



# CS773-2025-Spring: Computer Architecture for Performance and Security

## Lecture 2: Catch the Cache



**ON SILENT MODE PLEASE**

# Microarchitecture 101: World with no caches

*North pole ☹️*

**Core**

32-bit Address

Data

200 to 300 cycles

Minimizing costly DRAM accesses  
is critical for performance

**Costly DRAM  
accesses ☹️**

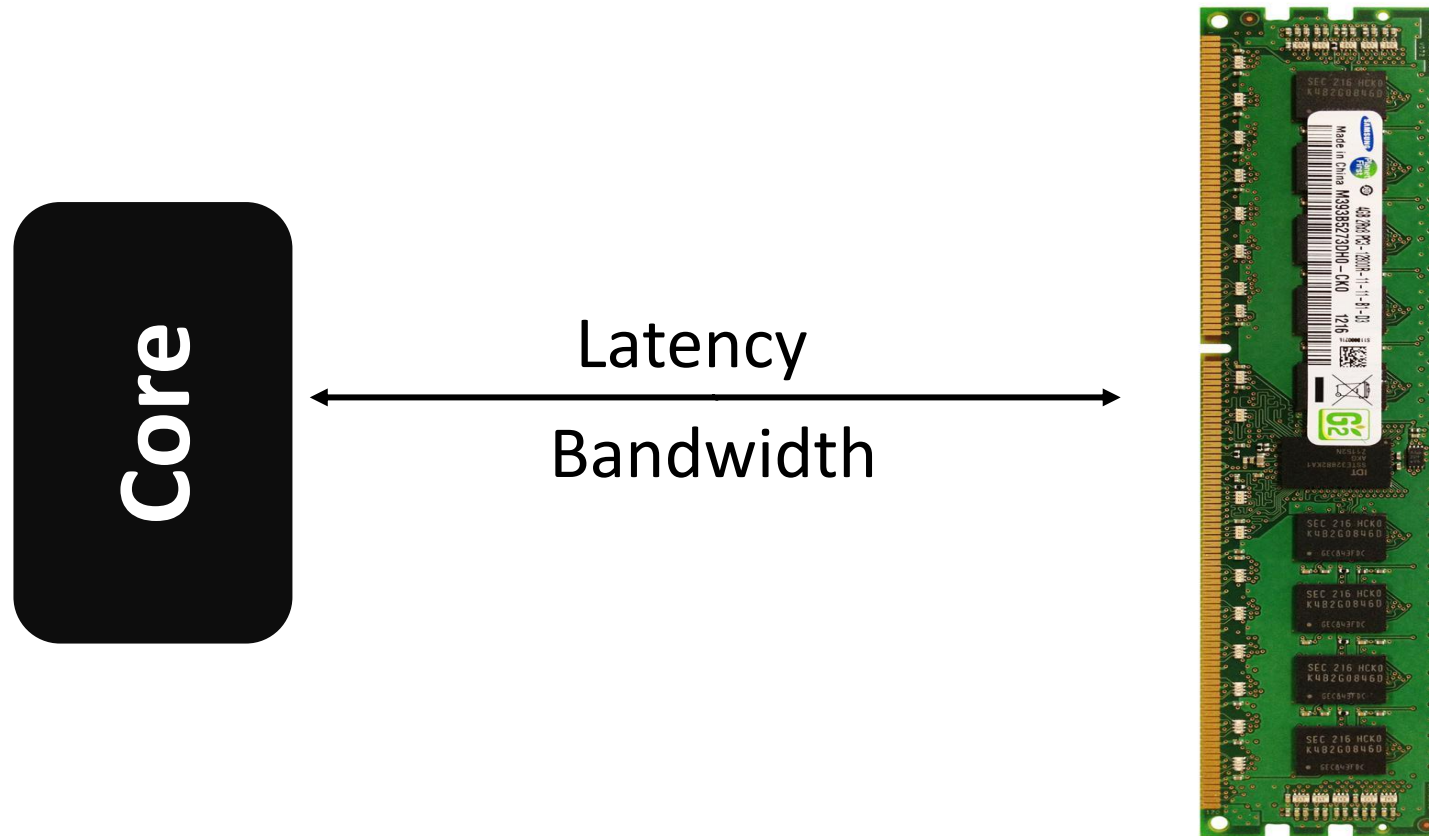


4 GB DRAM

*South pole ☹️*

CASPER

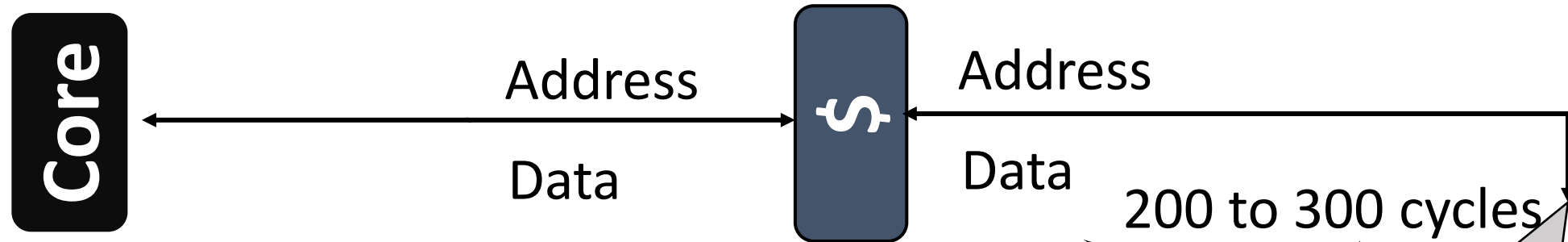
# Remember Latency and Bandwidth



CASPER

# Caching: 10K Feet View

North pole ☺



Caching is a *speculation* technique ☺  
Works – if locality

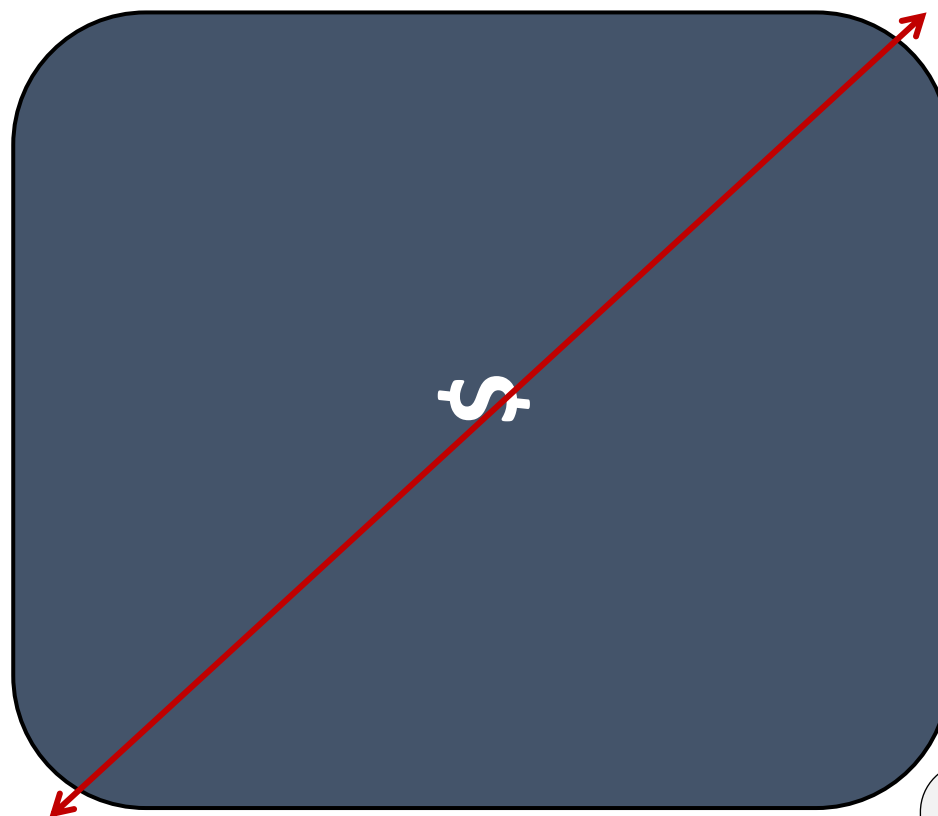


