# CASS: Exercise session 7

Caches and Microarchitectural Timing Attacks

**Processor**

**CPU**

SUPER FAST
SUPER EXPENSIVE
TINY CAPACITY

PROCESSOR
REGISTER

**CPU CACHE**

FASTER
EXPENSIVE
SMALL CAPACITY

LEVEL 1 (L1) CACHE

LEVEL 2 (L2) CACHE

LEVEL 3 (L3) CACHE

EDO, SD-RAM, DDR-SDRAM, RD-RAM
and More...

**PHYSICAL MEMORY**

FAST
PRICED REASONABLY
AVERAGE CAPACITY

RAMDOM ACCESS MEMORY (RAM)

SSD, Flash Drive

**SOLID STATE MEMORY**

AVERAGE SPEED
PRICED REASONABLY
AVERAGE CAPACITY

NON-VOLATILE FLASH-BASED MEMORY

Mechanical Hard Drives

**VIRTUAL MEMORY**

SLOW
CHEAP
LARGE CAPACITY

FILE-BASED MEMORY

▲ Simplified Computer Memory Hierarchy
Illustration: Ryan J. Leng

Processor Core | Processor Core

L1 Instruction Cache | L1 Data Cache | L1 Instruction Cache | L1 Data Cache

L2 Cache | L2 Cache

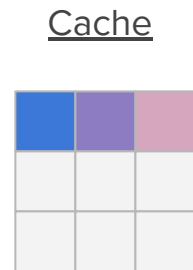L3 (Last-Level) Cache
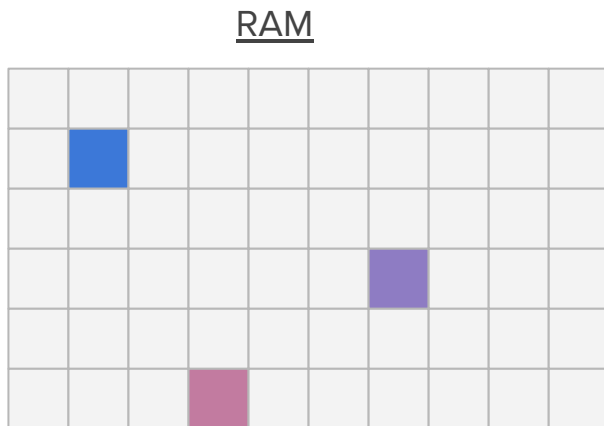
Memory

# Caches are small.

How can they make a difference?

# "Locality principle"

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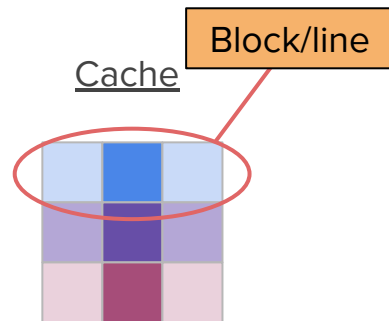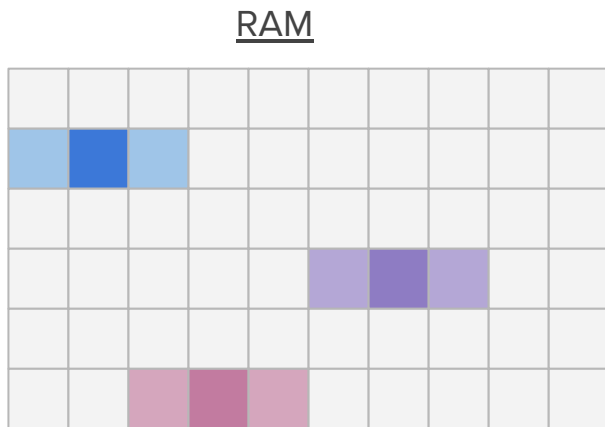**address space** at a time

