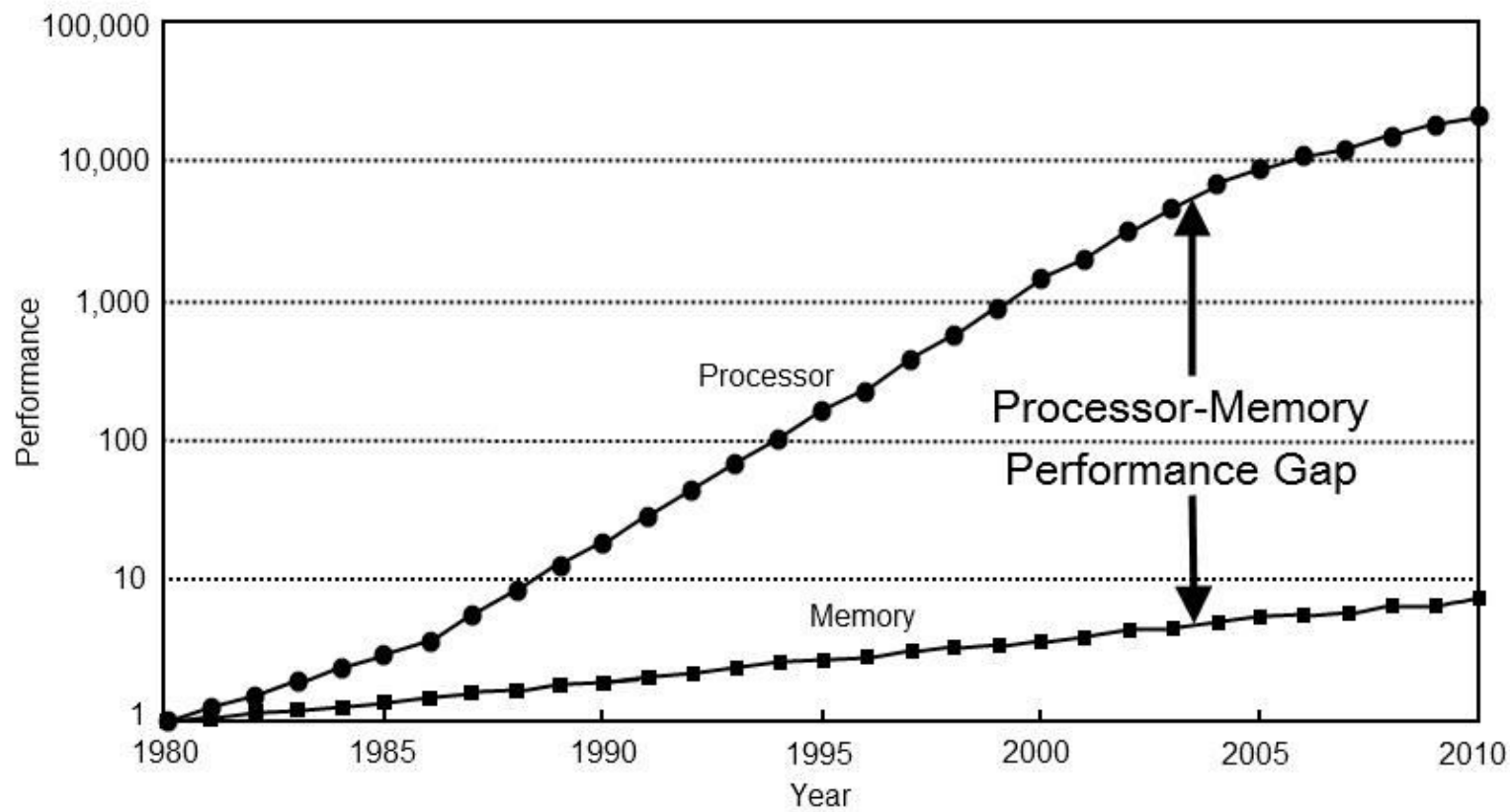
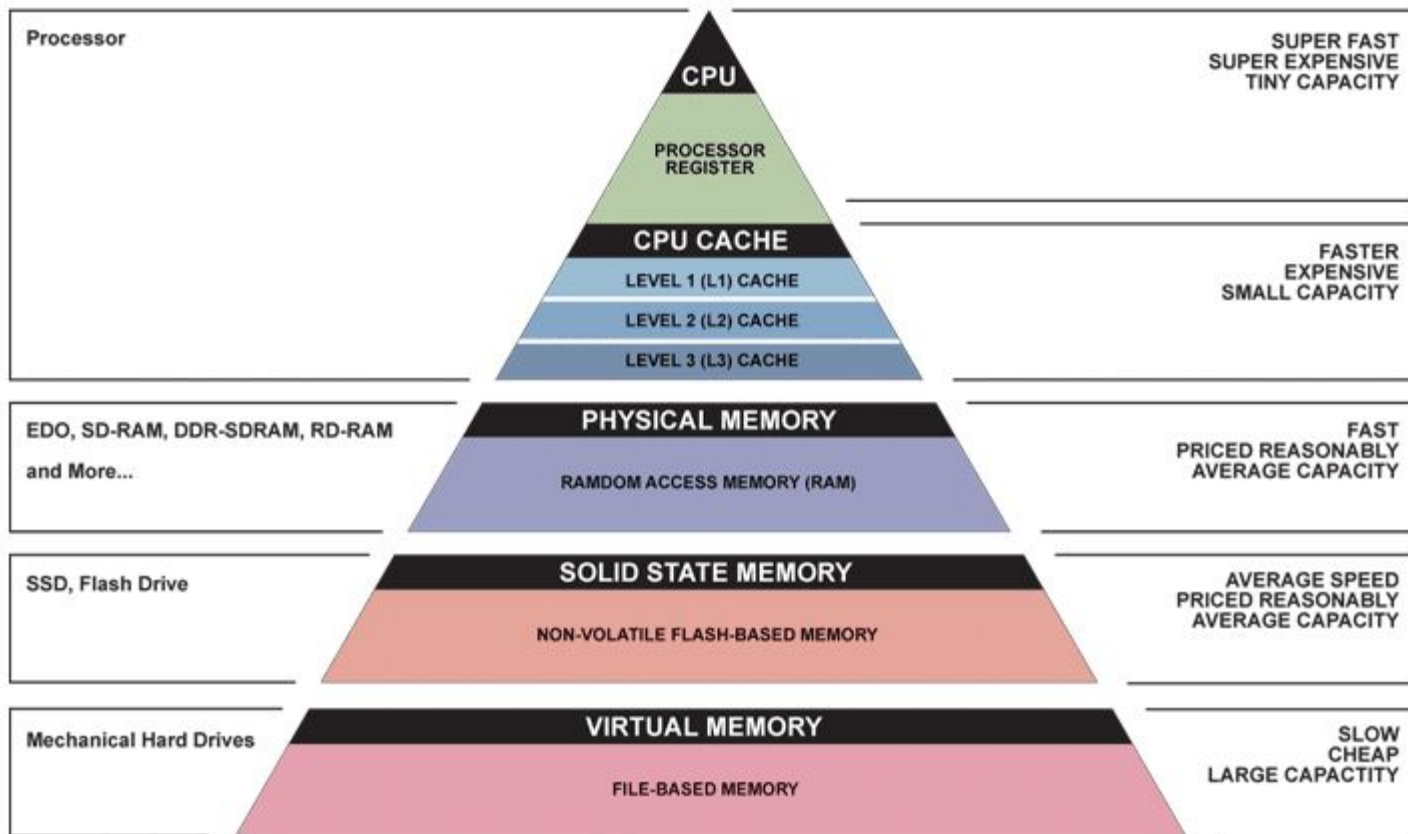


# CASS: Exercise session 7

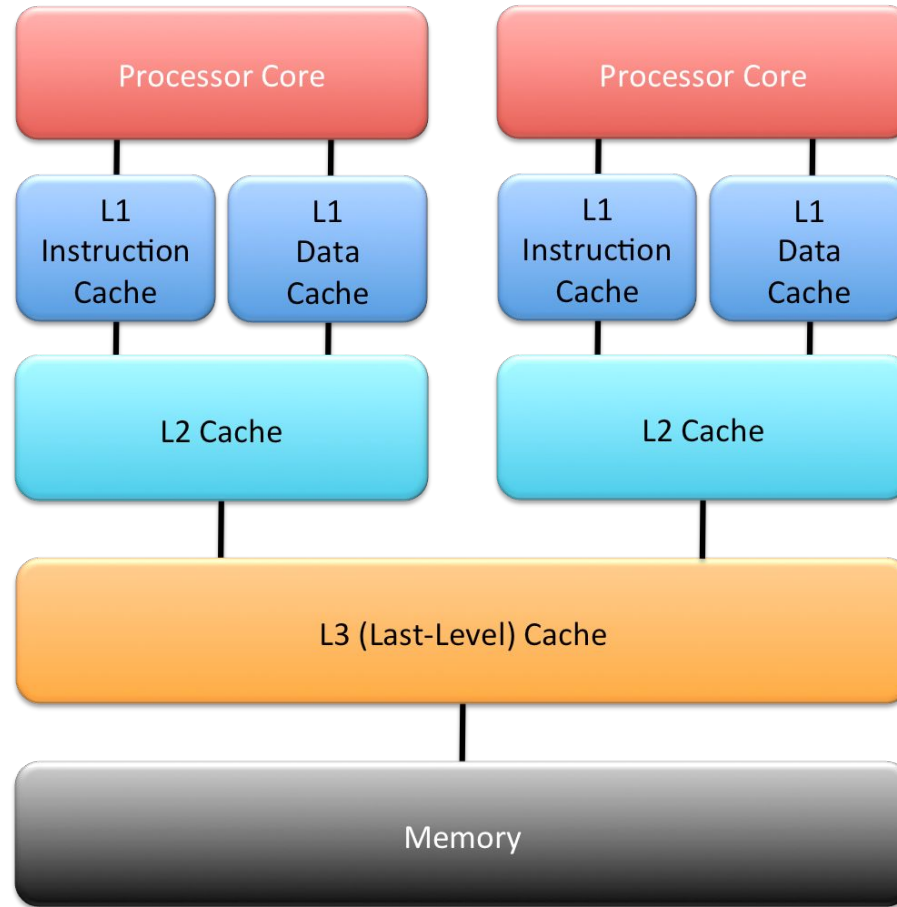
Caches and Microarchitectural Timing Attacks







▲ Simplified Computer Memory Hierarchy  
Illustration: Ryan J. Leng



Caches are small.