# How big/small?

Core



Latency: low  
Area: low  
Capacity: low

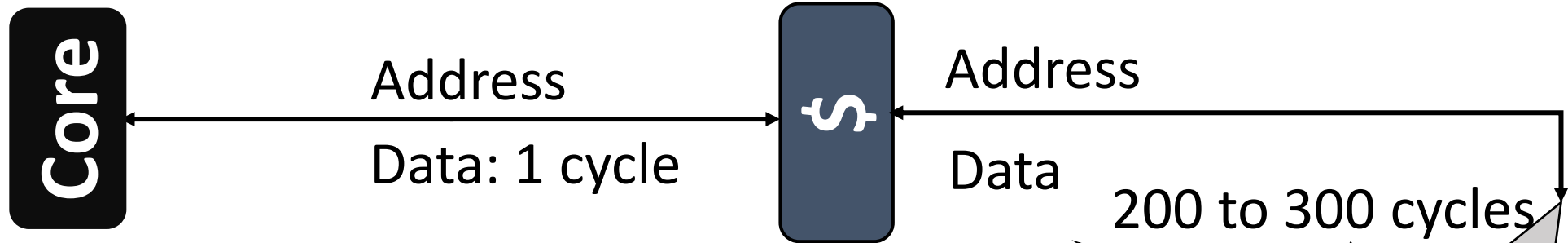


CASPER

Latency: high  
Area: high  
Capacity: high

# Cache with latency

*North pole ☺*



32 to 64KB \$ will be available in one to four cycles ☹

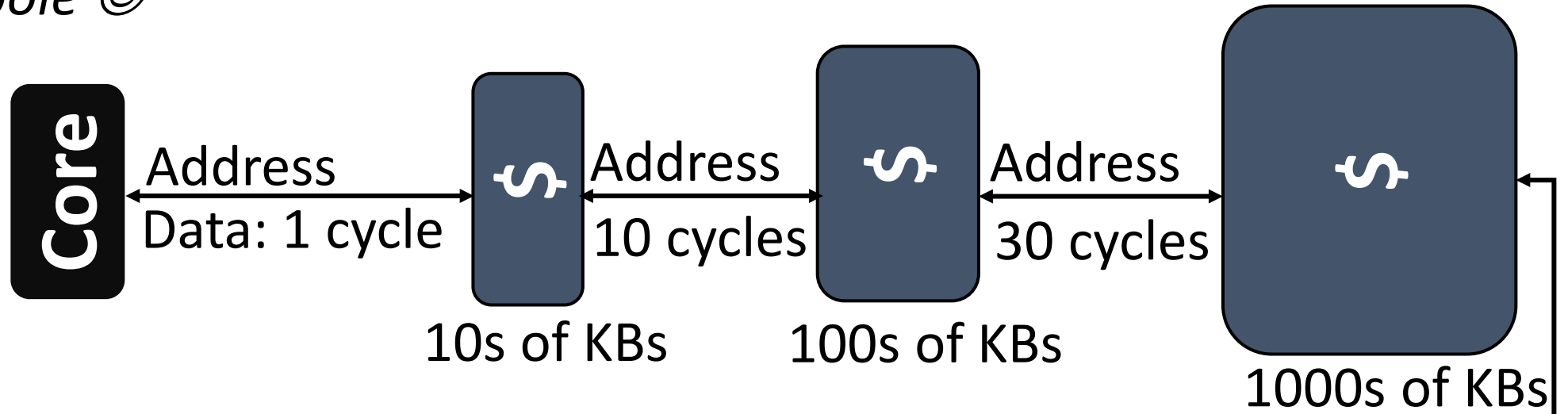
**Costly DRAM accesses ☹**



*South pole ☺*

# Cache hierarchy with latency

*North pole ☺*



Multi-level cache hierarchy



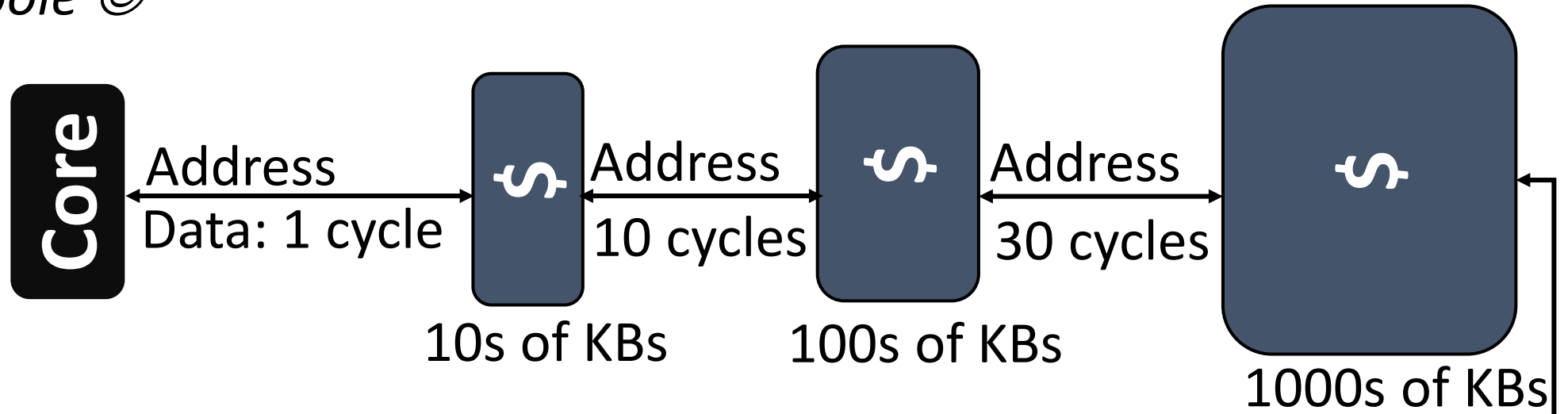
*South pole ☺*  
8

CASPER



# Cache hierarchy with latency

*North pole ☺*



Multi-level cache hierarchy

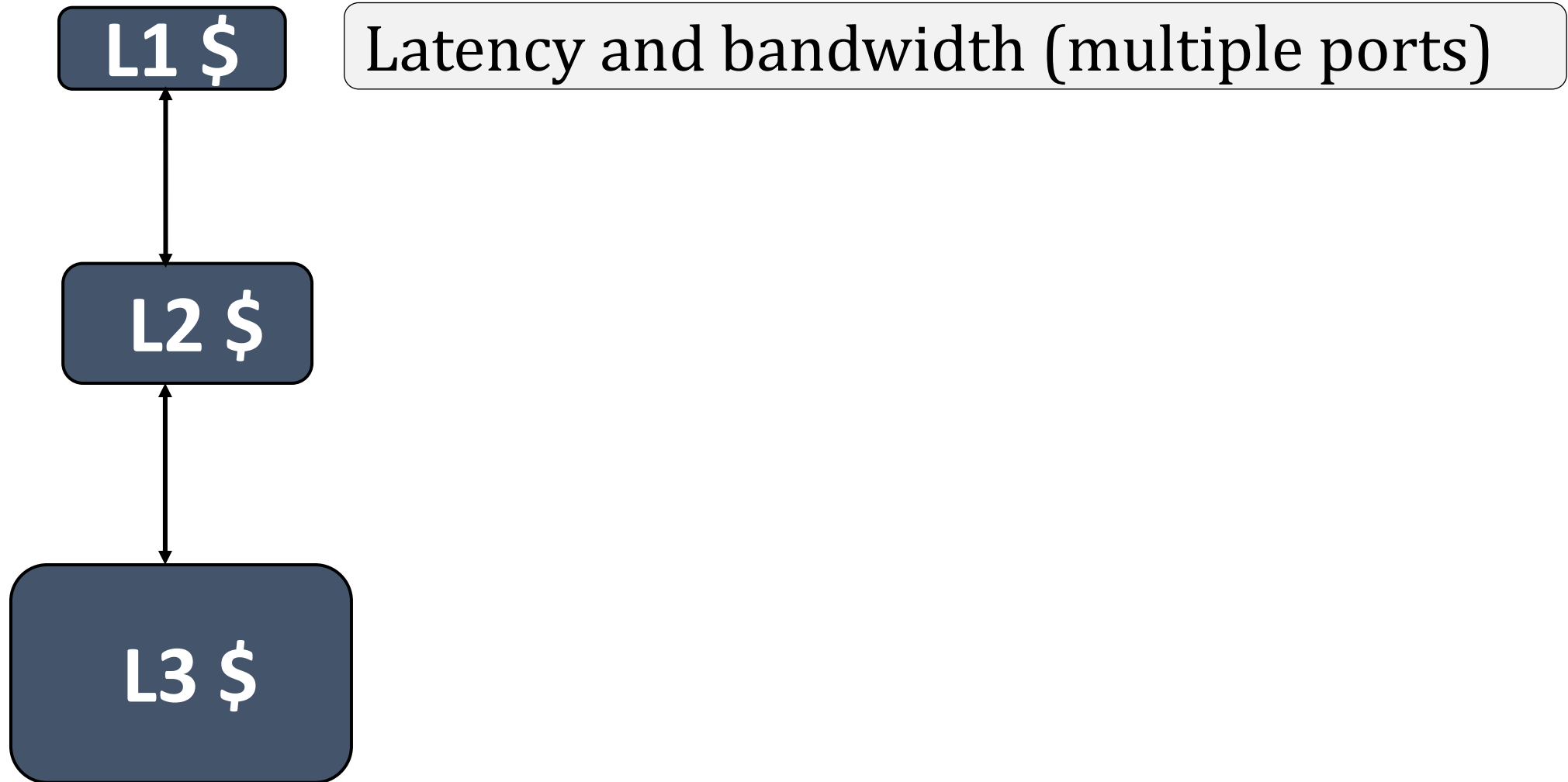
How many levels ?

