# 20.  Caches: Set Associative

**EECS 370 – Introduction to Computer Organization – Fall 2020**

Satish Narayanasamy

**EECS Department**
**University of Michigan in Ann Arbor, USA**

# Announcements

Upcoming deadlines:

Assignment Project Exam Help

HW4          due Nov 10th

Project 3       due Nov. 12th

https://powcoder.com

Add WeChat powcoder

Grading policy: Best of two

# Fully-associative caches

**Memory**

**tag   data**

**A block can go
to any location**

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Address:

| tag | block offset |
|---|---|
| **3 bits** | **1 bit** |

| | |
|---|---|
| 0 | 100 |
| 1 | 110 |
| 2 | 120 |
| 3 | 130 |
| 4 | 140 |
| 5 | 150 |
| 6 | 160 |
| 7 | 170 |
| 8 | 180 |
| 9 | 190 |
| 10 | 200 |
| 11 | 210 |
| 12 | 220 |
| 13 | 230 |
| 14 | 240 |
| 15 | 250 |

3

# Fully-associative cache: Placement & Access

Address

| Tag | Block offset |

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Tag array

Data array

v    tag

Compare all tags

MUX ← Block offset

Hit?

Data

4

# Direct-mapped caches

## Cache

**Memory**

A block can go to **one** location

**tag   data**

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |

| | |
|---|---|
| 0 | 78 |
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 71 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

**Address:**

| tag | line index | block offset |
|---|---|---|
| **1 bit** | **2 bits** | **1 bit** |

5

# Direct-mapped cache: Placement & Access

Address

| Tag | Line Index | Block offset |
|-----|-----------|--------------|

Tag array

Data array

| v | tag |
|---|-----|

=?

Hit?

MUX — Block offset

Data

# This lecture

**Set Associative Caches**

Assignment Project Exam Help

**Idea**

**Illustration** https://powcoder.com

**3C problem**

Add WeChat powcoder

# The middle ground...

Set associative caches

Partition memory into regions,like direct mapped but fewer partitions

Associate a region to a set of cache lines

Check tags for all lines in a set to determine a HIT

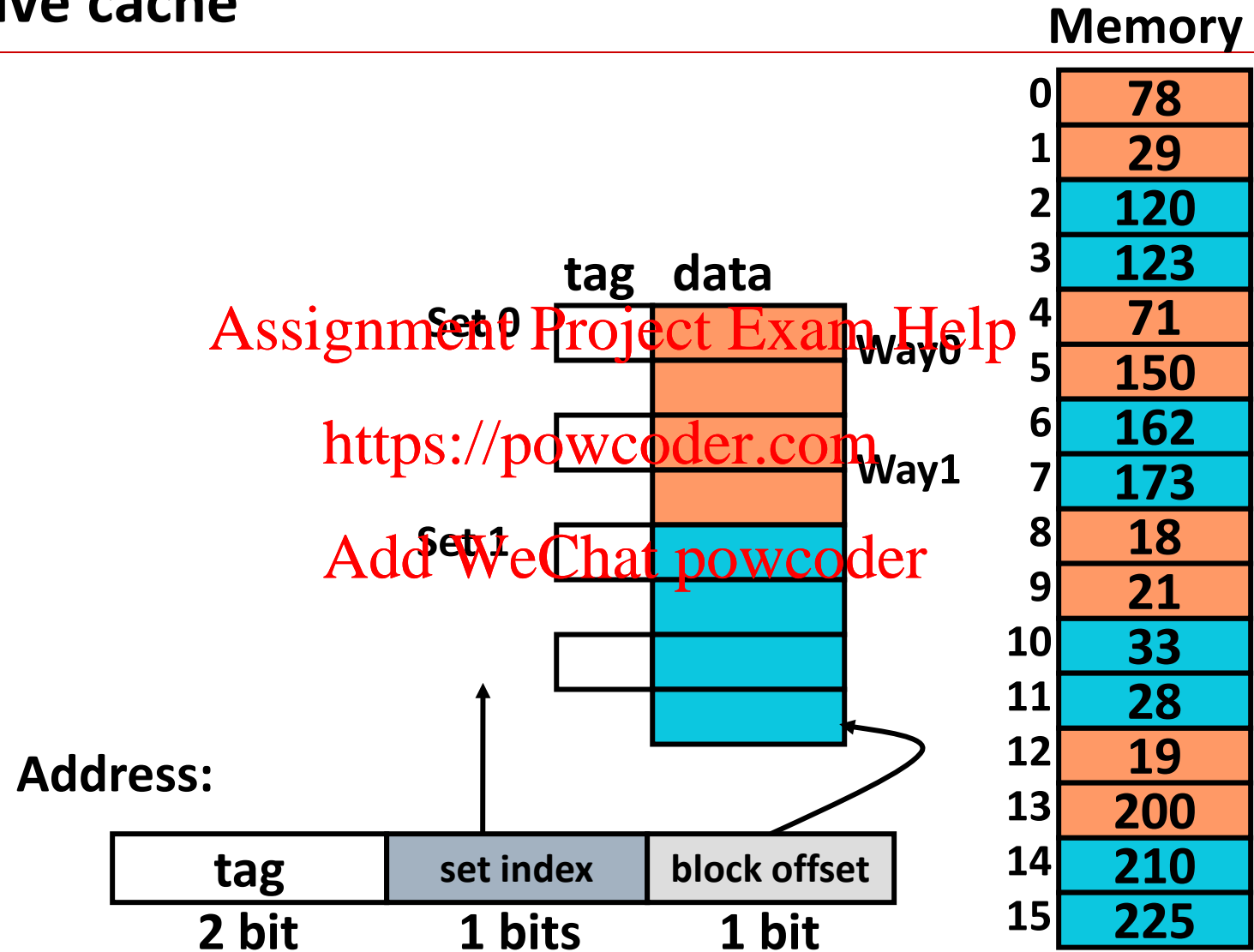Treat each line in a set like a small fully associative cache