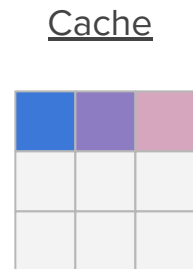
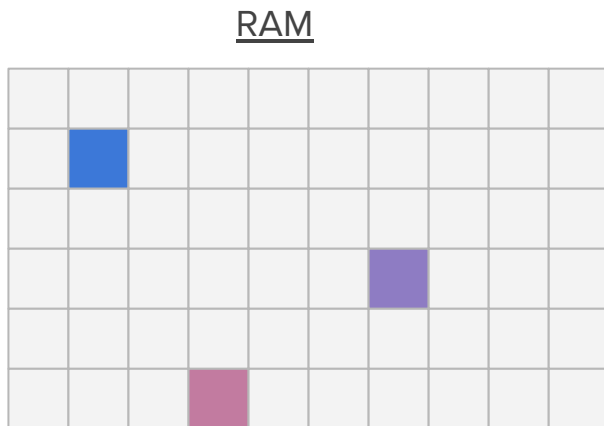
How can they make a difference?

# “Locality principle”

Programs access a relatively **small portion** of  
**address space** at a time

# Temporal locality

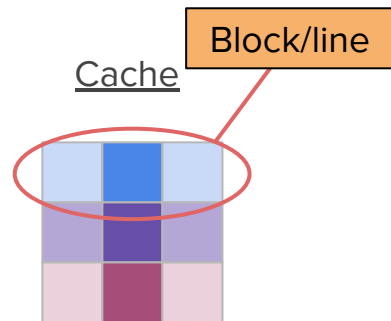
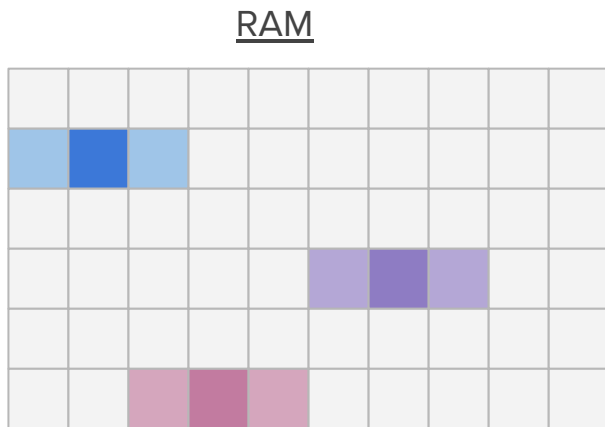
- Memory location accessed?
  - Same location *likely* accessed again soon!
  - Add it to cache!





# Spatial locality

- Memory location accessed?
  - Nearby locations *likely* accessed soon!
- Transfer **block** to cache
  - Cache line = block = #bytes transferred at once

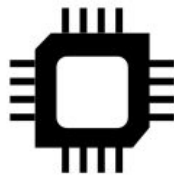


# CPU cache operation



**Cache principle:** CPU speed  $\gg$  DRAM latency  $\rightarrow$  *cache code/data*

```
while true do  
  maccess(&a);  
endwh
```



**CPU + cache**

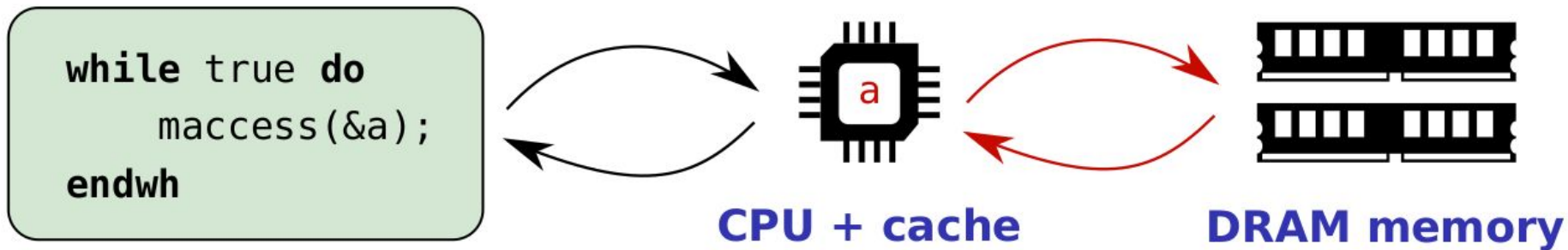


**DRAM memory**

# CPU cache operation



**Cache miss:** Request data from (slow) DRAM upon first use

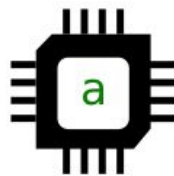
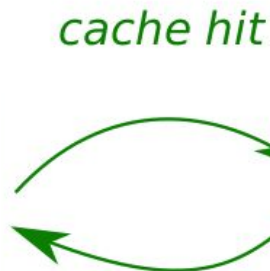


# CPU cache operation



**Cache hit:** No DRAM access required for subsequent uses

```
while true do  
  maccess(&a);  
endwh
```



**CPU + cache**



**DRAM memory**

# Cache performance

$$\text{miss rate} = (1 - \text{hit rate})$$

# Conclusion:

Small amount of memory, used in smart way






***Can we abuse timing differences as an attacker to infer secret information? :-)***



CPU's offer **instructions to**  
***flush*** the cache!

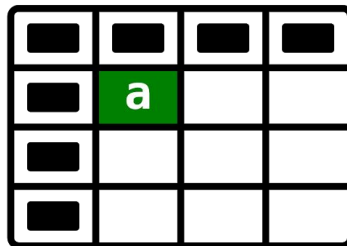


# Flush+Reload: Cache timing attacks on shared memory



```
if secret do
    maccess(&a);
else
    maccess(&b);
endif
```


