# **Caches IV**

CSE 351 Summer 2024

**Instructor:** Ellis Haker

#### **Teaching Assistants:**

Naama Amiel Micah Chang Shananda Dokka Nikolas McNamee Jiawei Huang The cache when you ask for something that was just evicted:



# Administrivia

- Nothing due tonight!
- Sunday, 7/28
  - Lab3 due (11:59pm)
- Monday, 7/29
  - HW15-16 due (11:59pm)
  - Quiz 2 due (11:59pm)

I made some mistakes in the last lecture (I was... very tired)

• Ed slides have been fixed, and we'll go over the last problem again today

#### **Code Analysis**

• Assuming cache starts **cold** (i.e. all blocks invalid), and sum, i, and j are all stored in registers, calculate the **miss rate**.

• 
$$m = 10$$
 bits,  $C = 64B$ ,  $K = 8B$ ,  $E = 2$   
#define SIZE 8  
short ar[SIZE][SIZE], sum = 0; // &ar=0x200 = cr [0]  
for (int i = 0; i < SIZE; i++) {  
for (int j = 0; j < SIZE; j++)  
sum += ar[j][i];  
}  
code visite code element of arm, in columns

#### **Code Analysis: relevant values**

- *m* = 10 bits, *C* = 64B, *K* = 8B, *E* = 2
   *k* = 3, *s* = 2
- 8B blocks = 4 shorts per block
- Starting address = 0b10000 00 000
  - Block stored in set 0, tag = 0b10000 = 0x10

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```



- Access ar [0] [0]
  - **0b10000 00 000** 
    - Miss! compations
  - Load block into set 0 with tag 0x10

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```

Set	Tag	Data	Tag	Data
0	10	a[0][0] a[0][3]		
1				
2				
3				

- Access ar [1][0]
  - 0 0b10000 10 000
    - Miss! Comput sory
  - Load block into set 2 with tag 0x10

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```

$$i = 0, j = 1$$
 Misses: 2  
Hits: 0

Set	Tag	Data	Tag	Data
0	10	a[0][0] … a[0][3]		
1				
2	10	a[1][0] … a[1][3]		
3				

- Access ar [2] [0]
  - 0 0b10001 00 000
    - Miss! Comphilsory
  - Load block into set 0 with tag 0x11
    - Can store both blocks in set 0 because of associativity

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```

$$i = 0, j = 2$$
 Misses: 3  
Hits: 0

Set	Tag	Data	Tag	Data 🖌
0	10	a[0][0] … a[0][3]	11	a[2][0] … a[2][3]
1				
2	10	a[1][0] … a[1][3]		
3				

- Access ar [3] [0]
  - **0b10001 10 000** 
    - Miss! Compulsing

• Load block into set 1 with tag 0x11

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```

$$i = 0, j = 3$$
 Misses: 4  
Hits: 0

Set	Tag	Data	Tag	Data
0	10	a[0][0] … a[0][3]	11	a[2][0] … a[2][3]
1				
2	10	a[1][0] … a[1][3]	11	a[3][0] … a[3][3]
3				

#### 9

- Access ar [4] [0]
  - **0b10010 00 000** 
    - Miss! Compelsiss
  - Load block into set 0 with tag 0x12
    - Evicts least recently used block  $(\sim [\circ] [\circ] \sim [\circ] [3]$

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```

$$i = 0, j = 4$$
 Misses: 5  
Hits: 0

Set	Tag	Data	Tag	Data
0	12	a[4][0] a[4][3]	11	a[2][0] … a[2][3]
) 1				
2	10	a[1][0] a[1][3]	11	a[3][0] … a[3][3]
3				

- Same as step 5
  - Accesses to a[5][0], a[6][0], and
     a[7][0] will kick out the old blocks in the cache
     3x completes miss
- So for i = 0:
  - $\circ$  8 accesses total (j = 0...7), 8 misses

