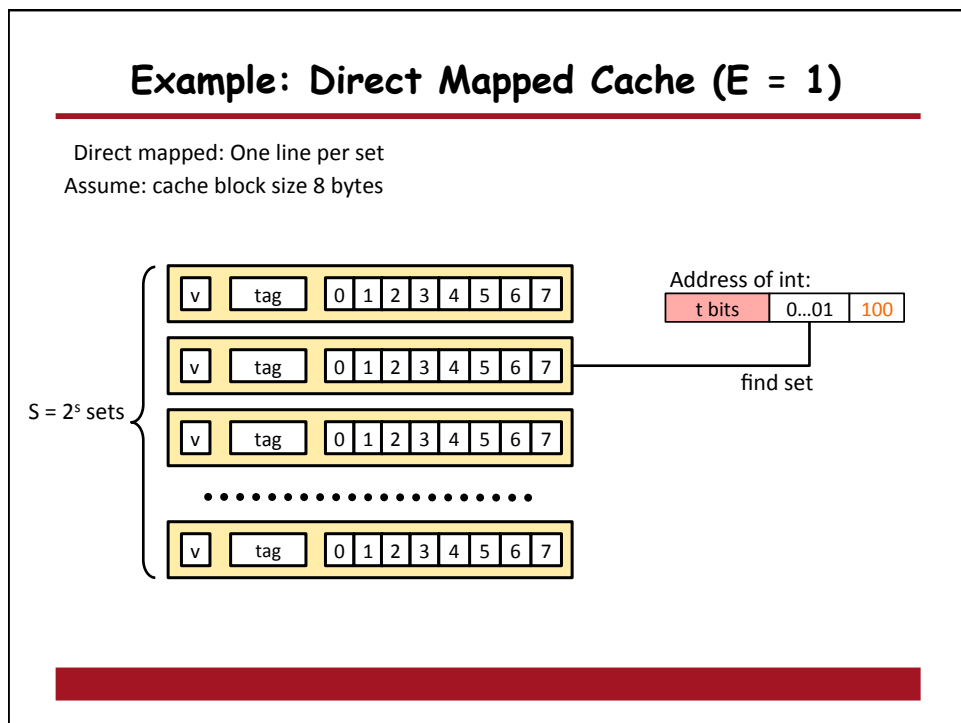
# The 1,2,3,4 of Caching

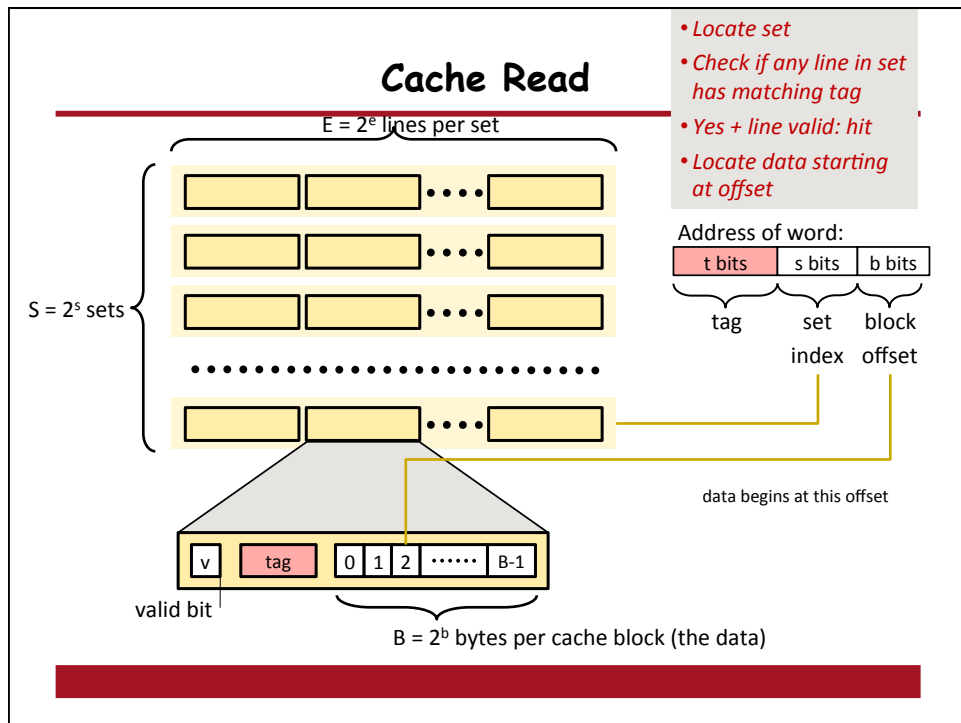1. Where can a block be placed in the upper level?
2. How is a block found if it is in the upper level?
3. Which block should be replaced on a miss?
4. What happens on a write?