*'a' is accessible  
to attacker*



**CPU cache**

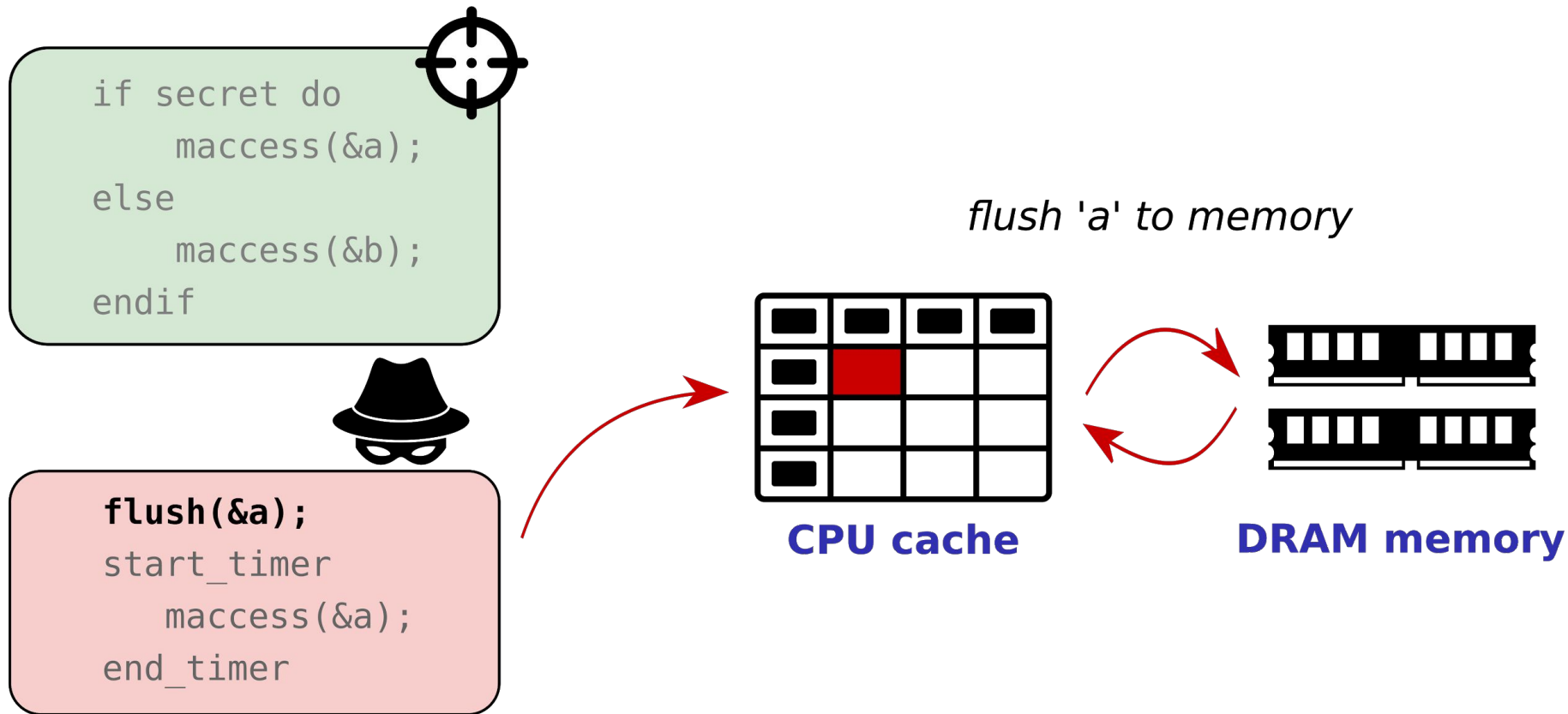


**DRAM memory**

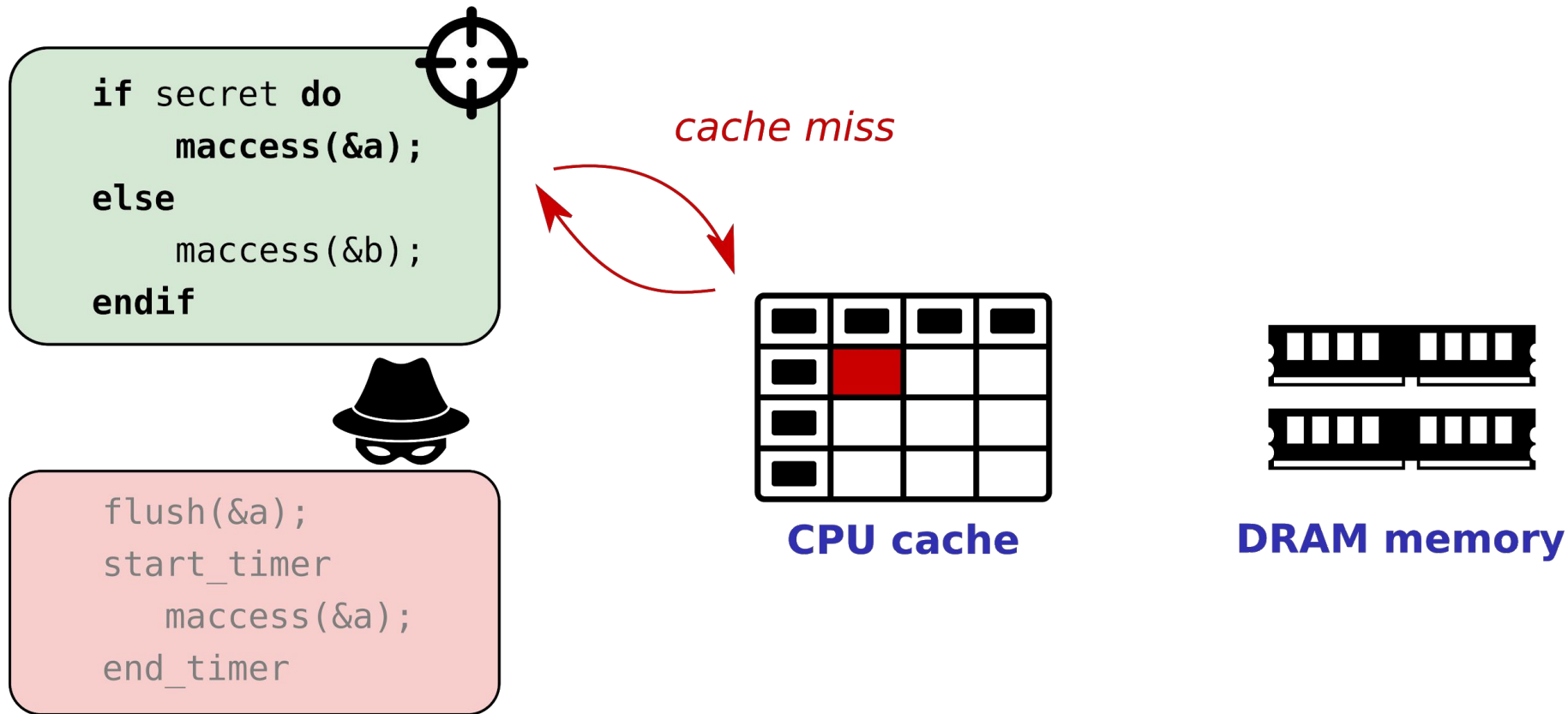


```
flush(&a);
start_timer
    maccess(&a);
end_timer
```

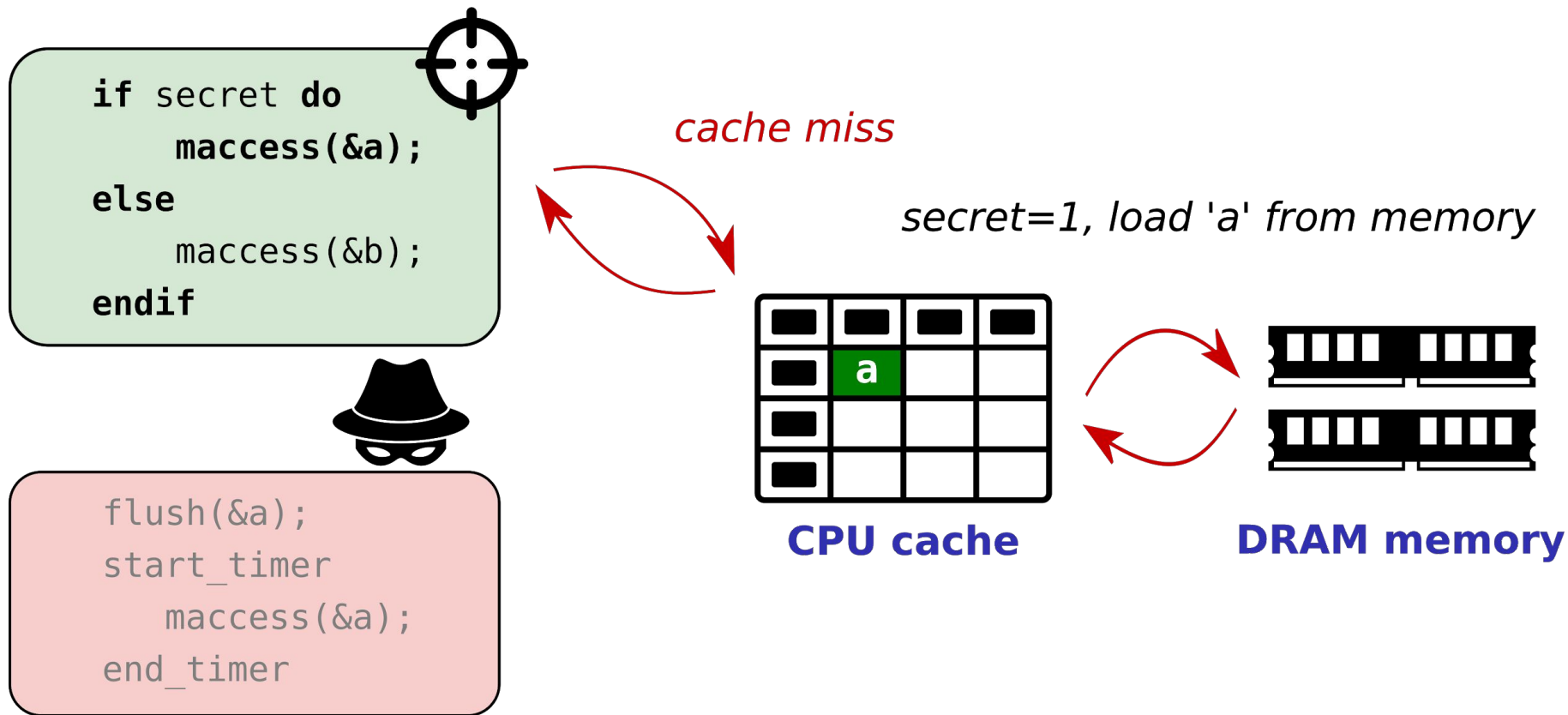
# Flush+Reload: Cache timing attacks on shared memory



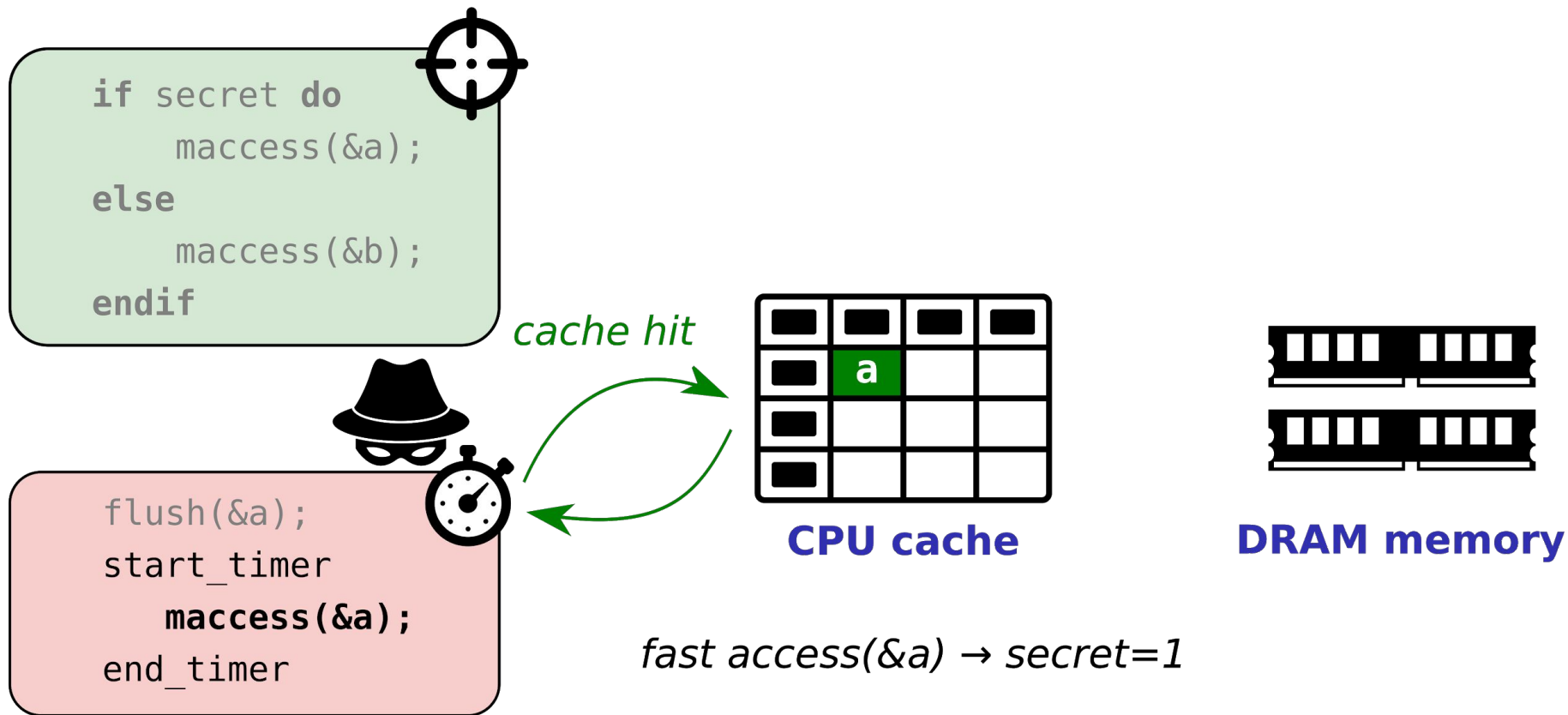
# Flush+Reload: Cache timing attacks on shared memory