LRU (or LRU-like) policy generally used

# Set-associative cache

**Memory**



tag    data

Set 0

Way0

Way1

Set 1

Address:

| tag | set index | block offset |
|-----|-----------|--------------|
| 2 bit | 1 bits | 1 bit |

| | |
|---|---|
| 0 | 78 |
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 71 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

# Set-associative cache: Placement & Access

# Cache Organization Comparison
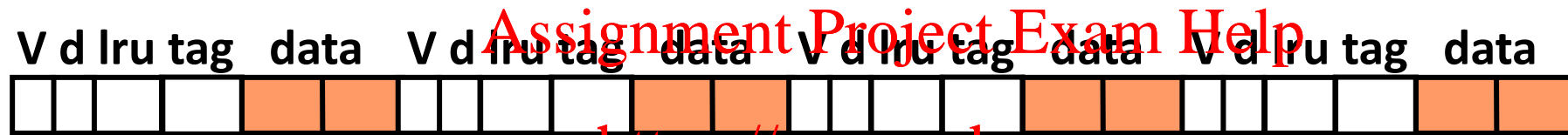
Cache size = 8 bytes (for all caches)
Block size  = 2 bytes
#blocks      = 4

Fully associative

# blocks per set = all blocks = 4 in this example;
so, also correct to view this cache as 4-way associative
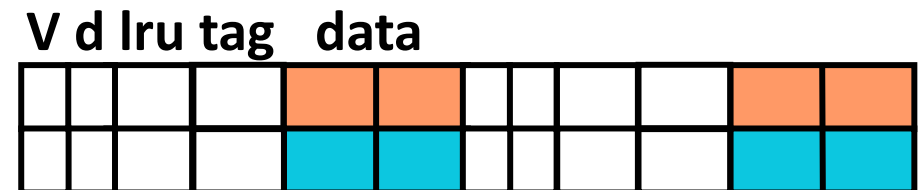
**V d lru tag   data   V d lru tag   data   V d lru tag   data   V d lru tag   data**



Direct mapped: (#blocks per set = 1)          2-way associative   (#blocks per set = 2)

**V d  tag   data**



**V d lru tag   data**



11

# Cache Organization: Equations

Block

     #blocks = cache size / block_size             #cache lines = #blocks

     block_offset_size = log2(#block_size)

Assignment Project Exam Help

Set

           #sets =  #lines / #ways    = #lines / (lines per set)

https://powcoder.com

Direct-mapped:      #sets =  #lines/1

2-way associative:    #sets =  #lines/2  Add WeChat powcoder

n-way associative:    #sets =  #lines/n

fully-associative:     #sets =   1           (all lines are in 1 set)


     set_index_size = log2(#sets)


Tag size = address size − set_index_size − block_offset_size

## Class Problem 1

For a 32-bit address and 16KB cache with 64-byte blocks, show the breakdown of the address for the following cache configuration:

Assignment Project Exam Help

**A) fully associative cache**                    **B) 4-way set associative cache**

https://powcoder.com

Add WeChat powcoder

**C) Direct-mapped cache**

# Class Problem 1 (Solution)

For a 32-bit address and 16KB cache with 64-byte blocks, show the breakdown of the address for the following cache configuration:

Assignment Project Exam Help

#lines = 16KB / 64 byte = 256

Block_offset_size = log(block_size) = log(64)= 6 bits

https://powcoder.com

**A) fully associative cache**

**B) 4-way set associative cache**

Add WeChat powcoder

Tag = 32 - 6 = 26 bits

#sets = #lines / ways = #lines/ 4 =  64

Set Index size = log(64) = 6 bits

**C) Direct-mapped cache**

Tag = address_size – set_index_size – block_offset_size

= 32 - 6 - 6

= 20 bits

#sets = #lines / ways = #lines/1 = 256

Set_index_size = log(#sets) = log(256) = 8 bits

Tag = 32 – 6 – 8 = 18 bits

## Set Associative Caches: Illustration

# Set-associative cache example (Write-back, write allocate)

| | Processor | Cache | Memory | Tag | Index | Block_offset |
|---|---|---|---|---|---|---|

**Processor**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

R0
R1
R2
R3

**Cache**

V d tag data

0
0
0
0

Misses:  0

Hits:  0

**Memory**

| | | Tag | Index | Block_offset |
|---|---|---|---|---|
| 0 | 78 | 00 | 0 | 0 |
| 1 | 29 | 00 | 0 | 1 |
| 2 | 120 | 00 | 1 | 0 |
| 3 | 123 | 00 | 1 | 1 |
| 4 | 71 | 01 | 0 | 0 |
| 5 | 150 | 01 | 0 | 1 |
| 6 | 162 | 01 | 1 | 0 |
| 7 | 173 | 01 | 1 | 1 |
| 8 | 18 | 10 | 0 | 0 |
| 9 | 21 | 10 | 0 | 1 |
| 10 | 33 | 10 | 1 | 0 |
| 11 | 28 | 10 | 1 | 1 |
| 12 | 19 | 11 | 0 | 0 |
| 13 | 200 | 11 | 0 | 1 |
| 14 | 210 | 11 | 1 | 0 |
| 15 | 225 | 11 | 1 | 1 |