# General Cache Organization (S, E, B)

$E = 2^e$ lines per set

set

line

$S = 2^s$ sets

*Cache size:*
*C = S x E x B data bytes*

| v | tag | 0 | 1 | 2 | ...... | B-1 |

valid bit

$B = 2^b$ bytes per cache block (the data)

# Cache Read

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

$E = 2^e$ lines per set

$S = 2^s$ sets

Address of word:

| t bits | s bits | b bits |
|--------|--------|--------|
| tag | set index | block offset |

data begins at this offset

valid bit

$B = 2^b$ bytes per cache block (the data)

v | tag | 0 | 1 | 2 | ...... | B-1

---

# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size 8 bytes

$S = 2^s$ sets

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

Address of int:

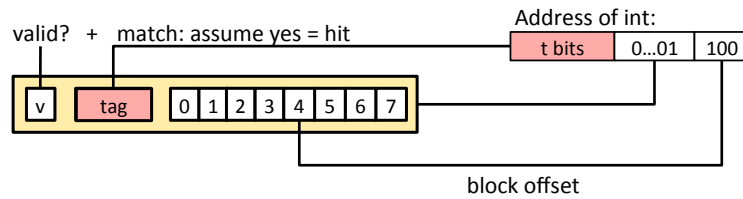| t bits | 0...01 | 100 |
|--------|--------|-----|

find set

6

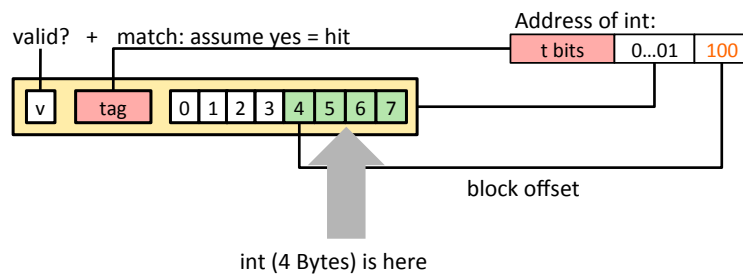## Example: Direct Mapped Cache (E = 1)

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



int (4 Bytes) is here

No match: old line is evicted and replaced

# Direct-Mapped Cache Simulation

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

| t=1 | s=2 | b=1 |
|---|---|---|
| x | xx | x |

Address trace (reads, one byte per read):

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

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

---

# Example

Cache block: 2 bytes
Associativity: Directly Mapped
Cache size: 16 bytes

cache

| Data | Tag |
|---|---|
| | |
| | |
| | |

| Data | Tag |
|---|---|
| | |
| | |
| | |

memory

| | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|---|---|---|---|---|---|---|---|---|
| 000 | 3 | 5 | 7 | 9 | 11 | 13 | 15 | 17 |
| 001 | 19 | 21 | 23 | 25 | 27 | 29 | 31 | 33 |
| 010 | 35 | 37 | 39 | 41 | 43 | 45 | 47 | 49 |
| 011 | 51 | 53 | 55 | 57 | 59 | 61 | 63 | 65 |
| 100 | 67 | 69 | 71 | 73 | 75 | 77 | 79 | 81 |
| 101 | 83 | 85 | 87 | 89 | 91 | 93 | 95 | 97 |
| 110 | 99 | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

