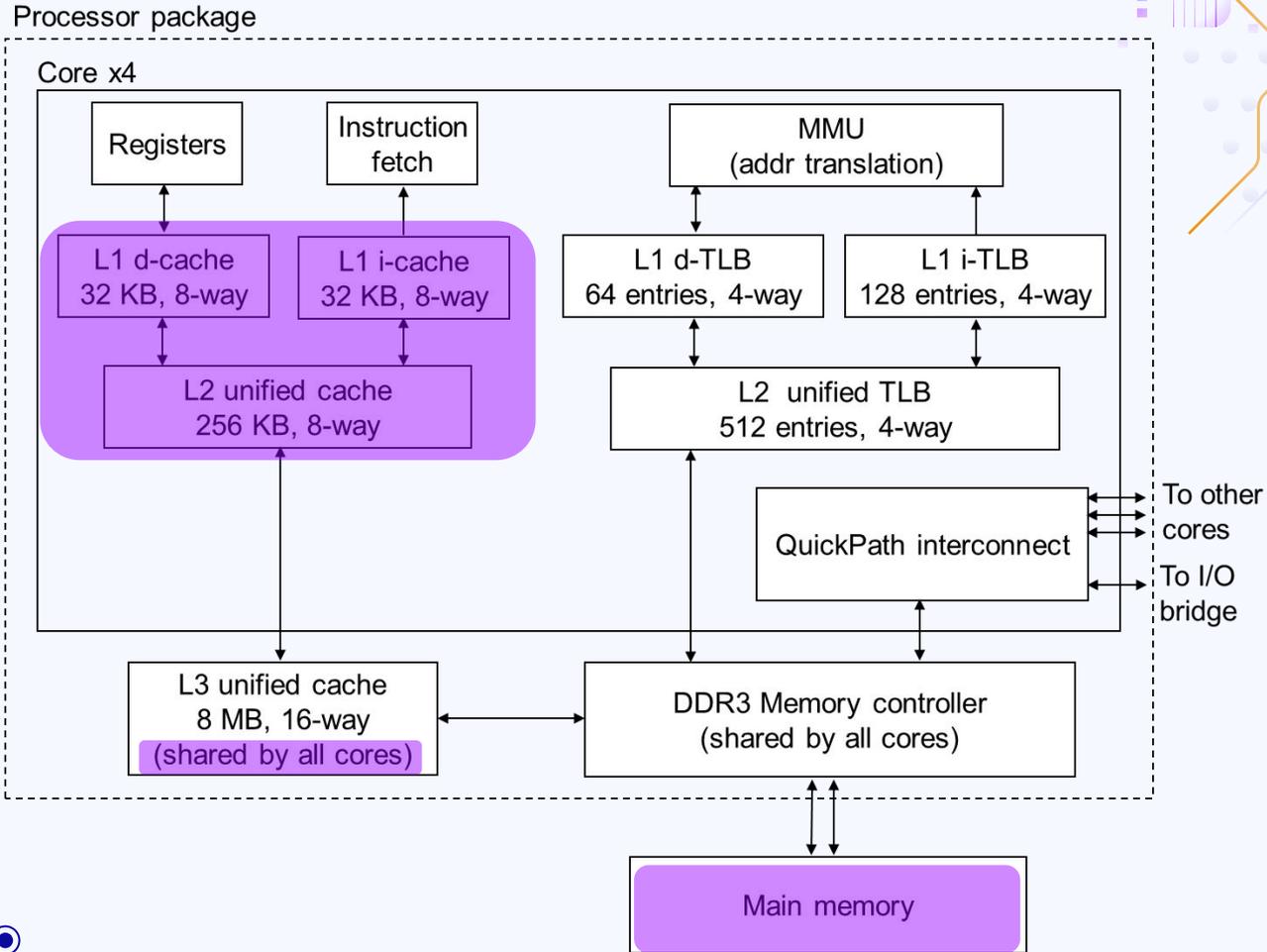


The perils of shared and unshared caches (coherence, side channels)

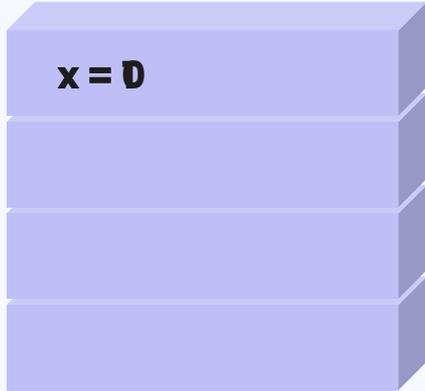
Intel i7

Source
(Bryant & O'Hallaron)