16

# Set-associative cache (REF 1)

| | Processor | Cache | Memory | | Tag | Index | Block_offset |
|---|---|---|---|---|---|---|---|

**Processor**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

R0
R1
R2
R3

**Cache**

V d tag data

0
0
0
0

Misses:  0

Hits:  0

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

**Memory**

| | | Tag | Index | Block_offset |
|---|---|---|---|---|
| 0 | 78 | 00 | 0 | 0 |
| 1 | 29 | 00 | 0 | 1 |
| 2 | 120 | 00 | 1 | 0 |
| 3 | 123 | 00 | 1 | 1 |
| 4 | 71 | 01 | 0 | 0 |
| 5 | 150 | 01 | 0 | 1 |
| 6 | 162 | 01 | 1 | 0 |
| 7 | 173 | 01 | 1 | 1 |
| 8 | 18 | 10 | 0 | 0 |
| 9 | 21 | 10 | 0 | 1 |
| 10 | 33 | 10 | 1 | 0 |
| 11 | 28 | 10 | 1 | 1 |
| 12 | 19 | 11 | 0 | 0 |
| 13 | 200 | 11 | 0 | 1 |
| 14 | 210 | 11 | 1 | 0 |
| 15 | 225 | 11 | 1 | 1 |

17

# Set-associative cache (REF 1)

| | Processor | Cache | Memory | | Tag | Index | Block_offset |
|---|---|---|---|---|---|---|---|

**Processor**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

| R0 | |
|---|---|
| R1 | 29 |
| R2 | |
| R3 | |

**Cache**

V d tag data

| 1 | 0 | 0 | 78 |
|---|---|---|---|
| | | | 29 |

lru

| 0 | | | |
|---|---|---|---|

| 0 | | | |
|---|---|---|---|

| 0 | | | |
|---|---|---|---|

Misses: 1

Hits: 0

**Memory**

| 0 | 78 |
|---|---|
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 71 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

| Tag | Index | Block_offset |
|---|---|---|
| 00 | 0 | 0 |
| 00 | 0 | 1 |
| 00 | 1 | 0 |
| 00 | 1 | 1 |
| 01 | 0 | 0 |
| 01 | 0 | 1 |
| 01 | 1 | 0 |
| 01 | 1 | 1 |
| 10 | 0 | 0 |
| 10 | 0 | 1 |
| 10 | 1 | 0 |
| 10 | 1 | 1 |
| 11 | 0 | 0 |
| 11 | 0 | 1 |
| 11 | 1 | 0 |
| 11 | 1 | 1 |

18

# Set-associative cache (REF 2)

| Processor | Cache | Memory |
|---|---|---|

Assignment Project Exam Help

V d tag data

Ld R1 ← M[ 1 ]    `1 0 0  78`

https://powcoder.com

Ld R2 ← M[ 5 ]    `29`

St R2 → M[ 7 ]    `0`
lru

St R1 → M[ 4 ]

Ld R3 ← M[ 0 ]    Add WeChat powcoder

Ld R2 ← M[ 8 ]    `0`

`0`

| R0 | |
|---|---|
| R1 | 29 |
| R2 | |
| R3 | |

**Misses: 1**

**Hits:     0**

**Memory**

| | Tag | Index | Block_offset |
|---|---|---|---|
| 0 | 78 | 00 | 0 | 0 |
| 1 | 29 | 00 | 0 | 1 |
| 2 | 120 | 00 | 1 | 0 |
| 3 | 123 | 00 | 1 | 1 |
| 4 | 71 | 01 | 0 | 0 |
| 5 | 150 | 01 | 0 | 1 |
| 6 | 162 | 01 | 1 | 0 |
| 7 | 173 | 01 | 1 | 1 |
| 8 | 18 | 10 | 0 | 0 |
| 9 | 21 | 10 | 0 | 1 |
| 10 | 33 | 10 | 1 | 0 |
| 11 | 28 | 10 | 1 | 1 |
| 12 | 19 | 11 | 0 | 0 |
| 13 | 200 | 11 | 0 | 1 |
| 14 | 210 | 11 | 1 | 0 |
| 15 | 225 | 11 | 1 | 1 |

# Set-associative cache (REF 2)

| | Processor | Cache | Memory | Tag | Index | Block_offset |
|---|---|---|---|---|---|---|

**Processor**

Ld  R1 ← M[ 1 ]
Ld  R2 ← M[ 5 ]
St  R2 → M[ 7 ]
St  R1 → M[ 4 ]
Ld  R3 ← M[ 0 ]
Ld  R2 ← M[ 8 ]

| R0 | |
|---|---|
| R1 | 29 |
| R2 | 150 |
| R3 | |

**Cache**

V d tag    data

| lru | 1 | 0 | 0 | 78 |
|---|---|---|---|---|
| | | | | 29 |
| | 1 | 0 | 1 | 71 |
| | | | | 150 |
| | 0 | | | |
| | | | | |
| | 0 | | | |
| | | | | |