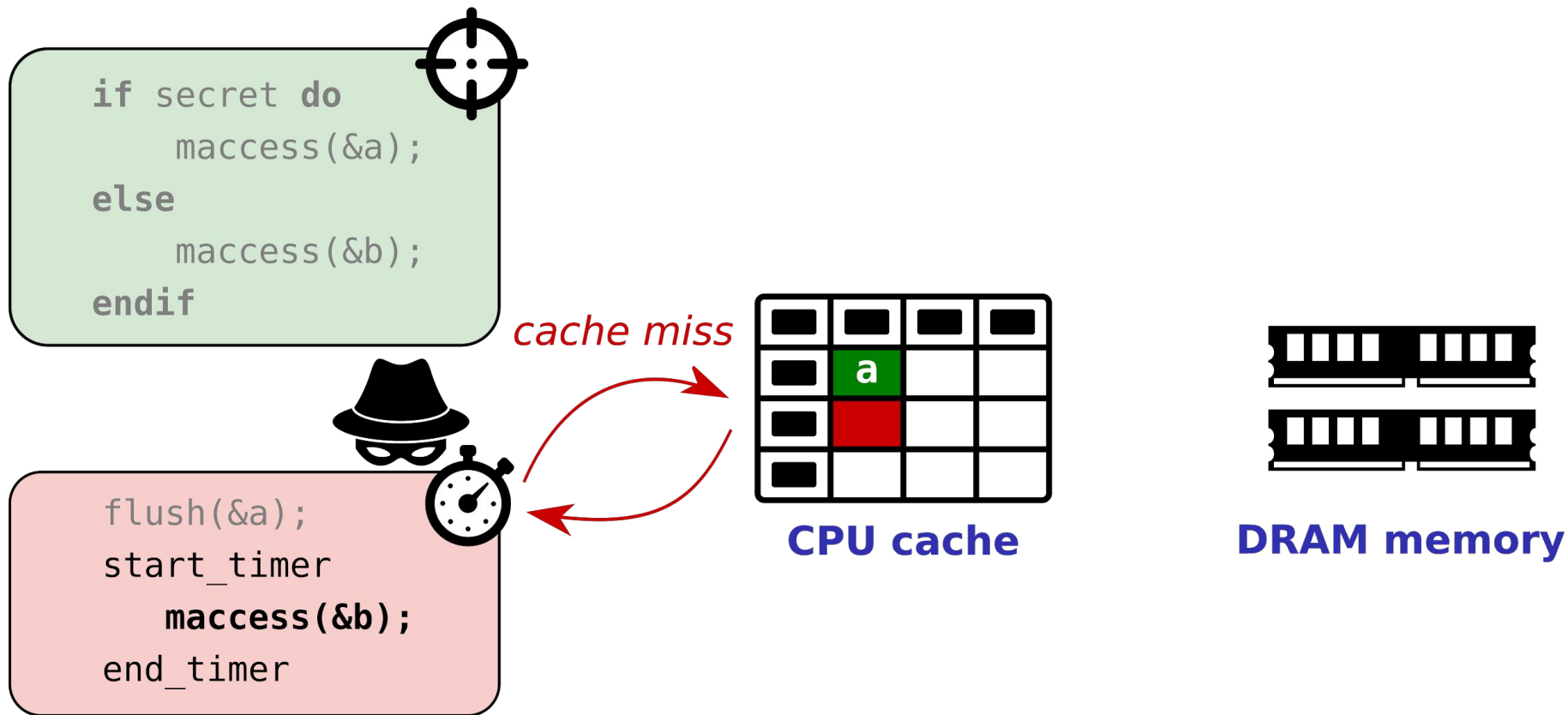
# Flush+Reload: Cache timing attacks on shared memory



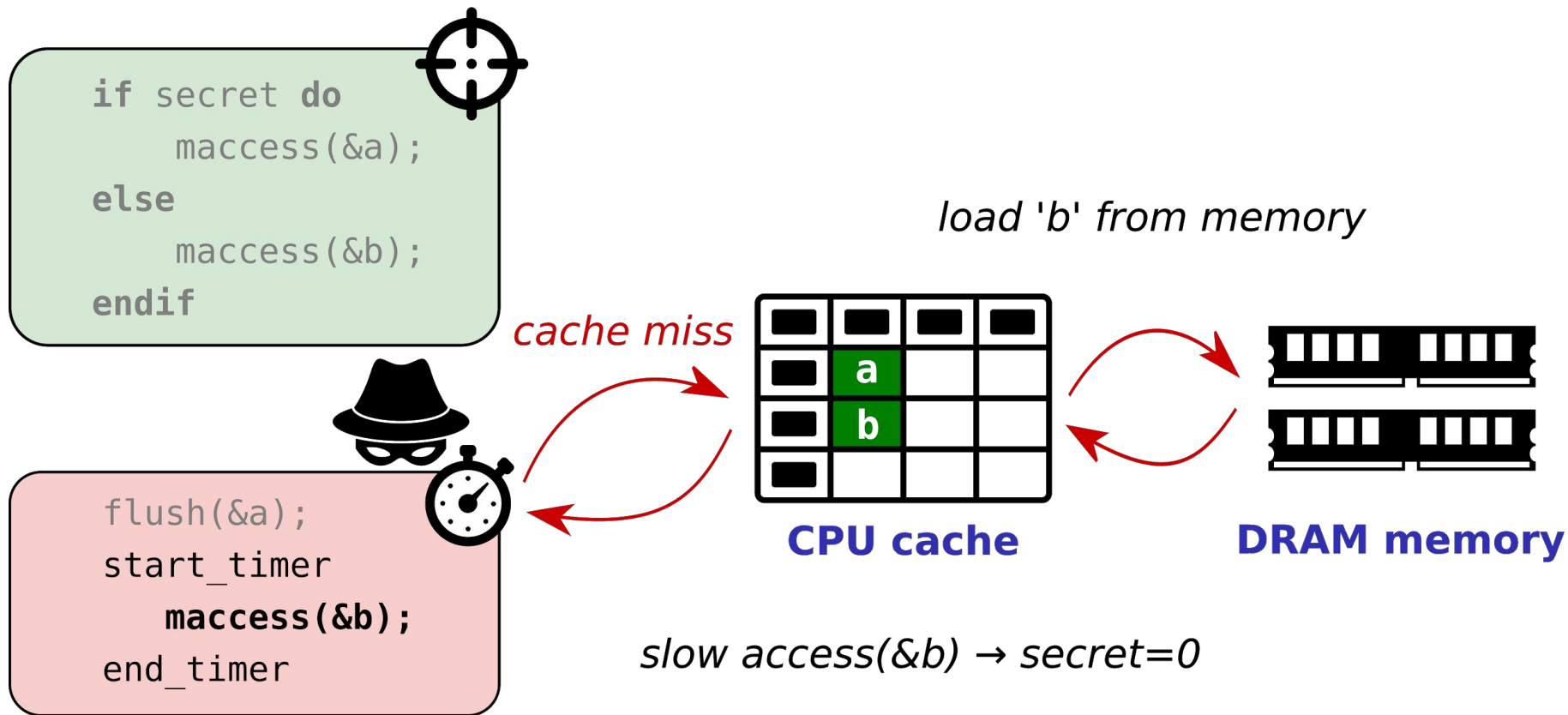
# Flush+Reload: Cache timing attacks on shared memory



# Flush+Reload: Cache timing attacks on shared memory



# Flush+Reload: Cache timing attacks on shared memory





# Flush+Reload Limitations

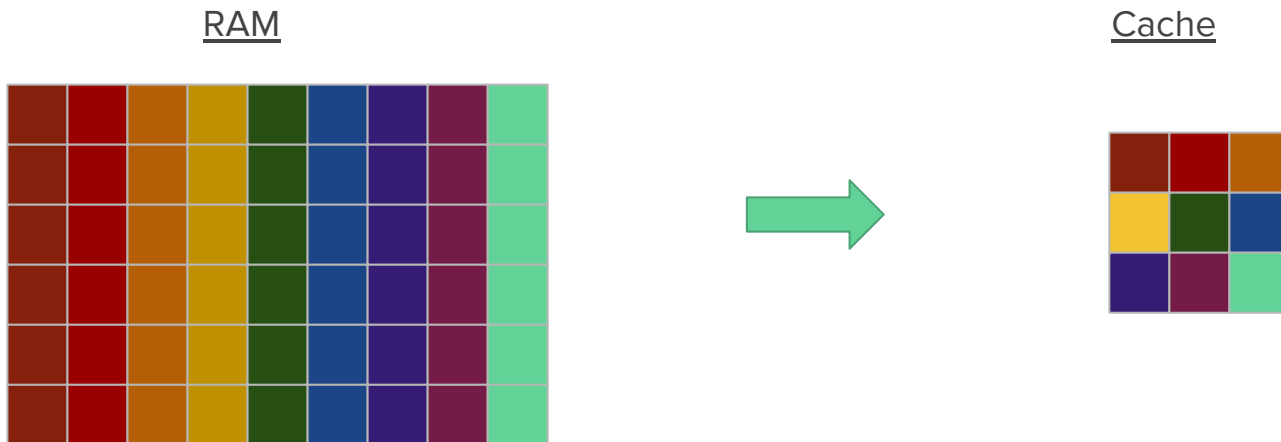
- Very **reliable and easy** attack
- But requires **shared memory** between victim and attacker apps
  - otherwise attacker cannot *flush* victim cache lines

→ **Generic** attack: does *not require knowledge* of internal cache organization, but only applicable in *limited scenarios* (when victim and attacker share memory)

*Where* in the cache should we store the contents of a memory location?