# Temporal locality

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

RAM

Cache

# Spatial locality

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

RAM

Cache

Block/line

# 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

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

*cache miss*

**CPU + cache**

**DRAM memory**

# CPU cache operation

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

*cache hit*

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

a

**CPU + cache**

**DRAM memory**

# Cache performance

**miss rate** = (1 – hit rate)

# Conclusion:

Small amount of memory, used in smart way

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

# Flush+Reload: Cache timing attacks on shared memory

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

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

'a' *is accessible
to attacker*

**CPU cache**

**DRAM memory**

# Flush+Reload: Cache timing attacks on shared memory

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

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

*flush 'a' to memory*

**CPU cache**

**DRAM memory**

# Flush+Reload: Cache timing attacks on shared memory

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

*cache miss*

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

**CPU cache**

**DRAM memory**

# Flush+Reload: Cache timing attacks on shared memory

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

*cache miss*

*secret=1, load 'a' from memory*

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

**CPU cache**

**DRAM memory**

# Flush+Reload: Cache timing attacks on shared memory

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

*cache hit*

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

**a**

**CPU cache**

**DRAM memory**

*fast access(&a) → secret=1*

# Flush+Reload: Cache timing attacks on shared memory

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

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

*cache miss*

**CPU cache**

**DRAM memory**

# Flush+Reload: Cache timing attacks on shared memory

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

*load 'b' from memory*

*cache miss*
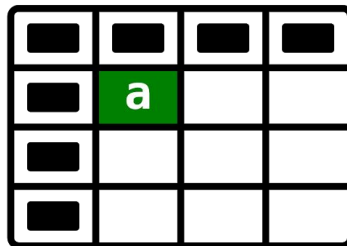
```
flush(&a);
start_timer
    maccess(&b);
end_timer
```

**CPU cache**

**DRAM memory**

*slow access(&b) → secret=1*

# 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

RAM

Cache

# Direct mapping: concept

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

Each cache block has multiple addresses!

RAM

Cache

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

Why?
Take n = 3
=> 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 |
|---|---|---|
| **0** | [ 0, 7 ] | [**00**000 - **00**111] |
| **1** | [ 8, 15 ] | [**01**000 - **01**111] |
| **2** | [ 16, 23 ] | [**10**000 - **10**111] |
| ... | ... | ... |

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

Why?
Take n = 3
=> 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 |
|---|---|---|
| **0** | [ 0, 7 ] | [**00**000 - **00**111] |
| **1** | [ 8, 15 ] | [**01**000 - **01**111] |
| **2** | [ 16, 23 ] | [**10**000 - **10**111] |
| ... | ... | ... |

binary(0xabc) =

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

| Index | 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

Address (showing bit positions)

31 30 · · · 13 12 11 · · · 2 1 0

Byte offset

Hit

Tag

20

Index

10

Data

Index   Valid   Tag                    Data

0
1
2
…

…
…
1021
1022
1023

20

32

=

# Set-associativity: concept

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

RAM

4-way Cache

# Set-associativity: concept

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

RAM

4-way Cache

Notice:
Index determines where to look
BUT: multiple (#ways) possibilities

Check all:
Tag determines hit or miss

# Fully associative cache

- n-way set associative cache
  - n = #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:

| Set | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | | | | | | | | | | | | | | | | | | | | | | | | |

# Fully associative cache

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

Solution: search in different ways in parallel.
Need **n** comparators for any associative cache.

| Set | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data |
|-----|---|-----|------|---|-----|------|---|-----|------|---|-----|------|---|-----|------|---|-----|------|---|-----|------|---|-----|------|
| 0   |   |     |      |   |     |      |   |     |      |   |     |      |   |     |      |   |     |      |   |     |      |   |     |      |

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

# Prime+Probe: Cache timing attacks across protection domains

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

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

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

**CPU cache**

**DRAM memory**

# Prime+Probe: Cache timing attacks across protection domains



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

'a' and 'c' share the same cache line (!)

cache miss

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

**CPU cache**

**DRAM memory**

# Prime+Probe: Cache timing attacks across protection domains

# Prime+Probe: Cache timing attacks across protection domains

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

*cache miss*

*load 'a' from memory (write back 'c')*

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



**CPU cache**

**DRAM memory**

# Prime+Probe: Cache timing attacks across protection domains

# Prime+Probe: Cache timing attacks across protection domains

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

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

*cache miss*

*load 'c' from memory (write back 'a')*

**CPU cache**

**DRAM memory**

*slow access(&c) → secret=1*

*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 | 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 | 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 | 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 | 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 | V | Tag | 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 | V | Tag | 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 | V | Tag | 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 | V | Tag | 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 | 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 | 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 | 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 | 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 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 | 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 | 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 | 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:
**L**east **R**ecently **U**sed

# 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  | 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  | 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    | 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 | 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 | 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 | 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