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```

Set	Tag	Data	Tag	Data
0	12	a[4][0] a[4][3]	13	a[6][0] … a[6][3]
1				
2	12	a[5][0] … a[5][3]	13	a[7][0] … a[7][3]
3				

- Access ar[0][1]
  - - Same block that we loaded in in step
       1, but it got evicted in step 5!
    - Miss! Conflict
  - Load block back into set 0

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```

$$i = 1, j = 0$$
 Misses: 9  
Hits: 0

Set	Tag	Data	Tag	Data
0	10	a[0][0] a[0][3]	13	a[6][0] … a[6][3]
1				
2	12	a[5][0] … a[5][3]	13	a[7][0] … a[7][3]
3				

- All future accesses will continue to follow this pattern
  - Each block is loaded in, then kicked out of the cache before it's accessed again
- Miss rate 100%!
- How can we fix this?

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[j][i];
}</pre>
```

# **Improving Cache Performance**

- Reduce stride
  - I.e. access data that's closer together

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[i][j];
}</pre>
```

Everything else is the some

Access each element in rows

#### Improving Cache Performance Example: step 1 Misses: 1 Hits: 3

- First 4 accesses:
  - ar [0] [0]: miss, load block into the cache
  - ar[0][1]: hit!
  - ar[0][2]: hit!
  - ar[0][3]: hit!

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[i][j];
}</pre>
```

Set	Tag	Data	Tag	Data
0	10	ar[0][0] … ar[0][3]		
1				
2				
3				

#### Improving Cache Performance Example: step 2 Misses: 2 Hits: 6

- Next 4 accesses:
  - ar [0] [4]: miss, load block into the cache
  - ar[0][5]: hit!
  - ar[0][6]: hit!
  - ar[0][7]: hit!

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[i][j];
}</pre>
```

Set	Tag	Data	Tag	Data
0	10	ar[0][0] … ar[0][3]		
1	10	ar[0][4] ar[0][7]		
2				
3				

#### Improving Cache Performance Example: step 3+

- All accesses follow this pattern
  - Because we use a whole block before moving on, we will miss 1 out of every 4 accesses
- Miss rate: 25%

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++) {
    for (int j = 0; j < SIZE; j++)
        sum += ar[i][j];
}</pre>
```

### Caches

- Cache basics
- Principle of locality
- Memory hierarchies
- Cache organization
  - Direct-mapped (sets; index + tag)
  - Associativity (ways)
  - Replacement policy
  - Handling writes
- Program optimizations that consider caches

# **Write-Hit Policies**

i.e. memory or no ther 1 cache level

What to do if the data is already in the cache?

- Write-through: immediately write to the next level
- Write-back: don't write to next level until we have to
  - Keep track of dirty bit for each block
  - $\circ$   $\,$  On eviction, if dirty bit is set, write contents back to memory



### **Write-Miss Policies**

- What to do if the block we want to write to is not in the cache?
- No-write-allocate ("write around"): don't load into the cache, just write to the next level
- Write-allocate ("fetch on write") load data into the cache before writing



### **Ex: Write-Back, Write-Allocate**

- Single-block mini cache
  - Tag includes the entire block number
  - Not a realistic example





#### Cache



# Ex: Write-Back, Write-Allocate (pt 2)

**1.** mov \$0×B0BA, (F)

Write miss

Not valid x86. Assume we mean an address within block F. Write-back: defer write to next level until line is evicted Write-allocate: on a miss, bring the data into cache



#### Ex: Write-Back, Write-Allocate (pt 3)

mov \$0xB0BA, (F)
 Write miss
 Secouse of write-allocate policy
 Bring F into cache
 Cache
 V D Tag Data
 I 0 F 0xDEAD

Write-back: defer write to next level until line is evicted Write-allocate: on a miss, bring the data into cache





#### Ex: Write-Back, Write-Allocate (pt 5)

2. mov \$0xF00D, (F)

Write hit

Write-back: defer write to next level until line is evicted Write-allocate: on a miss, bring the data into cache



#### Cache



#### Ex: Write-Back, Write-Allocate (pt 6)



Write-back: defer write to next level until line is evicted Write-allocate: on a miss, bring the data into cache



Data

0xF00D

#### Ex: Write-Back, Write-Allocate (pt 7)

**3.** mov (G), %ax

Read miss

Write-back: defer write to next level until line is evicted Write-allocate: on a miss, bring the data into cache



#### Cache

#### Ex: Write-Back, Write-Allocate (pt 8)

3. mov (G), %ax

Read miss

Write-back: defer write to next level until line is evicted Write-allocate: on a miss, bring the data into cache



### Ex: Write-Back, Write-Allocate (pt 9)

3. mov (G), %ax

Read miss

Write-back: defer write to next level until line is evicted Write-allocate: on a miss, bring the data into cache



- Write F back to memory since it is dirty
- 2. Bring G into cache

#### Cache



# **Common Policies**

- Write-back + Write-allocate
  - (most common)
- Write-through + No-write-allocate
  - When would this be used?

Ex: L2 code may be write - through so that other corres con access Lote through share 2 L3

#### Processor package



# **Polling Question**

- 1. Which of the following cache statements is FALSE?
  - A) A write-through cache will always match data with the memory hierarchy level below it.

B)We can reduce compulsory misses by decreasing our block size.

C) A write-back cache will save time for code with good temporal locality on writes.

D) We can reduce conflict misses by increasing associativity