Misses:   2

Hits:      0

**Memory**

| 0 | 78 |
|---|---|
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 71 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

| Tag | Index | Block_offset |
|---|---|---|
| 00 | 0 | 0 |
| 00 | 0 | 1 |
| 00 | 1 | 0 |
| 00 | 1 | 1 |
| 01 | 0 | 0 |
| 01 | 0 | 1 |
| 01 | 1 | 0 |
| 01 | 1 | 1 |
| 10 | 0 | 0 |
| 10 | 0 | 1 |
| 10 | 1 | 0 |
| 10 | 1 | 1 |
| 11 | 0 | 0 |
| 11 | 0 | 1 |
| 11 | 1 | 0 |
| 11 | 1 | 1 |

20

# Set-associative cache (REF 3)



| | Processor | Cache | Memory | Tag | Index | Block_offset |

**Processor**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

| R0 | |
| R1 | 29 |
| R2 | 150 |
| R3 | |

**Cache**

V d tag    data

| lru | 1 | 0 | 0 | 78 |
| | | | | 29 |
| | 1 | 0 | 1 | 71 |
| | | | | 150 |
| 0 | | | |
| | | | |
| 0 | | | |
| | | | |

Misses:  2

Hits:    0

**Memory**

| 0 | 78 |
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 71 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

| Tag | Index | Block_offset |
|-----|-------|--------------|
| 00 | 0 | 0 |
| 00 | 0 | 1 |
| 00 | 1 | 0 |
| 00 | 1 | 1 |
| 01 | 0 | 0 |
| 01 | 0 | 1 |
| 01 | 1 | 0 |
| 01 | 1 | 1 |
| 10 | 0 | 0 |
| 10 | 0 | 1 |
| 10 | 1 | 0 |
| 10 | 1 | 1 |
| 11 | 0 | 0 |
| 11 | 0 | 1 |
| 11 | 1 | 0 |
| 11 | 1 | 1 |

21

# Set-associative cache (REF 3)

| | Processor | Cache | Memory | Tag | Index | Block_offset |
|---|---|---|---|---|---|---|

**Processor**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

| R0 | |
|---|---|
| R1 | 29 |
| R2 | 150 |
| R3 | |

**Cache**

V d tag data

| lru | 1 | 0 | 0 | 78 |
|---|---|---|---|---|
| | | | | 29 |
| | 1 | 0 | 1 | 71 |
| | | | | 150 |
| | 1 | 1 | 1 | 162 |
| | | | | 150 |
| lru | 0 | | | |
| | | | | |

Misses: 3

Hits: 0

**Memory**

| 0 | 78 | 00 | 0 | 0 |
| 1 | 29 | 00 | 0 | 1 |
| 2 | 120 | 00 | 1 | 0 |
| 3 | 123 | 00 | 1 | 1 |
| 4 | 71 | 01 | 0 | 0 |
| 5 | 150 | 01 | 0 | 1 |
| 6 | 162 | 01 | 1 | 0 |
| 7 | 173 | 01 | 1 | 1 |
| 8 | 18 | 10 | 0 | 0 |
| 9 | 21 | 10 | 0 | 1 |
| 10 | 33 | 10 | 1 | 0 |
| 11 | 28 | 10 | 1 | 1 |
| 12 | 19 | 11 | 0 | 0 |
| 13 | 200 | 11 | 0 | 1 |
| 14 | 210 | 11 | 1 | 0 |
| 15 | 225 | 11 | 1 | 1 |

22

# Set-associative cache (REF 4)

| Processor | Cache | Memory | Tag | Index | Block_offset |
|---|---|---|---|---|---|

**Processor**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
→ St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

| R0 | |
|---|---|
| R1 | 29 |
| R2 | 150 |
| R3 | |

**Cache**

V d tag data

lru | 1 | 0 | 0 | 78 |
| | | | 29 |
| 1 | 0 | 1 | 71 |
| | | | 150 |
| 1 | 1 | 1 | 162 |
| | | | 150 |
lru | 0 | | | |

Misses: 3

Hits: 0

**Memory**

| 0 | 78 |
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 71 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

| Tag | Index | Block_offset |
|---|---|---|
| 00 | 0 | 0 |
| 00 | 0 | 1 |
| 00 | 1 | 0 |
| 00 | 1 | 1 |
| 01 | 0 | 0 |
| 01 | 0 | 1 |
| 01 | 1 | 0 |
| 01 | 1 | 1 |
| 10 | 0 | 0 |
| 10 | 0 | 1 |
| 10 | 1 | 0 |
| 10 | 1 | 1 |
| 11 | 0 | 0 |
| 11 | 0 | 1 |
| 11 | 1 | 0 |
| 11 | 1 | 1 |

23

# Set-associative cache (REF 4)

| Processor | Cache | Memory |
|---|---|---|

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

**Processor instructions:**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
→ St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

**Cache** — V d tag data

lru | 1 | 0 | 0 | 78 |
| | | | 29 |
| 1 | 1 | 1 | 29 |
| | | | 150 |
| 1 | 1 | 1 | 162 |
| | | | 150 |
lru | 0 | | | |