Total latency < DRAM latency  
CASPER

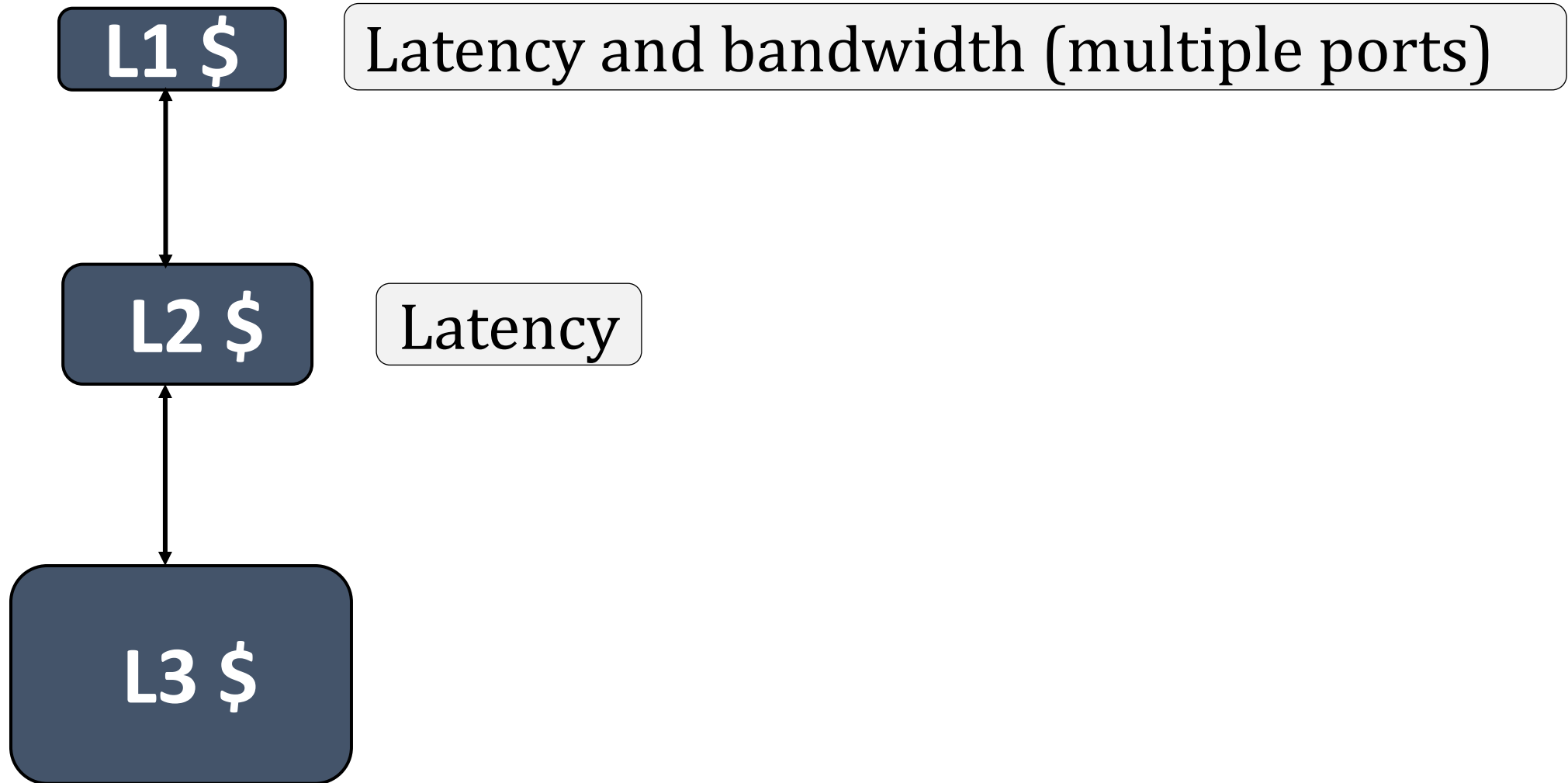


*South pole ☺*

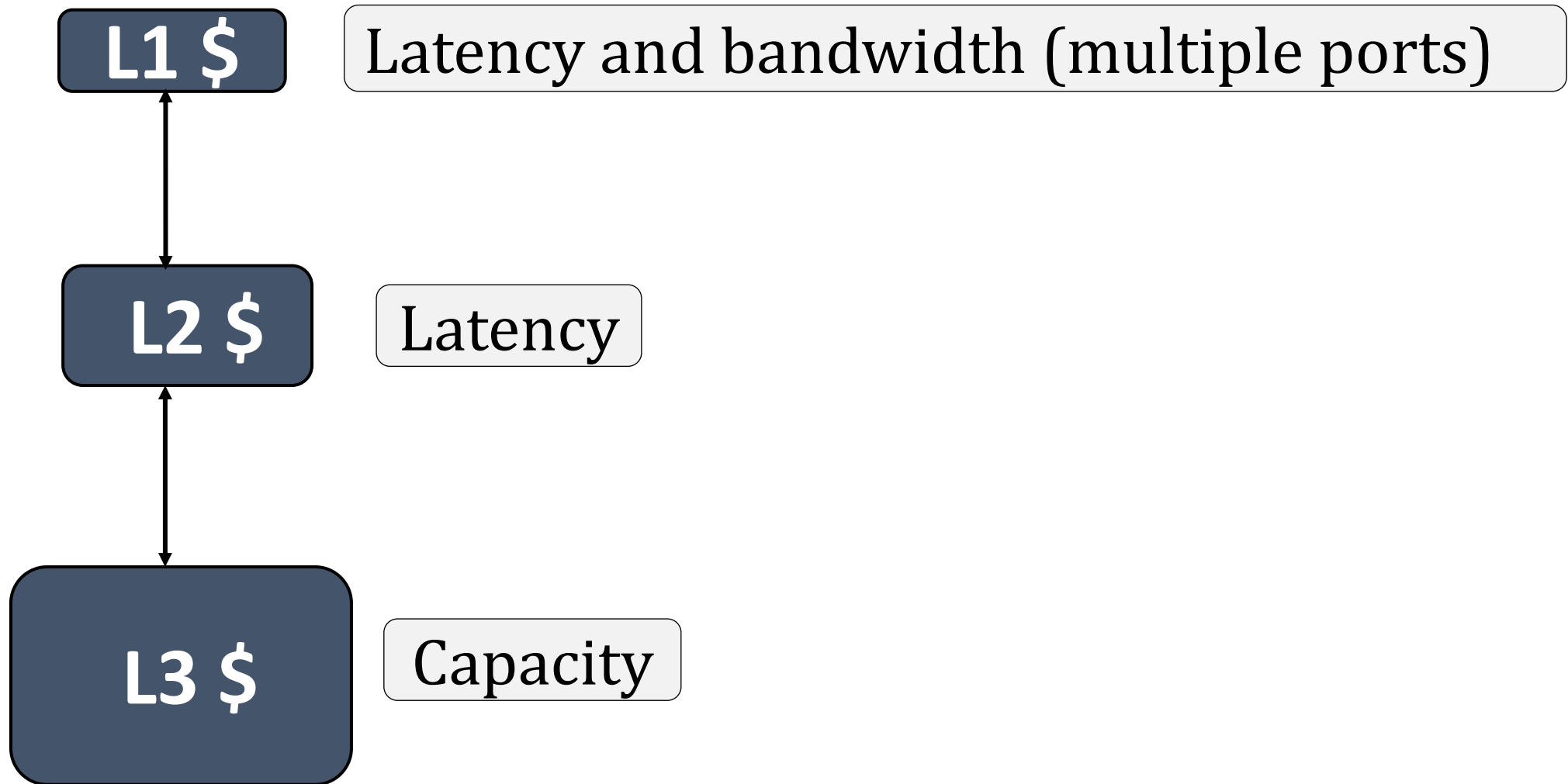