64 bytes

16

| 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|
|   |   |   |   |   |   |

Tag　　　　　　Block index　　Block offset

## Example

for(I = 0; I < 100; i++)  {
   … = a[i];
}

P

Compulsory miss
Capacity miss

c a c h e

| Data | Tag | Valid | Data | Tag | Valid |
|------|-----|-------|------|-----|-------|
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |

Load the following

010000
010001
010010
010011
010100
010101
010110
010111
011000
011001
011010
011011
011100
011101
011110
011111
10000

m e m o r y

|     | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 000 | 3   | 5   | 7   | 9   | 11  | 13  | 15  | 17  |
| 001 | 19  | 21  | 23  | 25  | 27  | 29  | 31  | 33  |
| 010 | 35  | 37  | 39  | 41  | 43  | 45  | 47  | 49  |
| 011 | 51  | 53  | 55  | 57  | 59  | 61  | 63  | 65  |
| 100 | 67  | 69  | 71  | 73  | 75  | 77  | 79  | 81  |
| 101 | 83  | 85  | 87  | 89  | 91  | 93  | 95  | 97  |
| 110 | 99  | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

## Example

### Slide 19

```
for(I = 0; I < 100; i++)  {
    … = a[i];
}
```

Load the following

Compulsory miss
Capacity miss

memory

|  | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|---|---|---|---|---|---|---|---|---|
| 000 | 3 | 5 | 7 | 9 | 11 | 13 | 15 | 17 |
| 001 | 19 | 21 | 23 | 25 | 27 | 29 | 31 | 33 |
| 010 | 35 | 37 | 39 | 41 | 43 | 45 | 47 | 49 |
| 011 | 51 | 53 | 55 | 57 | 59 | 61 | 63 | 65 |
| 100 | 67 | 69 | 71 | 73 | 75 | 77 | 79 | 81 |
| 101 | 83 | 85 | 87 | 89 | 91 | 93 | 95 | 97 |
| 110 | 99 | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

010000
010001
010010
010011
010100
010101
010110
010111
011000
011001
011010
011011
011100
011101
011110
011111
10000

cache

| Data | Tag | Valid | Data | Tag | Valid |
|---|---|---|---|---|---|
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |

19

### Slide 20

## Example

```
for(I = 0; I < 100; i++)  {
    … = a[i];
    … = b[i];
}
```

P

Load the following

Conflict miss

cache

| Data | Tag | Valid | Data | Tag | Valid |
|---|---|---|---|---|---|
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |

010000
110000
010001
110001
010010
110010

memory

|  | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|---|---|---|---|---|---|---|---|---|
| 000 | 3 | 5 | 7 | 9 | 11 | 13 | 15 | 17 |
| 001 | 19 | 21 | 23 | 25 | 27 | 29 | 31 | 33 |
| 010 | 35 | 37 | 39 | 41 | 43 | 45 | 47 | 49 |
| 011 | 51 | 53 | 55 | 57 | 59 | 61 | 63 | 65 |
| 100 | 67 | 69 | 71 | 73 | 75 | 77 | 79 | 81 |
| 101 | 83 | 85 | 87 | 89 | 91 | 93 | 95 | 97 |
| 110 | 99 | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

20

# Example

```
for(I = 0; I < 100; i++)  {
    … = a[i];
    … = b[i];
}
```

Conflict miss

010000
110000
010001
110001
010010
110010

memory

|     | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 000 | 3   | 5   | 7   | 9   | 11  | 13  | 15  | 17  |
| 001 | 19  | 21  | 23  | 25  | 27  | 29  | 31  | 33  |
| 010 | 35  | 37  | 39  | 41  | 43  | 45  | 47  | 49  |
| 011 | 51  | 53  | 55  | 57  | 59  | 61  | 63  | 65  |
| 100 | 67  | 69  | 71  | 73  | 75  | 77  | 79  | 81  |
| 101 | 83  | 85  | 87  | 89  | 91  | 93  | 95  | 97  |
| 110 | 99  | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

cache

| Data | Tag | Valid | Data | Tag | Valid |
|------|-----|-------|------|-----|-------|
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |

21

---

# Example

Load the following

P

```
for(I = 0; I < 100; i++)  {
    … = a[i];
    … = b[i];
    … = c[i];
}
```

Conflict miss

010000
110000
111000
010001
110001
111001
010010
110010
111010

cache

| Data | Tag | Valid | Data | Tag | Valid |
|------|-----|-------|------|-----|-------|
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |

memory

|     | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 000 | 3   | 5   | 7   | 9   | 11  | 13  | 15  | 17  |
| 001 | 19  | 21  | 23  | 25  | 27  | 29  | 31  | 33  |
| 010 | 35  | 37  | 39  | 41  | 43  | 45  | 47  | 49  |
| 011 | 51  | 53  | 55  | 57  | 59  | 61  | 63  | 65  |
| 100 | 67  | 69  | 71  | 73  | 75  | 77  | 79  | 81  |
| 101 | 83  | 85  | 87  | 89  | 91  | 93  | 95  | 97  |
| 110 | 99  | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

22

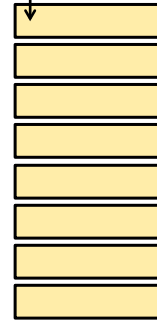# A Higher Level Example

```
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; i < 16; i++)
        for (i = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```
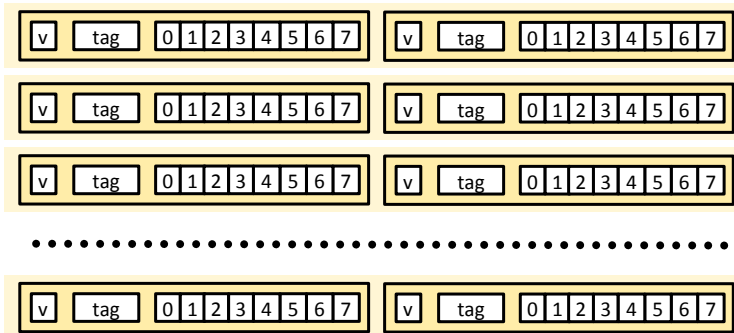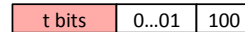
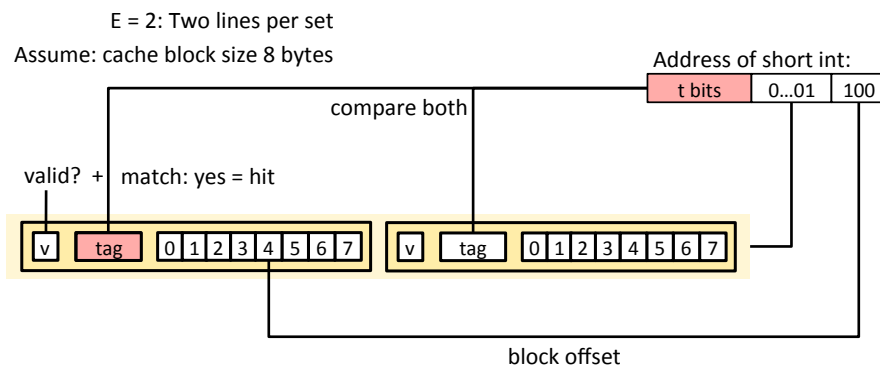assume: cold (empty) cache,
a[0][0] goes here

32 B = 4 doubles

---

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

E = 2: Two lines per set

Assume: cache block size 8 bytes

Address of short int:

| t bits | 0...01 | 100 |

find set

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

• • • • • • • • • • • • • • • • • • • • • • • • • • • • • • • • • • •

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

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

E = 2: Two lines per set

Assume: cache block size 8 bytes

Address of short 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

---

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

E = 2: Two lines per set

Assume: cache block size 8 bytes

Address of short 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

short int (2 Bytes) is here

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 | $[00\underline{0}0_2]$, | miss |
| 1 | $[00\underline{0}1_2]$, | hit |
| 7 | $[01\underline{1}1_2]$, | miss |
| 8 | $[10\underline{0}0_2]$, | miss |
| 0 | $[00\underline{0}0_2]$ | hit |

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

---

## Example



Cache block: 2 bytes
Associativity: 2
Cache size: 16 bytes

64 bytes

|     | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 000 | 3   | 5   | 7   | 9   | 11  | 13  | 15  | 17  |
| 001 | 19  | 21  | 23  | 25  | 27  | 29  | 31  | 33  |
| 010 | 35  | 37  | 39  | 41  | 43  | 45  | 47  | 49  |
| 011 | 51  | 53  | 55  | 57  | 59  | 61  | 63  | 65  |
| 100 | 67  | 69  | 71  | 73  | 75  | 77  | 79  | 81  |
| 101 | 83  | 85  | 87  | 89  | 91  | 93  | 95  | 97  |
| 110 | 99  | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

28

| 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|
|   |   |   |   |   |   |

Tag        Block index    Block offset

---

# Example