**Registers:**

| R0 | |
|---|---|
| R1 | 29 |
| R2 | 150 |
| R3 | |

Misses:  3

Hits:  1

**Memory:**

| Addr | Data | Tag | Index | Block_offset |
|---|---|---|---|---|
| 0 | 78 | 00 | 0 | 0 |
| 1 | 29 | 00 | 0 | 1 |
| 2 | 120 | 00 | 1 | 0 |
| 3 | 123 | 00 | 1 | 1 |
| 4 | 71 | 01 | 0 | 0 |
| 5 | 150 | 01 | 0 | 1 |
| 6 | 162 | 01 | 1 | 0 |
| 7 | 173 | 01 | 1 | 1 |
| 8 | 18 | 10 | 0 | 0 |
| 9 | 21 | 10 | 0 | 1 |
| 10 | 33 | 10 | 1 | 0 |
| 11 | 28 | 10 | 1 | 1 |
| 12 | 19 | 11 | 0 | 0 |
| 13 | 200 | 11 | 0 | 1 |
| 14 | 210 | 11 | 1 | 0 |
| 15 | 225 | 11 | 1 | 1 |

24

# Set-associative cache (REF 5)

| Processor | Cache | Memory |
|---|---|---|

**Processor**

Ld  R1 ← M[ 1 ]
Ld  R2 ← M[ 5 ]
St  R2 → M[ 7 ]
St  R1 → M[ 4 ]
Ld  R3 ← M[ 0 ]
Ld  R2 ← M[ 8 ]

R0
R1  29
R2  150
R3

**Cache**

V d tag  data

| lru | 1 | 0 | 0 | 78 |
| | | | | 29 |
| | 1 | 1 | 1 | 29 |
| | | | | 150 |
| | 1 | 1 | 1 | 162 |
| | | | | 150 |
| lru | 0 | | | |

Misses:  3

Hits:     1

**Memory**

| 0 | 78 |
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 71 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

| Tag | Index | Block_offset |
|---|---|---|
| 00 | 0 | 0 |
| 00 | 0 | 1 |
| 00 | 1 | 0 |
| 00 | 1 | 1 |
| 01 | 0 | 0 |
| 01 | 0 | 1 |
| 01 | 1 | 0 |
| 01 | 1 | 1 |
| 10 | 0 | 0 |
| 10 | 0 | 1 |
| 10 | 1 | 0 |
| 10 | 1 | 1 |
| 11 | 0 | 0 |
| 11 | 0 | 1 |
| 11 | 1 | 0 |
| 11 | 1 | 1 |

25

# Set-associative cache (REF 5)

| | Processor | | Cache | | Memory | | | Tag | Index | Block_offset |
|---|---|---|---|---|---|---|---|---|---|---|

**Processor**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

| R0 | |
|---|---|
| R1 | 29 |
| R2 | 150 |
| R3 | 78 |

**Cache**

V d tag data

lru | 1 | 0 | 0 | 78 |
| | | | 29 |
| 1 | 1 | 1 | 29 |
| | | | 150 |
| 1 | 1 | 1 | 162 |
| | | | 150 |
lru | 0 | | | |
| | | | |

Misses: 3

Hits: 2

**Memory**

| 0 | 78 |
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 71 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

| Tag | Index | Block_offset |
|---|---|---|
| 00 | 0 | 0 |
| 00 | 0 | 1 |
| 00 | 1 | 0 |
| 00 | 1 | 1 |
| 01 | 0 | 0 |
| 01 | 0 | 1 |
| 01 | 1 | 0 |
| 01 | 1 | 1 |
| 10 | 0 | 0 |
| 10 | 0 | 1 |
| 10 | 1 | 0 |
| 10 | 1 | 1 |
| 11 | 0 | 0 |
| 11 | 0 | 1 |
| 11 | 1 | 0 |
| 11 | 1 | 1 |