# Direct mapping: concept

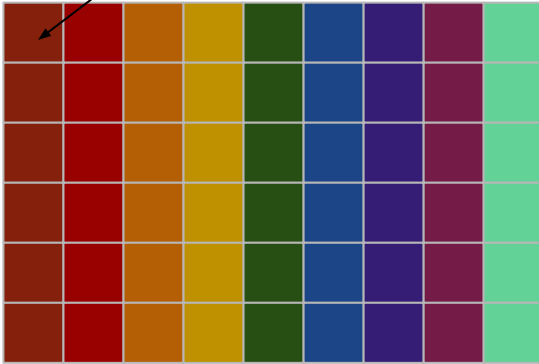
Notice: memory block index % 9 = cache block index



# Direct mapping: concept

Each cache block has multiple addresses!

RAM



size of a block:  $2^n$  bytes

=>

Address in same block share all but last n bits

Cache



size of a block:  $2^n$  bytes  $\Rightarrow$  Address in same block share all but last  $n$  bits

Why?

Take  $n = 3$

$\Rightarrow$  8 memory (byte)-addresses per block

**Notice #1: shared bits = *memory block index* in binary!!**

**Notice #2: *last  $n$  bits* are in index within the block**

Block index	Memory byte-address range	Binary range
<b>0</b>	[ 0, 7 ]	[ <b>00</b> 000 - <b>00</b> 111]
<b>1</b>	[ 8, 15 ]	[ <b>01</b> 000 - <b>01</b> 111]
<b>2</b>	[ 16, 23 ]	[ <b>10</b> 000 - <b>10</b> 111]
...	...	...

size of a block:  $2^n$  bytes  $\Rightarrow$  Address in same block share all but last  $n$  bits

Why?

Take  $n = 3$

$\Rightarrow$  8 memory (byte)-addresses per block

**Notice #1: shared bits = *memory block index* in binary!!**

**Notice #2: *last  $n$  bits* are in index within the block**

**Remember**  
memory block index % cache  
size (in blocks) = cache block  
index

Block index	Memory byte-address range	Binary range
<b>0</b>	[ 0, 7 ]	[ <b>00</b> 000 - <b>00</b> 111]
<b>1</b>	[ 8, 15 ]	[ <b>01</b> 000 - <b>01</b> 111]
<b>2</b>	[ 16, 23 ]	[ <b>10</b> 000 - <b>10</b> 111]
...	...	...

binary(0xabc) =

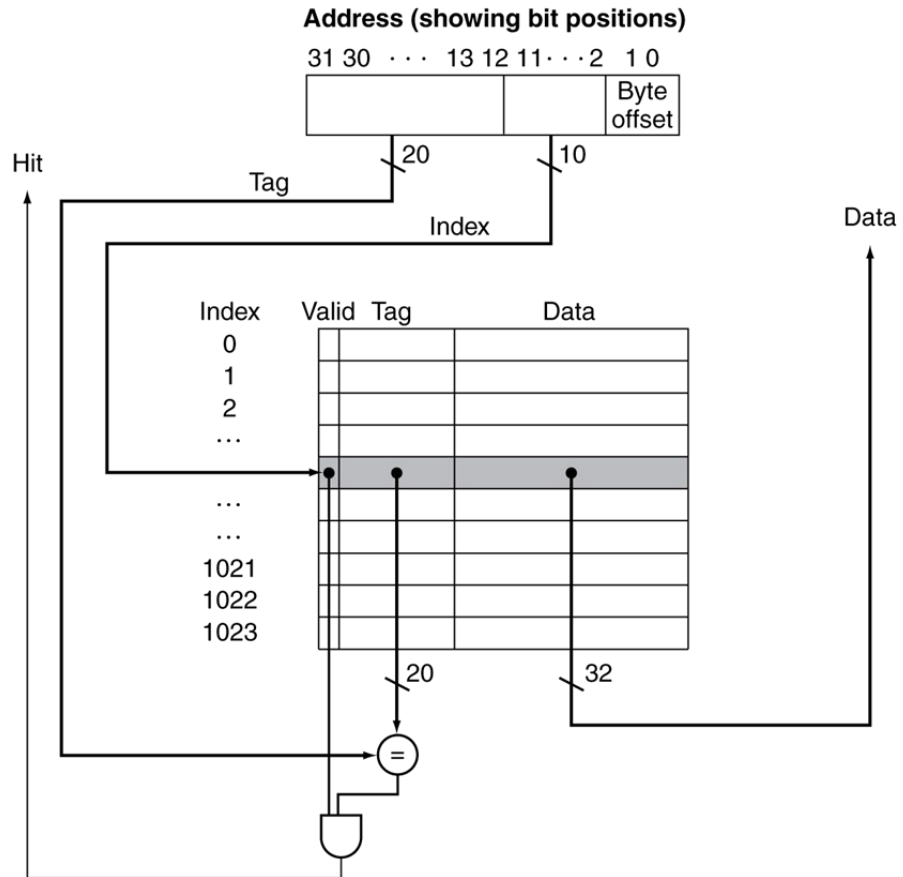
$\{1010\} \{10111\} \{100\}$   
= 10      = 23      = 4

Assume:  
Cache size: 32 blocks ( $2^5$ )  
Block size: 8 bytes ( $2^3$ )

<u>Index</u>	Valid	Tag	Data
0	0	???	???
...	...	...	
10111 (23)	1	1010	8 bytes (0xab8 - 0xabf)
...	...	...	
11111 (31)	???	???	???

Notice:

Index determines where to look  
Tag determines hit or miss

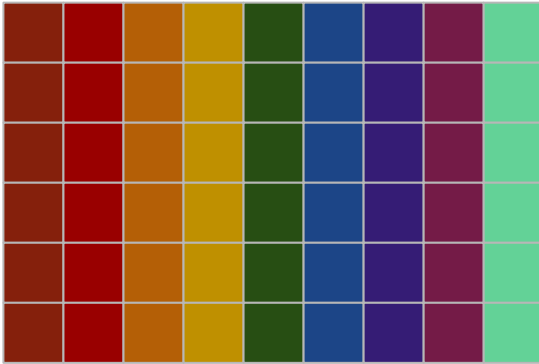




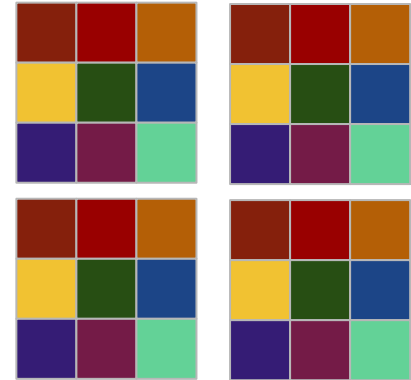
# Set-associativity: concept

Notice: each memory block can be located in  $n$  different **sets**

RAM



4-way Cache

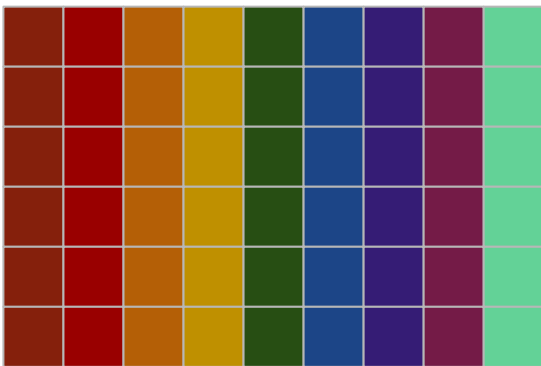


# Set-associativity: concept

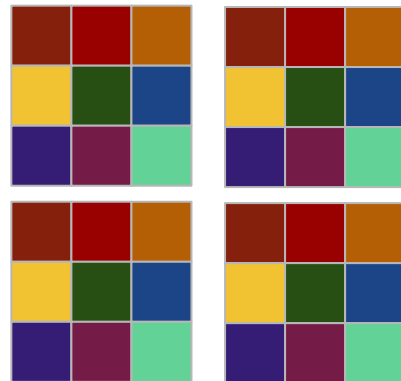
**Notice**  
 $n = 1$ : direct mapping!