# Takeaway



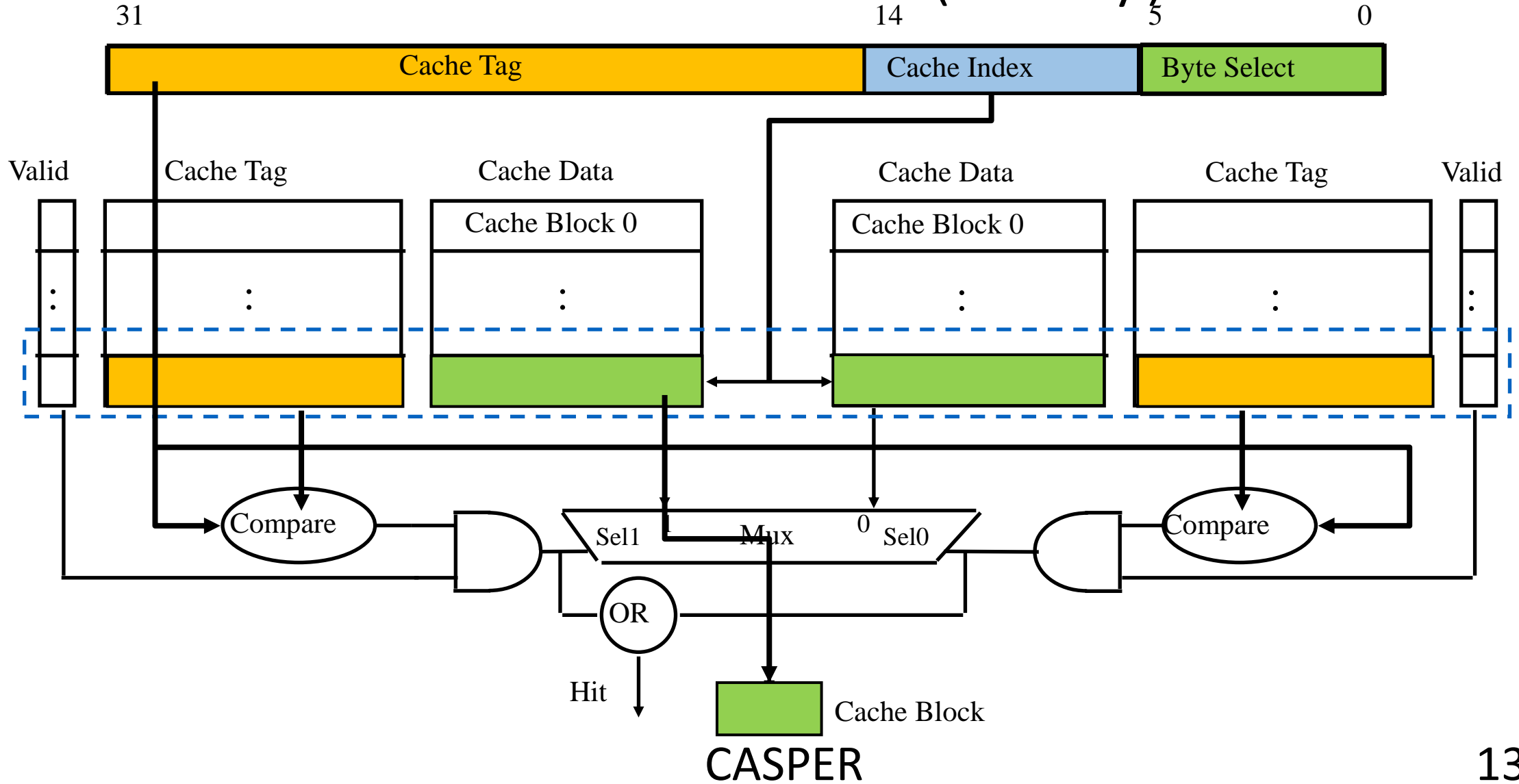
# Takeaway



Takeaway (Do not forget the word microarch.)