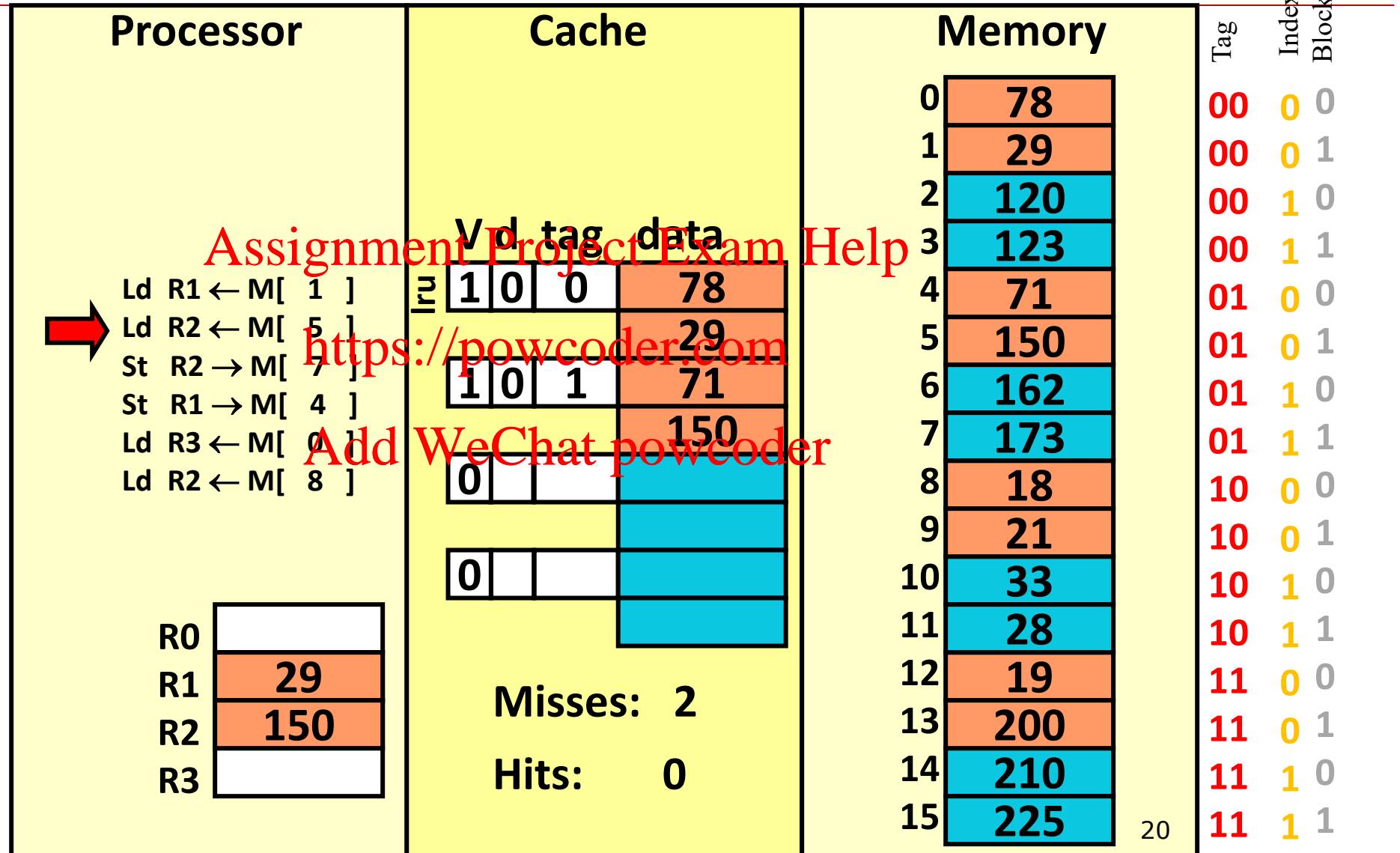Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

26

# Set-associative cache (REF 6)

| | Processor | Cache | Memory | | Tag | Index | Block_offset |
|---|---|---|---|---|---|---|---|

**Processor**

Ld  R1 ← M[ 1 ]
Ld  R2 ← M[ 5 ]
St  R2 → M[ 7 ]
St  R1 → M[ 4 ]
Ld  R3 ← M[ 0 ]
➡ Ld  R2 ← M[ 8 ]

| R0 | |
|---|---|
| R1 | 29 |
| R2 | 150 |
| R3 | 78 |

**Cache**

V d tag data

| 1 | 0 | 0 | 78 |
|---|---|---|---|
| | | | 29 |

lru

| 1 | 1 | 1 | 29 |
|---|---|---|---|
| | | | 150 |

| 1 | 1 | 1 | 162 |
|---|---|---|---|
| | | | 150 |

lru

| 0 | | | |
|---|---|---|---|

Misses:  3

Hits:   2

**Memory**

| 0 | 78 | 00 | 0 | 0 |
|---|---|---|---|---|
| 1 | 29 | 00 | 0 | 1 |
| 2 | 120 | 00 | 1 | 0 |
| 3 | 123 | 00 | 1 | 1 |
| 4 | 71 | 01 | 0 | 0 |
| 5 | 150 | 01 | 0 | 1 |
| 6 | 162 | 01 | 1 | 0 |
| 7 | 173 | 01 | 1 | 1 |
| 8 | 18 | 10 | 0 | 0 |
| 9 | 21 | 10 | 0 | 1 |
| 10 | 33 | 10 | 1 | 0 |
| 11 | 28 | 10 | 1 | 1 |
| 12 | 19 | 11 | 0 | 0 |
| 13 | 200 | 11 | 0 | 1 |
| 14 | 210 | 11 | 1 | 0 |
| 15 | 225 | 11 | 1 | 1 |

# Set-associative cache (REF 6)

| Processor | Cache | Memory | Tag | Index | Block_offset |
|---|---|---|---|---|---|

**Processor**

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

| R0 | |
|---|---|
| R1 | 29 |
| R2 | 150 |
| R3 | 78 |

**Cache**

V d tag data

| 1 | 0 | 0 | 78 |
|---|---|---|---|
| | | | 29 |

lru
| 1 | 1 | 1 | 29 |
|---|---|---|---|
| | | | 150 |

| 1 | 1 | 1 | 162 |
|---|---|---|---|
| | | | 150 |

lru
| 0 | | | |
|---|---|---|---|

Misses:   3

Hits:     2

**Memory**

| 0 | 78 |
|---|---|
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 29 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

| Tag | Index | Block_offset |
|---|---|---|
| 00 | 0 | 0 |
| 00 | 0 | 1 |
| 00 | 1 | 0 |
| 00 | 1 | 1 |
| 01 | 0 | 0 |
| 01 | 0 | 1 |
| 01 | 1 | 0 |
| 01 | 1 | 1 |
| 10 | 0 | 0 |
| 10 | 0 | 1 |
| 10 | 1 | 0 |
| 10 | 1 | 1 |
| 11 | 0 | 0 |
| 11 | 0 | 1 |
| 11 | 1 | 0 |
| 11 | 1 | 1 |