Notice: each memory block can be located in  $n$  different **sets**

RAM



4-way Cache

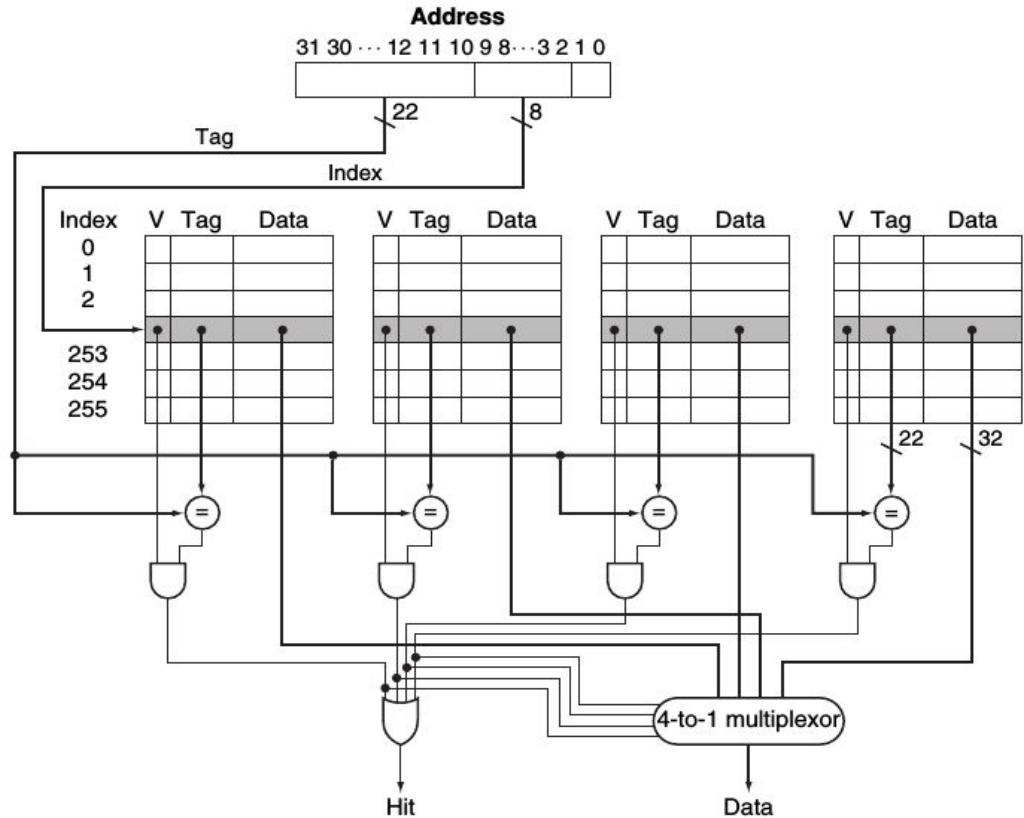


Notice:

Index determines where to look  
**BUT:** multiple (#sets) possibilities

*Check all:*

Tag determines hit or miss



# Fully associative cache

- n-way set associative cache
  - $n = \text{\#memory-blocks}$
- Best performance possible
  - only miss when a new unique address is accessed
  - any repeated request is a hit!

E.g., a cache with 8 blocks:

[illegible]

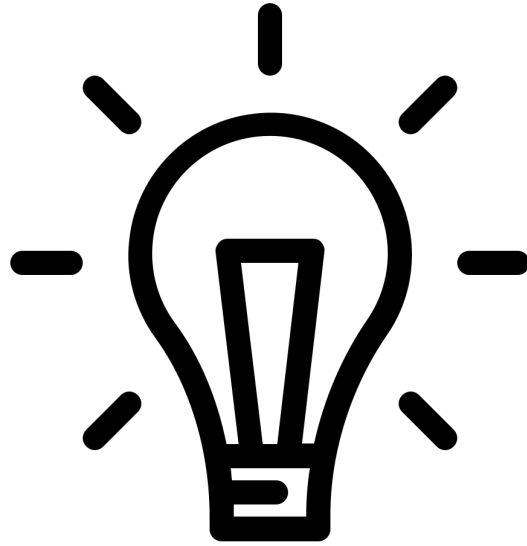
# Fully associative cache

Basically, searching the entire cache without any indexing  
– sequentially super slow!

Solution: search in different sets in parallel.  
Need  $n$  comparators for any associative cache.


[illegible]





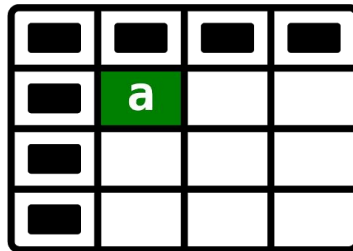
***Can we leverage detailed knowledge of mapping schemes and cache collisions to mount more advanced cache timing attack?***

# Prime+Probe: Cache timing attacks across protection domains

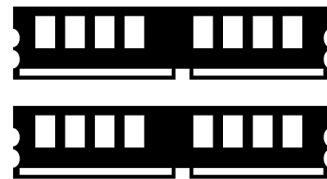


```
if secret do
    maccess(&a);
else
    maccess(&b);
endif
```