# Associative Cache in action (2-way)



# Knobs of interest

Line size, associativity, cache size

Tradeoff: latency, complexity, energy/power

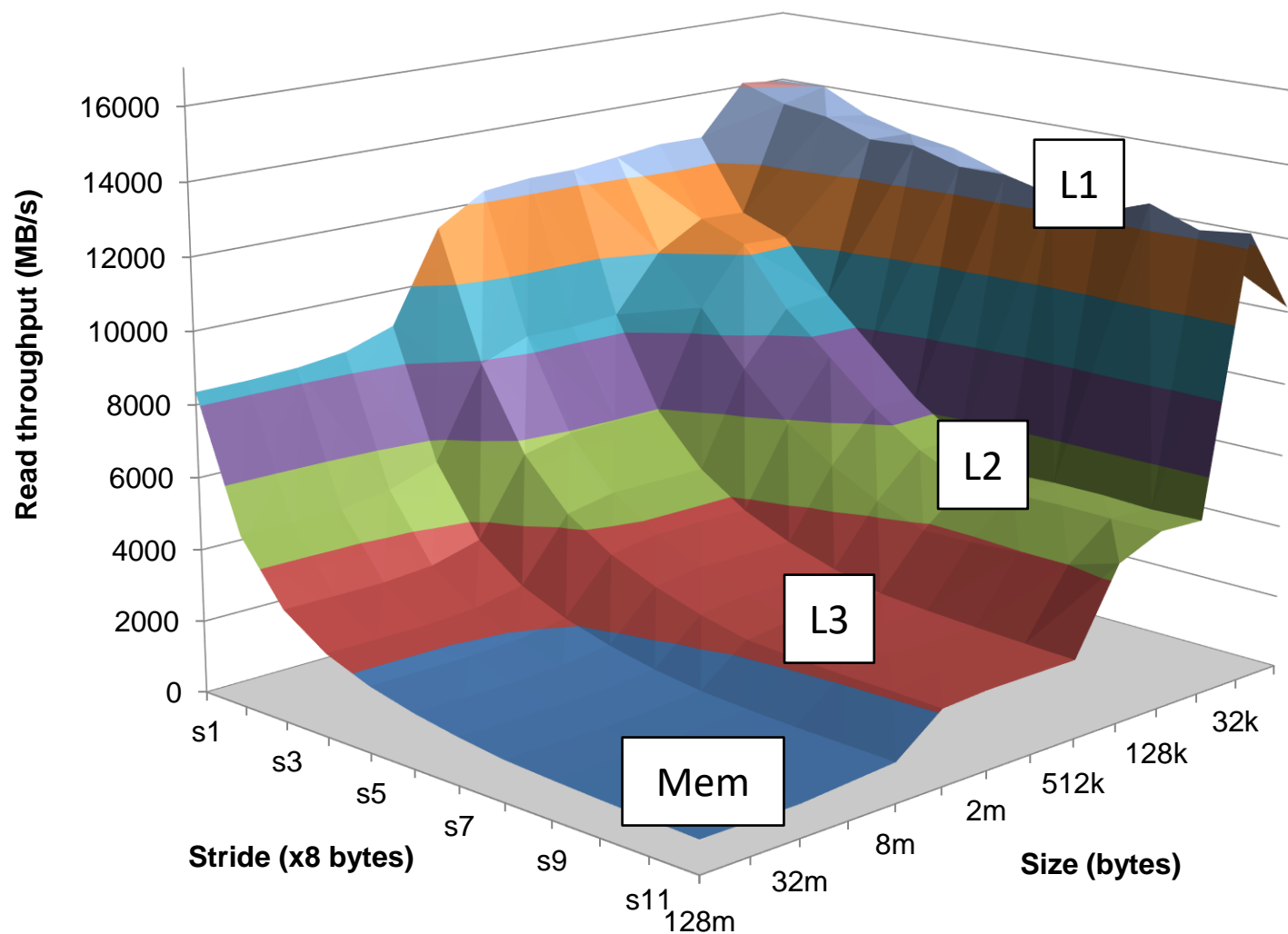
Tips: Think about the extremes:

**Line size** = one byte or cache size

**Associativity** = one or #lines

**Cache size** = Goal oriented: latency/bandwidth or capacity

# Memory Mountain



CASPER

# Into the Real World

sudo dmidecode -t cache

cat /proc/cpuinfo

getconf -a | grep CACHE

lscpu

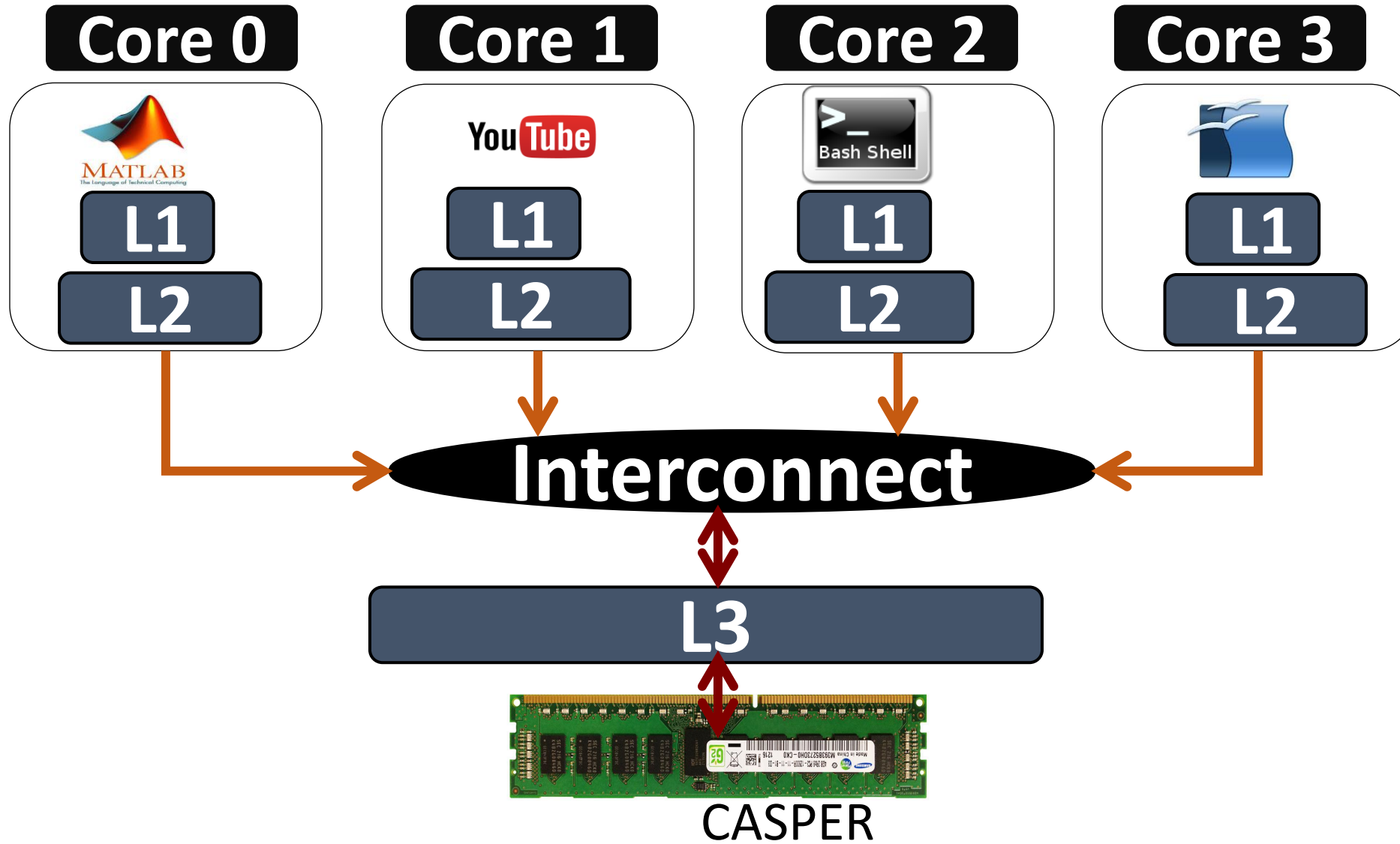
Wiki chip: <https://en.wikichip.org/wiki/WikiChip>