28

# Set-associative cache (REF 6)



| | Processor | Cache | Memory | Tag | Index | Block_offset |
|---|---|---|---|---|---|---|

Processor:

Ld R1 ← M[ 1 ]
Ld R2 ← M[ 5 ]
St R2 → M[ 7 ]
St R1 → M[ 4 ]
Ld R3 ← M[ 0 ]
Ld R2 ← M[ 8 ]

R0
R1  29
R2  18
R3  78

Cache:

V d  tag   data

| 1 | 0 | 0 | 78 |
| | | | 29 |
lru
| 1 | 0 | 2 | 18 |
| | | | 21 |
| 1 | 1 | 1 | 162 |
| | | | 150 |
lru
| 0 | | | |

Misses:  4

Hits:  2

Memory:

| 0 | 78 |
| 1 | 29 |
| 2 | 120 |
| 3 | 123 |
| 4 | 29 |
| 5 | 150 |
| 6 | 162 |
| 7 | 173 |
| 8 | 18 |
| 9 | 21 |
| 10 | 33 |
| 11 | 28 |
| 12 | 19 |
| 13 | 200 |
| 14 | 210 |
| 15 | 225 |

| Tag | Index | Block_offset |
|---|---|---|
| 00 | 0 | 0 |
| 00 | 0 | 1 |
| 00 | 1 | 0 |
| 00 | 1 | 1 |
| 01 | 0 | 0 |
| 01 | 0 | 1 |
| 01 | 1 | 0 |
| 01 | 1 | 1 |
| 10 | 0 | 0 |
| 10 | 0 | 1 |
| 10 | 1 | 0 |
| 10 | 1 | 1 |
| 11 | 0 | 0 |
| 11 | 0 | 1 |
| 11 | 1 | 0 |
| 11 | 1 | 1 |

29

**Reasons for cache misses a.k.a.**
**The 3C's of Cache Misses**

**Compulsory** miss

First reference to any block will always miss
Also sometimes called a **"cold start"** miss

Assignment Project Exam Help

**Capacity** miss

https://powcoder.com

Cache is too small to hold all the data
Would have had a hit with an infinite cache

Add WeChat powcoder

**Conflict** miss

Would have had a hit with a fully associative cache

# Classifying Cache Misses

Can we classify a cache miss into one of the following?

**Compulsory miss**

**Capacity miss**

**Conflict miss**

Yes! Simulate three different caches

Simulate with a cache of unlimited size (cache size = memory size)

– Any misses must be **compulsory misses**

Simulate again with a fully associative cache of the intended size

- Any new misses must be **capacity misses**

Simulate a third time, with the actual intended cache

- Any new misses must be **conflict misses**

# Fixing cache misses

**Compulsory** misses

    First reference to an address

    No way to completely avoid these

    Reduce by **increasing block size (spatial locality)**

    This reduces the total number of blocks

**Capacity** misses

    Would have a hit with a large enough cache

    Reduce by **building a bigger cache**

**Conflict** misses

    Would have had a hit with a fully associative cache

    Cache does not have enough associativity

    Reduce by **increasing associativity**

# 3 C's Sample Problem

Consider a cache with the following configuration:

write-allocate

Cache size = 64 bytes

Block size = 16 bytes

2-way associative.

16-bit byte-addressable ISA (address size is 16 bits)

LRU replacement policy.

Assume the cache is empty at the start.

For the following memory accesses, indicate whether the reference is a hit or miss, and the type of a miss (compulsory, conflict, capacity)

# 3 C's Practice Problem – Address sequence

**Address**

0x00

0x14

0x27

0x08

0x38

0x4A

0x18

0x27

0x0F

0x40

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# 3 C's Practice Problem – Simulate infinite cache

**Address**

| Tag | Block_offset |
| --- | --- |
| 0x0 | 0 |
| 0x1 | 4 |
| 0x2 | 7 |
| 0x0 | 8 |
| 0x3 | 8 |
| 0x4 | A |
| 0x1 | 8 |
| 0x2 | 7 |
| 0x0 | F |
| 0x4 | 0 |

## 3 C's Practice Problem – Simulate fully associative cache

**Address**

| | |
|---|---|
| 0x0 | 0 |
| 0x1 | 4 |
| 0x2 | 7 |
| 0x0 | 8 |
| 0x3 | 8 |
| 0x4 | A |
| 0x1 | 8 |
| 0x2 | 7 |
| 0x0 | F |
| 0x4 | 0 |

Tag

Block_offset

**tag**

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# 3 C's Practice Problem – Simulate given set associative cache

## Address

| Address | |
|---|---|
| 0x0 | 0 |
| 0x1 | 4 |
| 0x2 | 7 |
| 0x0 | 8 |
| 0x3 | 8 |
| 0x4 | A |
| 0x1 | 8 |
| 0x2 | 7 |
| 0x0 | F |
| 0x4 | 0 |

Tag

Set index

Block_offset

**tag  data**          **tag  data**

**Set 0**

**Set 1**

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

## 3 C's Practice Problem – 3 C's

| Address | Infinite | FA | SA | 3Cs |
|---------|----------|-----|-----|-----|
| 0x00 | M | M | M | |
| 0x14 | M | M | M | |
| 0x27 | M | M | M | |
| 0x08 | H | H | H | |
| 0x38 | M | M | M | |
| 0x4A | M | M | M | |
| 0x18 | H | M | H | |
| 0x27 | H | M | M | |
| 0x0F | H | M | M | |
| 0x40 | H | H | M | |