```
for(I = 0; I < 100; i++)  {
   … = a[i];
}
```

P

Load the following

Compulsory miss
Capacity miss

cache

| Data | Tag | Valid | Data | Tag | Valid |
|------|-----|-------|------|-----|-------|
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |

memory

|     | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 000 | 3   | 5   | 7   | 9   | 11  | 13  | 15  | 17  |
| 001 | 19  | 21  | 23  | 25  | 27  | 29  | 31  | 33  |
| 010 | 35  | 37  | 39  | 41  | 43  | 45  | 47  | 49  |
| 011 | 51  | 53  | 55  | 57  | 59  | 61  | 63  | 65  |
| 100 | 67  | 69  | 71  | 73  | 75  | 77  | 79  | 81  |
| 101 | 83  | 85  | 87  | 89  | 91  | 93  | 95  | 97  |
| 110 | 99  | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

010000
010001
010010
010011
010100
010101
010110
010111
011000
011001
011010
011011
011100
011101
011110
011111
10000

# Example

for(I = 0; I < 100; i++)  {
    ... = a[i];
    ... = b[i];
}

Conflict miss

P

c a c h e

| Data | Tag | Valid | Data | Tag | Valid |
|------|-----|-------|------|-----|-------|
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |

Load the following

010000
110000
010001
110001
010010
110010

m e m o r y

|     | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 000 | 3   | 5   | 7   | 9   | 11  | 13  | 15  | 17  |
| 001 | 19  | 21  | 23  | 25  | 27  | 29  | 31  | 33  |
| 010 | 35  | 37  | 39  | 41  | 43  | 45  | 47  | 49  |
| 011 | 51  | 53  | 55  | 57  | 59  | 61  | 63  | 65  |
| 100 | 67  | 69  | 71  | 73  | 75  | 77  | 79  | 81  |
| 101 | 83  | 85  | 87  | 89  | 91  | 93  | 95  | 97  |
| 110 | 99  | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

31

---

# Example

for(I = 0; I < 100; i++)  {
    ... = a[i];
    ... = b[i];
    ... = c[i];
}

Conflict miss

P

c a c h e

| Data | Tag | Valid | Data | Tag | Valid |
|------|-----|-------|------|-----|-------|
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |
|      |     |       |      |     |       |

Load the following

010000
110000
111000
010001
110001
111001
010010
110010
111010

m e m o r y

|     | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 000 | 3   | 5   | 7   | 9   | 11  | 13  | 15  | 17  |
| 001 | 19  | 21  | 23  | 25  | 27  | 29  | 31  | 33  |
| 010 | 35  | 37  | 39  | 41  | 43  | 45  | 47  | 49  |
| 011 | 51  | 53  | 55  | 57  | 59  | 61  | 63  | 65  |
| 100 | 67  | 69  | 71  | 73  | 75  | 77  | 79  | 81  |
| 101 | 83  | 85  | 87  | 89  | 91  | 93  | 95  | 97  |
| 110 | 99  | 101 | 103 | 105 | 107 | 109 | 111 | 113 |
| 111 | 115 | 117 | 119 | 121 | 123 | 125 | 127 | 129 |

32

16

## 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_col(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,
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)
    - 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)
    - 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**

# Intel Core i7 Cache Hierarchy

Processor package

Core 0

Regs

L1 d-cache | L1 i-cache

L2 unified cache

...

Core 3

Regs

L1 d-cache | L1 i-cache

L2 unified cache

L3 unified cache
(shared by all cores)

Main memory

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.

---

# Cache Performance

18

# Cache Performance Metrics

## Miss Rate

- Fraction of memory references not found in cache (misses/references)
- Typical numbers:
  - 3-10% for L1; can be quite small 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 clock cycle for L1; 3-8 clock cycles for L2

## Miss Penalty

- Additional time required because of a miss
  - Typically 25-100 cycles for accessing the main memory

---

# Average Memory Access Time

AMAT = Average Memory Access Time

$$AMAT = HitTime$$
$$+ MissRate \times MissPenalty$$

- Repeated references to variables are good (temporal locality)
- Stride-1 reference patterns are good (spatial locality)

## 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% hits:  1 cycle + 0.03 * 100 cycles = **4 cycles**
    99% hits:  1 cycle + 0.01 * 100 cycles = **2 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.

# Writing Cache Friendly Code

16-bit address space