Perf tool: [https://perf.wiki.kernel.org/index.php/Main\\_Page](https://perf.wiki.kernel.org/index.php/Main_Page)

sudo perf stat -e cache-misses ....



# Multicore





PAUSE

CASPER

18

# Information leakage

**CIA: Confidentiality, Integrity,  
Availability**

CASPER

# Information leakage

**CIA: Confidentiality, Integrity,  
Availability**

**Confidentiality:** was data being computed upon not revealed to an un-permitted party?

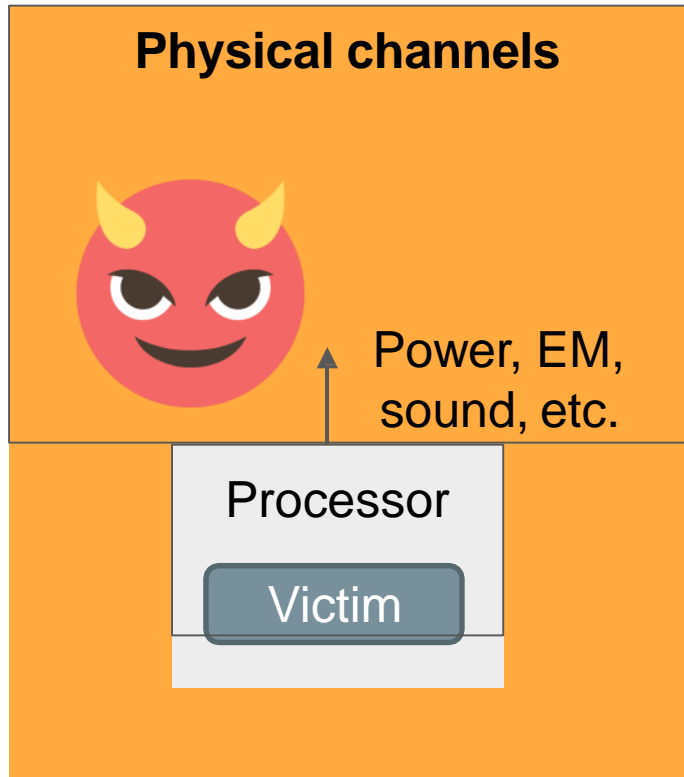
**Integrity:** was the computation performed correctly, returning the correct result?

**Availability:** did the computational resource carry out the task at all?

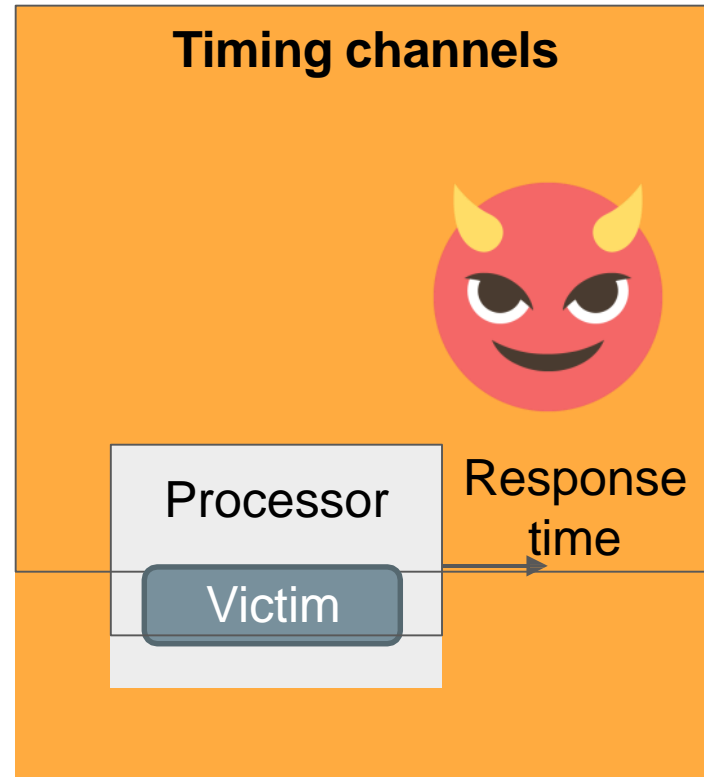


# Channels of Interest

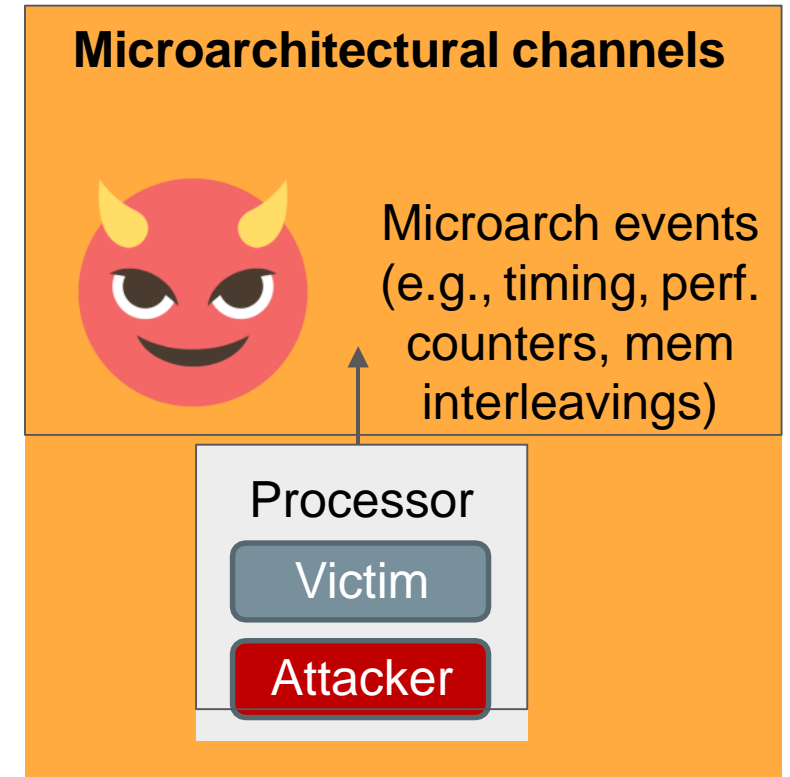
CS773



Attacker requires measurement equipment → physical access



Attacker may be remote (e.g., over an internet connection)