## 3 C's Practice Problem – 3 C's

| Address | Infinite | FA | SA | 3Cs |
|---------|----------|----|----|-----|
| 0x00 | M | M | M | Compulsory |
| 0x14 | M | M | M | Compulsory |
| 0x27 | M | M | M | Compulsory |
| 0x08 | H | H | H | --- |
| 0x38 | M | M | M | Compulsory |
| 0x4A | M | M | M | Compulsory |
| 0x18 | H | M | H | --- |
| 0x27 | H | M | M | Capacity |
| 0x0F | H | M | M | Capacity |
| 0x40 | H | H | M | Conflict |

# Cache Parameters vs. Miss Rate

Cache Size

Block Size

<span style="color:red">Assignment Project Exam Help</span>

Associativity

<span style="color:red">https://powcoder.com</span>

Replacement policy

<span style="color:red">Add WeChat powcoder</span>

## Questions to ask

Can block size be not power of 2?

Can number of sets be not power of 2?

Can number of ways be not power of 2?

Can we have 3-way set associative cache?

# Cache Size

Cache size in the total data (not including tag) capacity

bigger can exploit temporal locality better

not ALWAYS better

Too large a cache adversely affects hit & miss latency

smaller is faster => bigger is slower

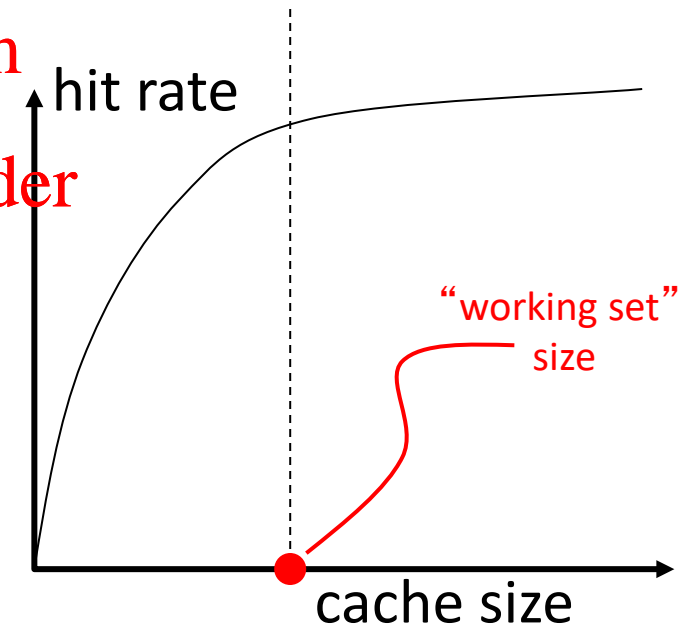access time may degrade critical path

Too small a cache

doesn't exploit temporal locality well

useful data replaced often

Working set: the whole set of data
executing application references

**Within a time interval**



hit rate

"working set"
size

cache size

# Block size (also called Line size)

Block size is the data that is associated with an address tag

    Sub-blocking: A block divided into multiple pieces (each with V bit)

        Can improve "write" performance

Too small blocks

    don't exploit spatial locality well

    have larger tag overhead

Too large blocks
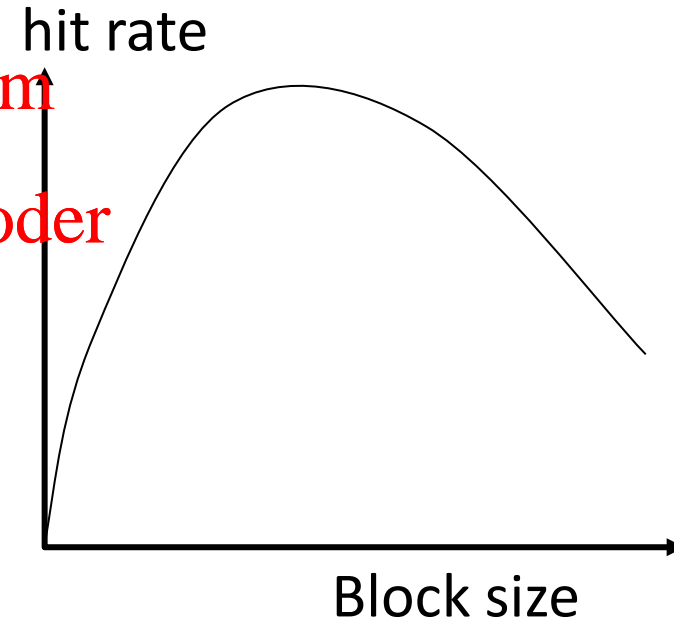
    too few total # of blocks

        likely-useless data transferred

        Extra bandwidth/energy consumed

hit rate

Block size

# Associativity

How many blocks map to the same set (same set index)?

Larger associativity

    lower miss rate, less variation among programs

    diminishing returns <span style="color:red">Assignment Project Exam Help</span>

<span style="color:red">https://powcoder.com</span>

hit rate

Smaller associativity

    lower cost <span style="color:red">Add WeChat powcoder</span>

    faster hit time

        Especially important for L1 caches

Power of 2 associativity?

associativity