'a' is **not** accessible  
to attacker



**CPU cache**



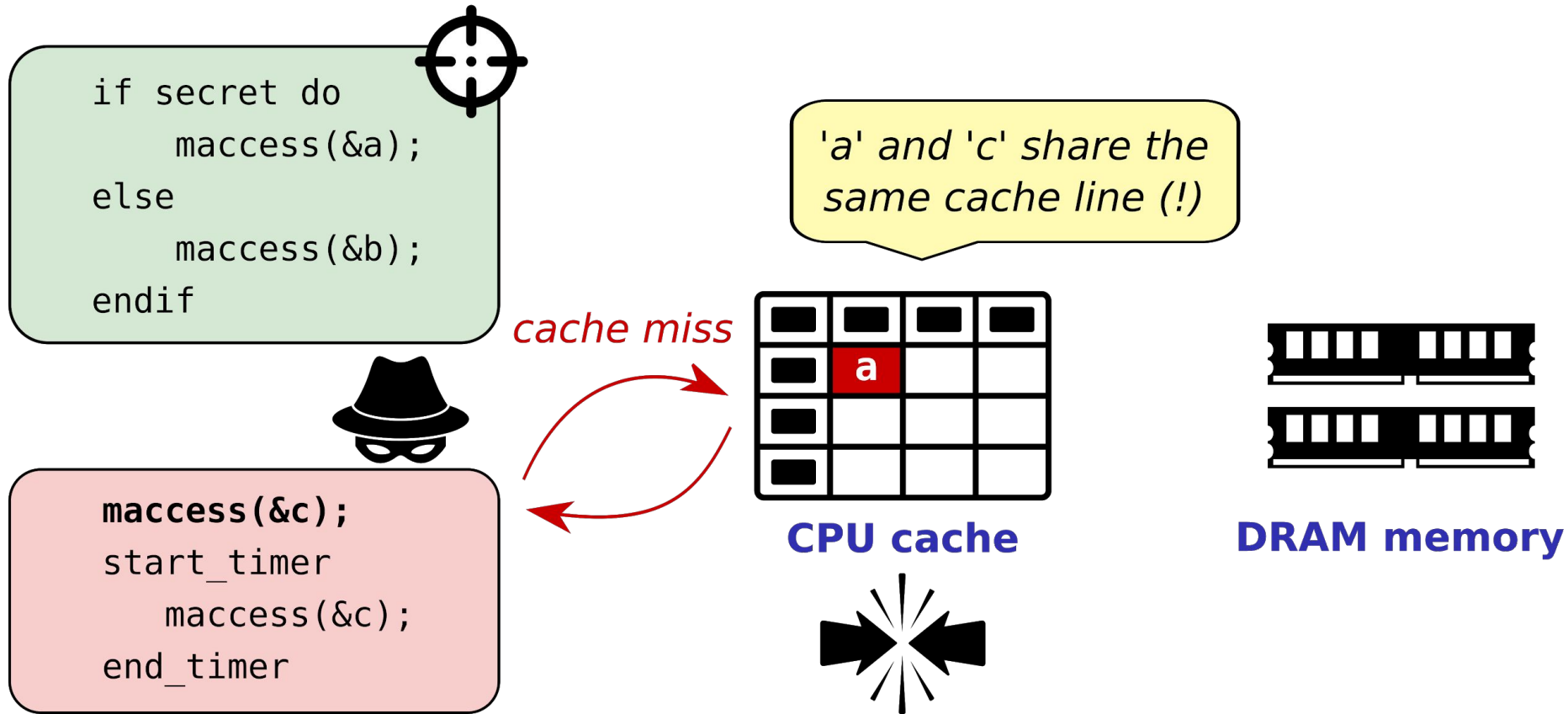
**DRAM memory**



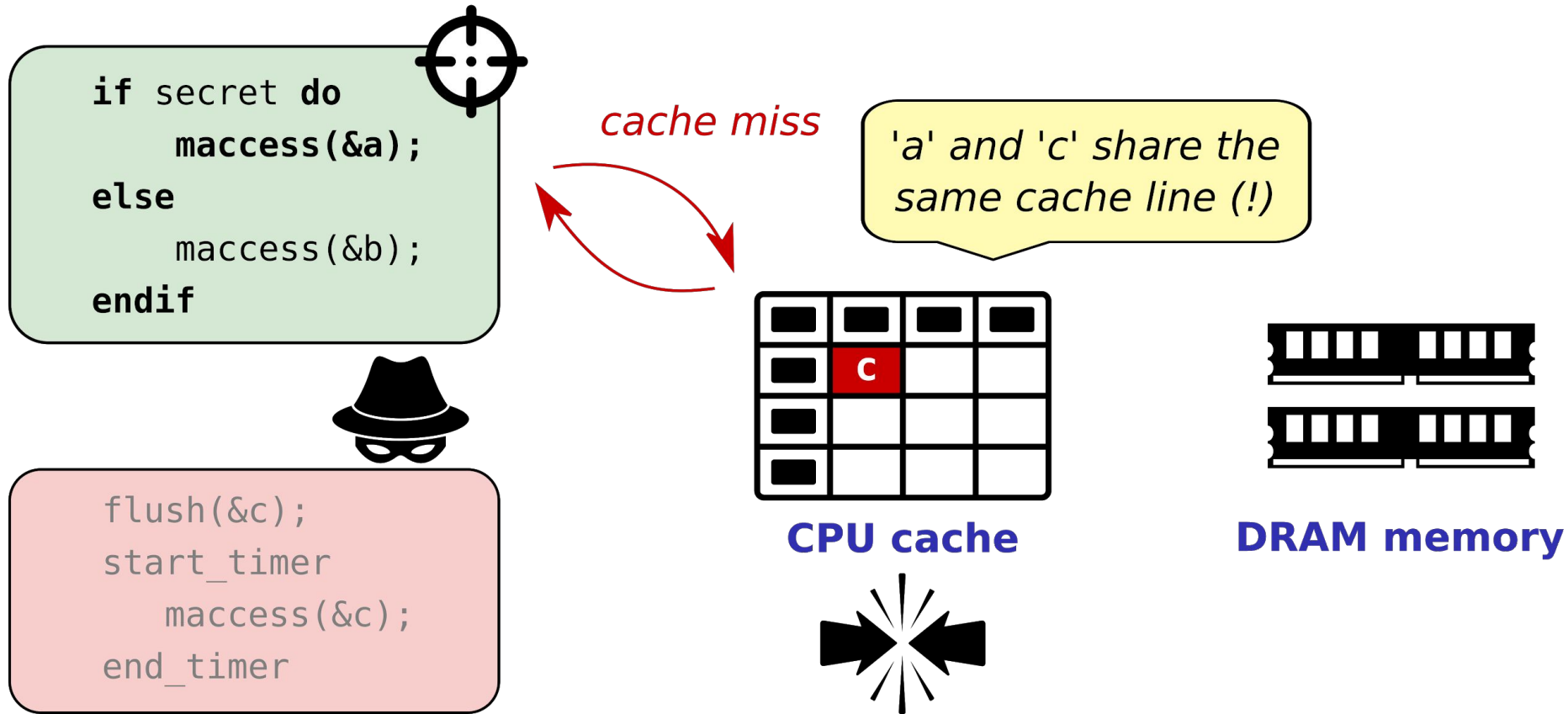
```
maccess(&c);
start_timer
    maccess(&c);
end_timer
```



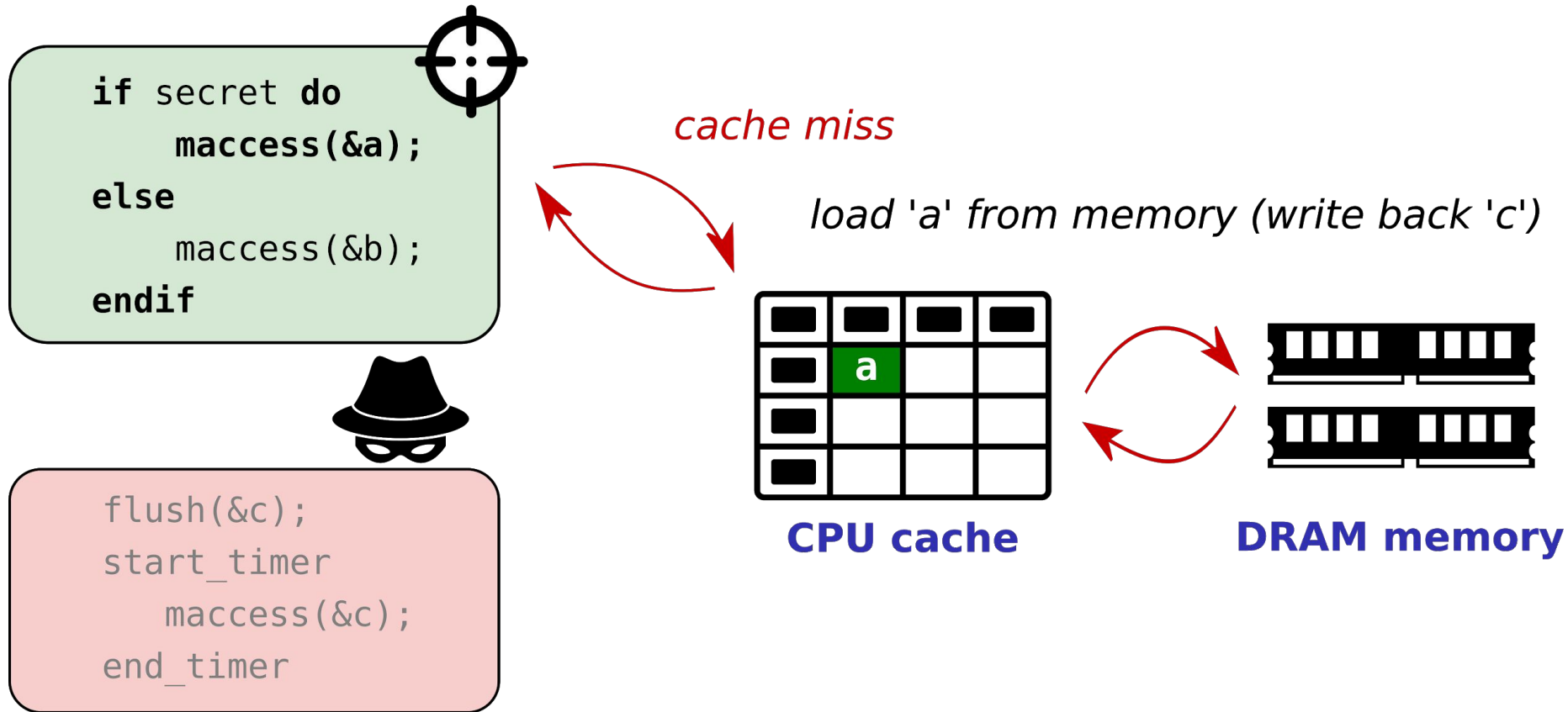
# Prime+Probe: Cache timing attacks across protection domains



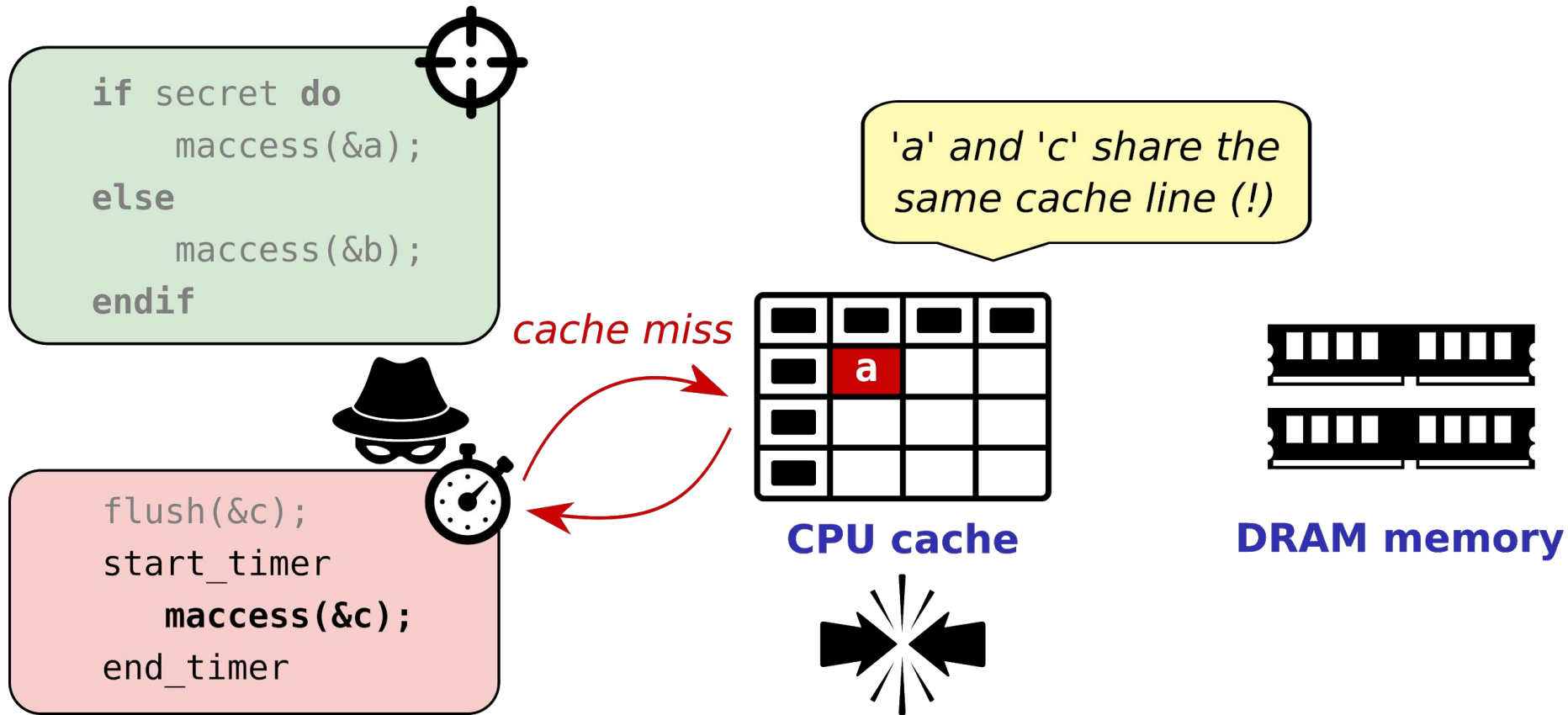
# Prime+Probe: Cache timing attacks across protection domains