```
actually the apposite!
I block size = need to book
in more blocks
```

con hold more blocks per set = fewer evictions

Lefinition of write - through

Louly need to write to memory once, instead of for each write

# **Cache Simulator**

- Want to play around with cache parameters and policies? Check out our cache simulator!
  - <u>https://courses.cs.washington.edu/courses/cse351/cachesim/</u>
- Way to use:
  - Take advantage of "explain mode" and navigable history to test your own hypotheses and answer your own questions
  - Self-guided Cache Sim Demo in section
  - Will be used in HW17 Lab 4 Preparation

### Caches

- Cache basics
- Principle of locality
- Memory hierarchies
- Cache organization
  - Direct-mapped (sets; index + tag)
  - Associativity (ways)
  - Replacement policy
  - o Handling writes
- Program optimizations that consider caches

# **Optimizations for the Memory Hierarchy**

- Write code that has locality
  - Spatial: access data contiguously
  - Temporal: make sure access to the same data is not too far apart in time
- How can you achieve locality?
  - Adjust memory accesses in code to improve miss rate (MR)
    - Requires knowledge of **both** how caches work as well as your system's parameters
  - Proper choice of algorithm
  - Loop transformations

#### **Example: Matrix Multiplication**

 $(C_{j}) = AC_{j} =$ 



# **Matrices in Memory**

- How do cache blocks fit into this scheme?
  - Row-major in memory
  - Column of matrix is spread across multiple cache blocks





# **Cache Miss Analysis (Naive)**

- Example parameters:
  - A and B are square  $(n \times n)$  matrices of doubles - Joing by rows gets ophinal miss rule Ο
  - Cache block size K = 64B $\cap$ 
    - 8 elements per block
  - Cache size *C* << *n* (much smaller) Ο
- For each element of C
- n/8 misses per row of A 0
- → *n* misses per column *B* 
  - 9n/8 misses per iteration =  $s + \frac{1}{5}$ Ο
  - $n^2$  elements in total
    - 9n<sup>3</sup>/8 misses =  $\frac{9}{5}$  ...<sup>2</sup>

Ignoring misses in matrix C

#### Linear Algebra To The Rescue

- Can get the same result of a matrix multiplication by splitting the matrices into smaller submatrices ("blocks")
- Example: multiply 2 4×4 matrices

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{pmatrix} = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}, \text{ with } B \text{ defined similarly}$$

$$Oon't \text{ week to know the behavis,}$$

$$AB = \begin{pmatrix} (A_{11}B_{11} + A_{12}B_{21}) & (A_{11}B_{12} + A_{12}B_{22}) \\ (A_{21}B_{11} + A_{22}B_{21}) & (A_{21}B_{12} + A_{22}B_{22}) \end{pmatrix}$$

$$Sholles \text{ constrained}$$

#### Linear Algebra To The Rescue (pt 2)



A <sub>11</sub>	A <sub>12</sub>	A <sub>13</sub>	A <sub>14</sub>