Attacker may be remote, or be co-located

# Side/Covert Channel

- A **side channel** is an *unintended* communication between two or more parties
- A **covert channel** is an *intended* communication between two or more parties (you upload a video on YouTube to communicate some information to your friends, if Gmail, whatsapp, call is not allowed)

In both cases:

- Communication *should not be possible*, following computing systems semantics
- The physical channel used for the communication can be the same

Side channels → unintended → need de-noising

Covert channels can show “best case” leakage

# Scope of these channels

- Inter-process(application) communication that can violate privilege boundaries
- Infer information from application's data-dependent HW resource usage

Side/covert channels not in any **interface specification (e.g. ISA)**.

Therefore stealthy

- Sophisticated mechanisms needed to detect channel
- No permanent indication one has been exploited

# Let's try to send a bit



Two processes can agree on “dead drops”

Cache:

**# ways**

**# sets**




# Let's try to send a bit

Two processes can agree on “dead drops”

Process 1  
(Sender)

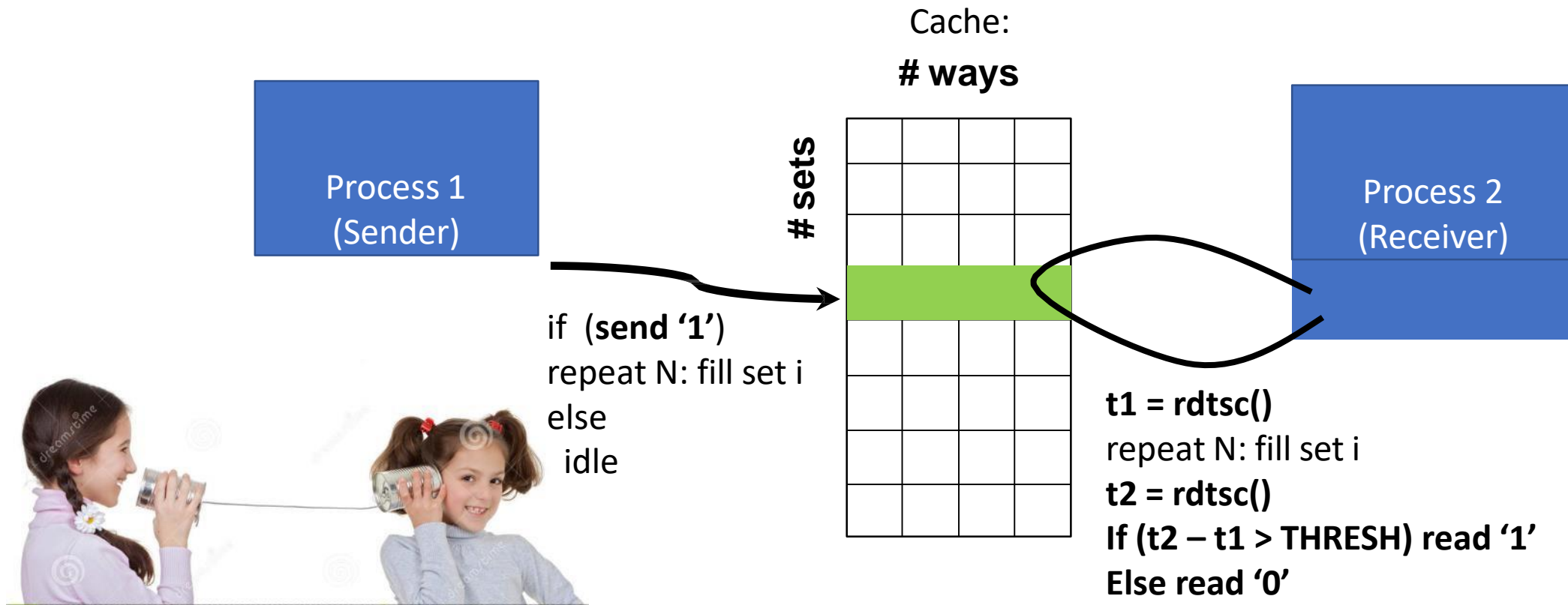
Cache:  
**# ways**

**# sets**


Process 2  
(Receiver)

# Let's try to send a bit

Two processes can agree on “dead drops”



# How is it different from legitimate communication

## Normal communication

```
include <socket.h>

void send(bit msg) {
    socket.send(msg);
}

bit recv() {
    return socket.recv(msg);
}
```

send(msg)



Channel



recv()