```
for (i = 0; i < 997; i++)
    for (j = 0; j < 997; j++)
        sum += a[i][j];
return sum;
```

Direct-mapped, 4-byte words, 4-word cache blocks, 1024 bytes in size

Cache

---

# Writing Cache Friendly Code

0x0000
0x0004
0x0008
0x000C
0x0010
0x0014
0x0018
0x001C

0x0000
0x0F94
0x1F28
0x2EBC
. . .
0x0004
0x0F98
0x1F2C

```
for (i = 0; i < 997; i++)
    for (j = 0; j < 997; j++)
        sum += a[i][j];
return sum;
```

```
for (j = 0; j < 997; j++)
    for (i = 0; i < 997; i++)
        sum += a[i][j];
return sum;
```

**Miss rate = 1/4 = 25%**

**Miss rate = 100%**

21

## How to Improve Cache Performance?

Hit Time = 1 Cycle                    Miss Penalty = 100 Cycle

*Average memory access time =*

$$HitTime + MissRate \times MissPenalty$$

1. Reduce the miss rate,

2. Reduce the miss penalty, or

3. Reduce the time to hit in the cache.

What if we have a bigger cache?

---

## Another Example: Matrix Multiplication

Description:

- Multiply N x N matrices
- O(N3) total operations

Accesses

- N reads per source element
- N values summed per destination

```
/* 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;
  }
}
```

*Variable sum held in register*

## Miss Rate Analysis for Matrix Multiply

Assume:

- Line size = 32B (big enough for 4 8-byte words)
- Matrix dimension (N) is very large
- Cache is not even big enough to hold multiple rows

Analysis Method:

- Look at access pattern of inner loop

---

## 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:



A — Row-wise
B — Column-wise
C — Fixed

- <u>Misses per Inner Loop Iteration:</u>

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

## 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:



| | | |
|---|---|---|
| A | B | C |
| Row-wise | Column-wise | Fixed |

Misses per Inner Loop Iteration:

| A | B | C |
|---|---|---|
| 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:



| | | |
|---|---|---|
| A | B | C |
| Fixed | Row-wise | Row-wise |

• Misses per Inner Loop Iteration:

| A | B | C |
|---|---|---|
| 0.0 | 0.25 | 0.25 |

# 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:



|  (i,k)  |  (k,*)  |  (i,*)  |
|  A  |  B  |  C  |
|  Fixed  |  Row-wise  |  Row-wise  |

- <u>Misses per Inner Loop Iteration:</u>

| <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:



|  (*,k)  |  (k,j)  |  (*,j)  |
|  A  |  B  |  C  |
|  Column-wise  |  Fixed  |  Column-wise  |

- <u>Misses per Inner Loop Iteration:</u>

| <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:

```
      (*,k)                    (*,j)
              (k,j)
   A          B            C
```

| Column-wise | Fixed | Column-wise |

- Misses per Inner Loop Iteration:

| A | B | C |
|---|---|---|
| 1.0 | 0.0 | 1.0 |

---

# Summary of Matrix Multiplication

Cache hit = 1 cycle;   Cache miss penalty = 100 cycle

| **ijk (& jik):** | **kij (& ikj):** | **jki (& kji):** |
|---|---|---|
| • 2 loads, 0 stores | • 2 loads, 1 store | • 2 loads, 1 store |
| • misses/iter = **1.25** | • misses/iter = **0.5** | • misses/iter = **2.0** |

```
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;
  }
}
```

```
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;
  }
}
```

## 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];
}
```
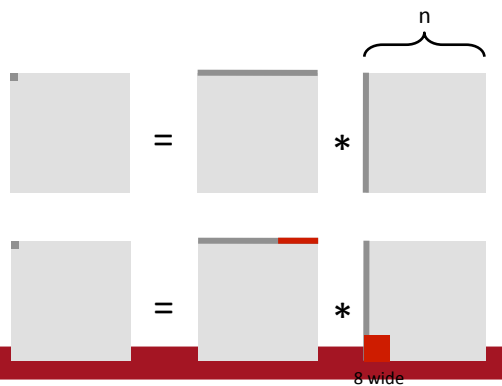
c = a * b

i

---

## Cache Miss Analysis

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

- First iteration:
  - $n/8 + n = 9n/8$ misses

  - Afterwards in cache: (schematic)

n

= *

= *

8 wide

# Cache Miss Analysis

- 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

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



n

= * 

8 wide

# 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];
}
```
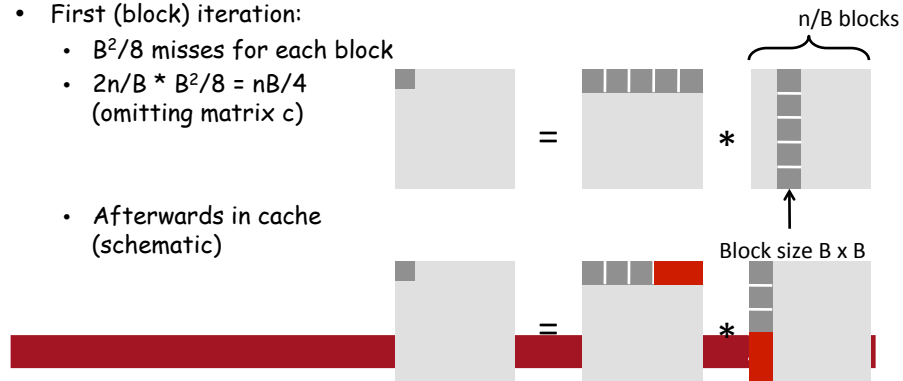
j1



c   =   a   *   b   +   c

i1

Block size B x B

# Cache Miss Analysis

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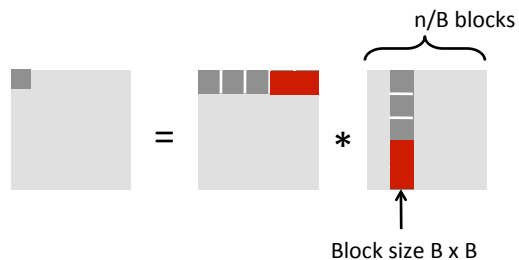  - Three blocks ▪ fit into cache: $3B^2 < C$

- First (block) iteration:
  - $B^2/8$ misses for each block
  - $2n/B * B^2/8 = nB/4$ (omitting matrix c)

  - Afterwards in cache (schematic)

n/B blocks

=       *

Block size B x B

=       *

---

# 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
  - $2n/B * B^2/8 = nB/4$

n/B blocks

=       *

- Total misses:
  - $nB/4 * (n/B)^2 = n^3/(4B)$

Block size B x B

## Summary

- No blocking: $(9/8) * n^3$
- Blocking: $1/(4B) * n^3$

- 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

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)