A <sub>21</sub>	A <sub>22</sub>	A <sub>23</sub>	A <sub>24</sub>
<b>A</b> <sub>31</sub>	A <sub>32</sub>	A <sub>33</sub>	A <sub>34</sub>
A <sub>41</sub>	A <sub>42</sub>	A <sub>43</sub>	A <sub>144</sub>

B <sub>11</sub>	B <sub>12</sub>	B <sub>13</sub>	B <sub>14</sub>
B <sub>21</sub>	B <sub>22</sub>	B <sub>23</sub>	B <sub>24</sub>
B <sub>32</sub>	B <sub>32</sub>	B <sub>33</sub>	B <sub>34</sub>
B <sub>41</sub>	B <sub>42</sub>	B <sub>43</sub>	B <sub>44</sub>

- Split up each matrix into smaller blocks
  - Choose a block size that will fit in the cache! 0
  - This technique is called cache blocking Ο
- Similar alg as earlier, but now each of yuese is a mini- ment Perform multiplication block-by-block (rather than row-by-row)  $\bullet$

$$C_{22} = A_{21}B_{12} + A_{22}B_{22} + A_{23}B_{32} + A_{24}B_{42} = \sum_{k=1}^{n} A_{2k} * Bk_2$$

# **Blocked Matrix Multiply**

• *r* = length of a block (assume *r* divides *n* evenly)

???

# **Cache Miss Analysis (Blocked)**

- Example parameters:
  - Cache block size K = 64B
    - 8 elements per block
  - Cache size *C* << *n* (much smaller)
  - Cache can hold 3 blocks ( $r \times r$ ):  $3r^2 < C$
- For each block of C
- "> For Ao r2/8 misses per block optime 1 miles rele
  - 8 elements per cache block
     6 is in the coch
  - $^{3}2n/r$  blocks accessed (row and col) ->  $2n/r * r^{2}/8 = nr/4$  misses
  - $(n/r)^2$  blocks in total
    - $nr/4 * (n/r)^2 = n^3/4r$  misses-

Compare to 9n<sup>3</sup>/8 from before Ignoring misses in matrix *C* 



# **Cache Images**

- Recap:
  - Blocks are contiguous chunks of memory
  - Contiguous blocks in memory map to 00
     contiguous cache sets 00
- Conclusion:
  - Let *B* = the number of blocks the cache can hold (*C* / *K*)
  - A contiguous, aligned chunk of B blocks (C bytes) can all fit into the cache at once!





# Cache Images (pt 2)

- Cache image: an aligned chunk of memory the same size as the cache
  - Guaranteed to all fit in the cache at once
    - Locality!
  - The offset of an address within its cache image tells you where in the cache it will map to
    - Two addresses with the same offset in different images will map to the same location in the cache



#### Memory

# Summary

- Programmer can optimize for cache performance
  - How data structures are organized
  - $\circ$   $\,$  How data is accessed
    - Cache blocking for 2D arrays
- Getting absolute optimum performance is very platform specific
  - Depends on cache size, block size, associativity, etc.
  - Generic guidelines:
    - Keep working set reasonably small (temporal locality)
    - Use small strides (spatial locality)
    - Focus on inner loop code (for nested loops)
- Cache images can help you quickly figure out where data will go in the cache