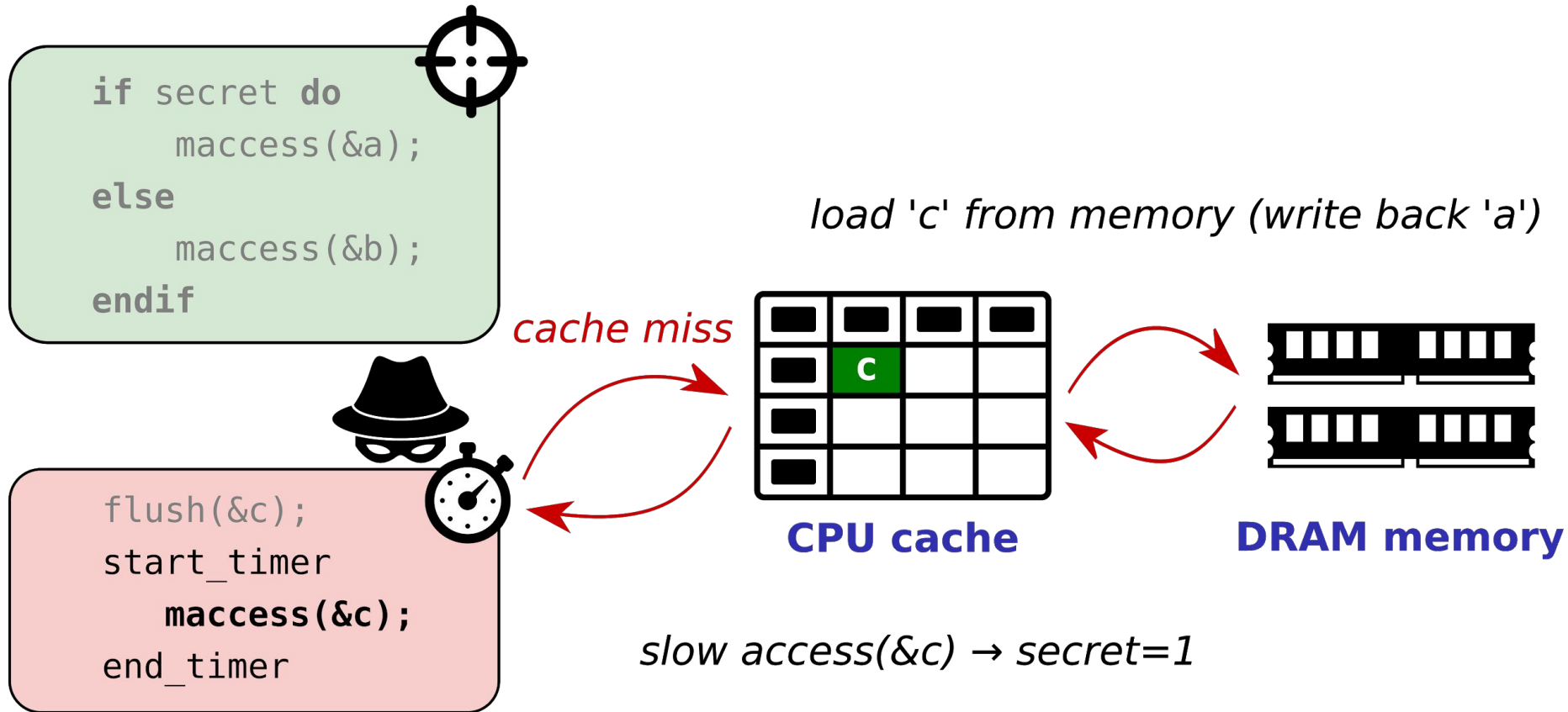
# Prime+Probe: Cache timing attacks across protection domains



# Prime+Probe: Cache timing attacks across protection domains



# Prime+Probe: Cache timing attacks across protection domains



# *Appendix*

# Example 1: direct mapping 1-word blocks

Word address references:  
22, 26, 22, 27, 3, 27, 53, 55

010110 → miss  
011010 → miss  
010110 → hit  
011011 → miss  
000011 → miss  
011011 → miss  
110101 → miss  
110111 → miss

Index	V	Tag	Data
000	0		
001	0		
010	0		
011	0		
100	0		
101	0		
110	1	010	22
111	0		

Index	V	Tag	Data
000	0		
001	0		
010	1	011	26
011	0		
100	0		
101	0		
110	1	010	22
111	0		

Index	V	Tag	Data
000	0		
001	0		
010	1	011	26
011	1	011	27
100	0		
101	0		
110	1	010	22
111	0		

Index	V	Tag	Data
000	0		
001	0		
010	1	011	26
011	1	000	3
100	0		
101	0		
110	1	010	22
111	0		

Index	V	Tag	Data
000	0		
001	0		
010	1	011	26
011	1	011	27
100	0		
101	0		
110	1	010	22
111	0		

Index	V	Tag	Data
000	0		
001	0		
010	1	011	26
011	1	011	27
100	0		
101	0	110	53
110	1	010	22
111	0		

Index	V	Tag	Data
000	0		
001	0		
010	1	011	26
011	1	011	27
100	0		
101	0	110	53
110	1	010	22
111	1	110	55

# Example 2: direct mapping 4-word blocks

Word address references:  
22, 26, 22, 27, 3, 27, 53, 55

010110 → 101 → miss  
011010 → 110 → miss  
010110 → 101 → hit  
011011 → 110 → hit  
000011 → 000 → miss  
011011 → 110 → hit  
110101 → 110 → miss  
110111 → 110 → hit

Index			Data			
	V	Tag	W1	W2	W3	W4
000	1					
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	0					
111	0					

Index			Data			
	V	Tag	W1	W2	W3	W4
000	0					
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	011	24	25	26	27
111	0					

Index			Data			
	V	Tag	W1	W2	W3	W4
000	1	000	0	1	2	3
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	011	24	25	26	27
111	0					

Index			Data			
	V	Tag	W1	W2	W3	W4
000	1	000	0	1	2	3
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	110	52	53	54	55
111	0					



# Example 2: direct mapping 4-word blocks

Word address references:  
22, 26, 22, 27, 3, 27, 53, 55

010110 → 101 → miss  
011010 → 110 → miss  
010110 → 101 → hit  
011011 → 110 → hit  
000011 → 000 → miss  
011011 → 110 → hit  
110101 → 110 → miss  
110111 → 110 → hit

Miss rate  
decreased.  
Why?

Index			Data			
			W1	W2	W3	W4
000	1					
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	0					
111	0					

Index			Data			
			W1	W2	W3	W4
000	1	000	0	1	2	3
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	011	24	25	26	27
111	0					

Index			Data			
			W1	W2	W3	W4
000	0					
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	011	24	25	26	27
111	0					

Index			Data			
			W1	W2	W3	W4
000	1	000	0	1	2	3
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	110	52	53	54	55
111	0					

# Example 2: direct mapping 4-word blocks

Word address references:  
22, 26, 22, 27, 3, 27, 53, 55

010110 → 101 → miss  
011010 → 110 → miss  
010110 → 101 → hit  
011011 → 110 → hit  
000011 → 000 → miss  
011011 → 110 → hit  
110101 → 110 → miss  
110111 → 110 → hit

**Miss rate  
decreased.  
Why?**

**But miss penalty  
increased.  
Why?**

Index			Data			
			W1	W2	W3	W4
000	1					
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	0					
111	0					

Index			Data			
			W1	W2	W3	W4
000	1	000	0	1	2	3
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	011	24	25	26	27
111	0					

Index			Data			
			W1	W2	W3	W4
000	0					
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	011	24	25	26	27
111	0					

Index			Data			
			W1	W2	W3	W4
000	1	000	0	1	2	3
001	0					
010	0					
011	0					
100	0					
101	1	010	20	21	22	23
110	1	110	52	53	54	55
111	0					

# Example 3: 2-way set-associative 1-word blocks

Word address references:  
22, 26, 22, 27, 3, 27, 53, 55

010110 → miss  
011010 → miss  
010110 → hit  
011011 → miss  
000011 → miss  
011011 → hit  
110101 → miss  
110111 → miss

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22		
11	0				

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	1	0110	27		


Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	1	1101	53		
10	1	0101	22	0110	26
11	1	0110	27	0000	3

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	0				

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	1	0110	27	0000	3

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	1	1101	53		
10	1	0101	22	0110	26
11	1	0110	27	1101	55

Replacement rule:  
**Least Recently Used**



# Example 3: 2-way set-associative 1-word blocks

Word address references:  
22, 26, 22, 27, 3, 27, 53, 55

010110 → miss  
011010 → miss  
010110 → hit  
011011 → miss  
000011 → miss  
011011 → hit  
110101 → miss  
110111 → miss

Miss rate  
decreased.  
Why?

Trade-off?

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22		
11	0				

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	1	0110	27		

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	1	1101	53		
10	1	0101	22	0110	26
11	1	0110	27	0000	3

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	0				

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	1	0110	27	0000	3

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	1	1101	53		
10	1	0101	22	0110	26
11	1	0110	27	1101	55

# Example 3: 2-way set-associative 1-word blocks

Word address references:  
22, 26, 22, 27, 3, 27, 53, 55

010110 → miss  
011010 → miss  
010110 → hit  
011011 → miss  
000011 → miss  
011011 → hit  
110101 → miss  
110111 → miss

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22		
11	0				

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	1	0110	27		

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	1	1101	53		
10	1	0101	22	0110	26
11	1	0110	27	0000	3

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	0				

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	0				
10	1	0101	22	0110	26
11	1	0110	27	0000	3

Set	V	Tag0	Data0	Tag1	Data1
00	0				
01	1	1101	53		
10	1	0101	22	0110	26
11	1	0110	27	1101	55

Miss rate  
decreased.  
Why?

Hit time  
increased!  
We search longer