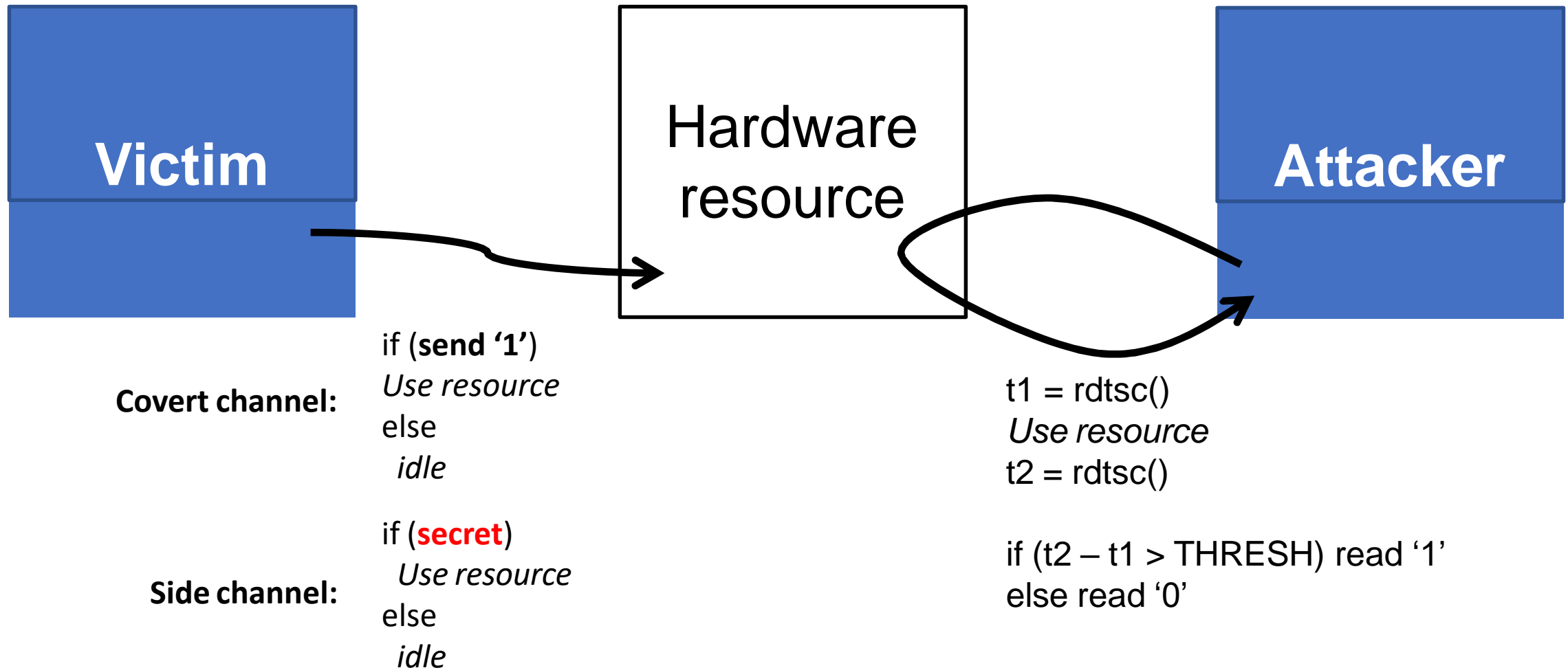


## Covert Channel communication

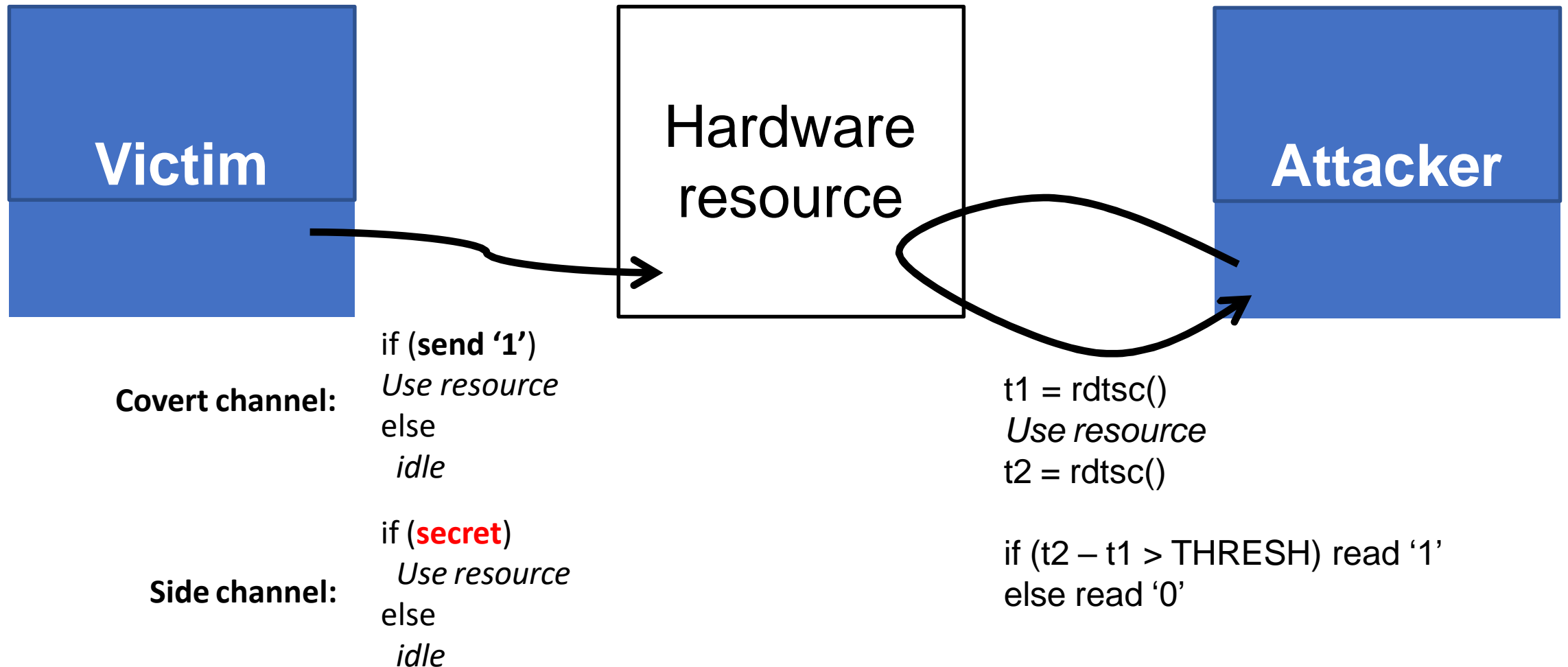
```
void send(bit msg) {
    // pressure on cache
}

bit recv() {
    st = time();
    // pressure on cache
    return time() - st > THRESH;
}
```

# From Covert to Side Channel

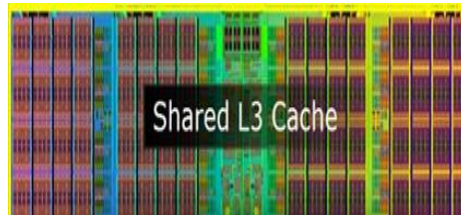
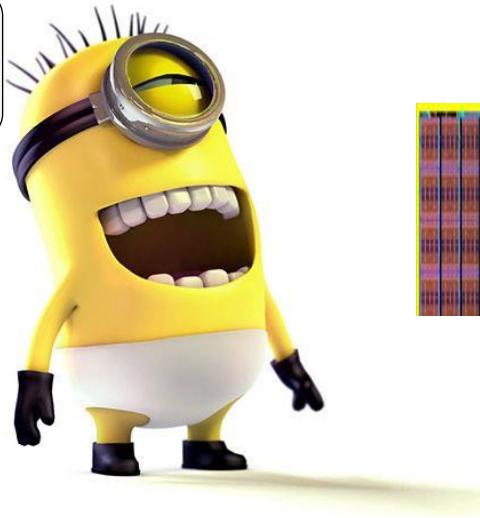


# From Covert to Side Channel



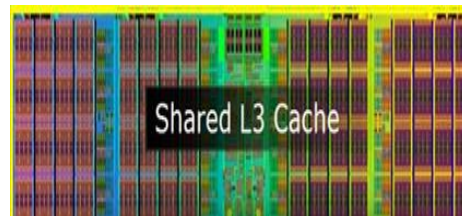
# Side/Covert Channel: Summary

Spy



Side-channel attacks

Let's  
play



Covert-channel attacks

Victim



Oh Yes!!

# Information leakage in the real world

$x \leftarrow 1$

Modular exponentiation,  $b^e \bmod n$

**for**  $i \leftarrow |e|-1$  **downto** 0 **do**

Exponent  $e$  is used for decryption

$x \leftarrow x^2 \bmod n$

*square*

**if** ( $e_i = 1$ ) **then**

*reduce*

$x = xb \bmod n$

**endif**

*multiply*

**done**

**return**  $x$

Attacker tries to get the  $e$

# Information leakage

$x \leftarrow 1$

Modular exponentiation,  $b^e \bmod n$

**for**  $i \leftarrow |e|-1$  **downto** 0 **do**

Exponent  $e$  is used for decryption

$x \leftarrow x^2 \bmod n$

*square*

**if** ( $e_i = 1$ ) **then**

*reduce*

$x = xb \bmod n$

**endif**

*multiply*

**done**

**return**  $x$

$e_i = 0$ , Square Reduce (SR)  
 $e_i = 1$ , SRMR

Attacker tries to get the  $e$



# Summary

- Latency
- Bandwidth
- Side/Covert Channels
- Bandwidth, Accuracy, Agility
- Attack detector
- More details next lecture