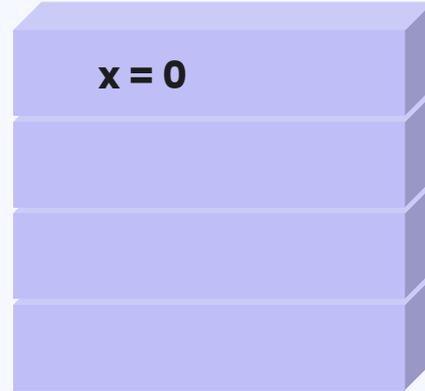


Cache coherence problem

Core A \$



Core B \$



Memory

$x = D$

Definitions

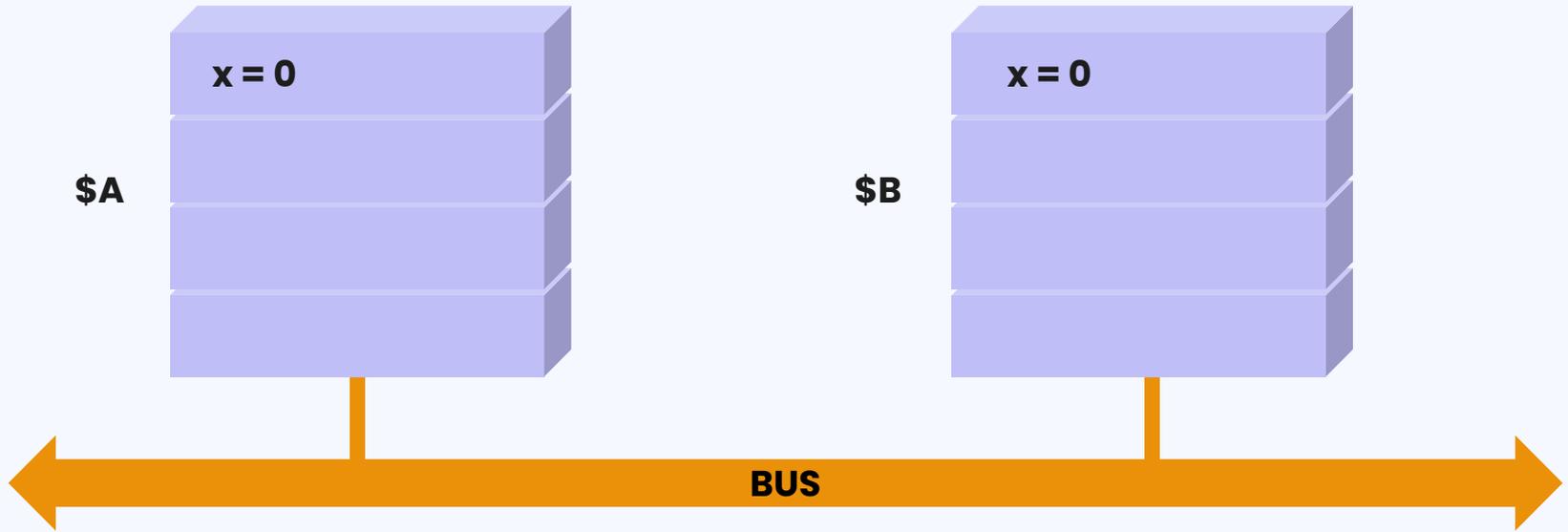
Intuitively: want any read of an item to return most recently written value to item

Coherence – what values can be returned by a read?

1. On a uniprocessor: reads after writes return written value
2. On a multiprocessor: reads by **B** after writes by **A** return written value **when given sufficient time**
3. Two writes to same location by one processor are seen in the same order by all processors
 - a. **A** writes 1 and then 2 to a memory location
 - b. We don't want **B** to decide the value is 1 and **C** to decide the value is 2

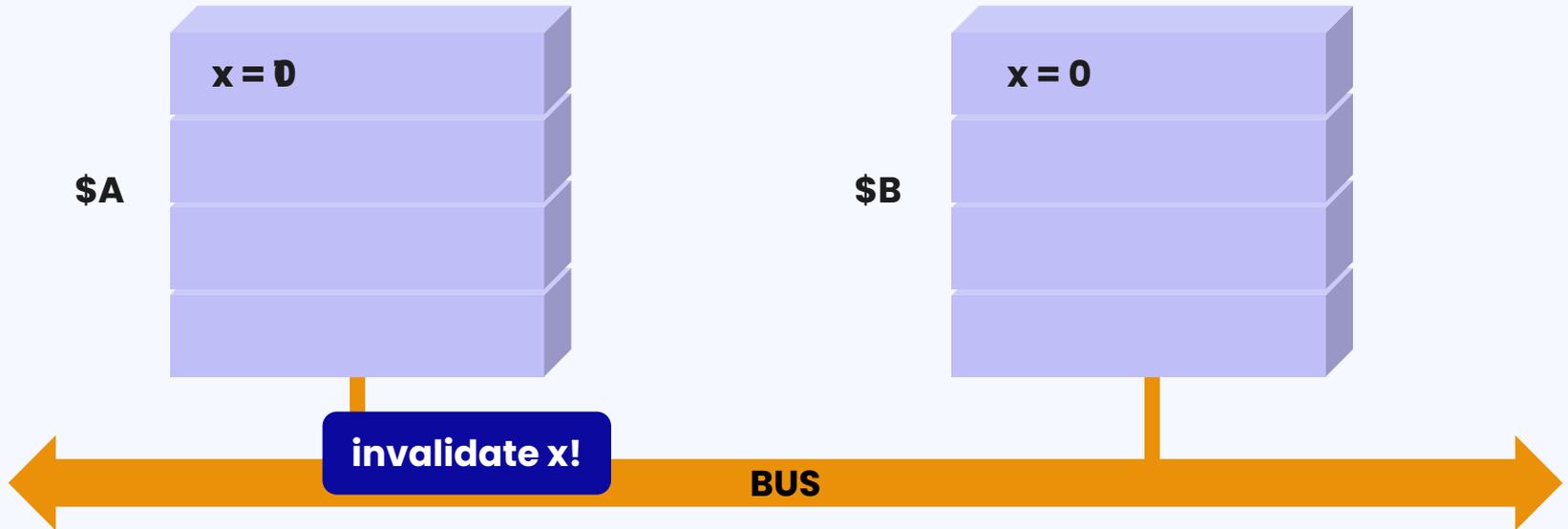
Snooping

1. A reads x
2. B reads x



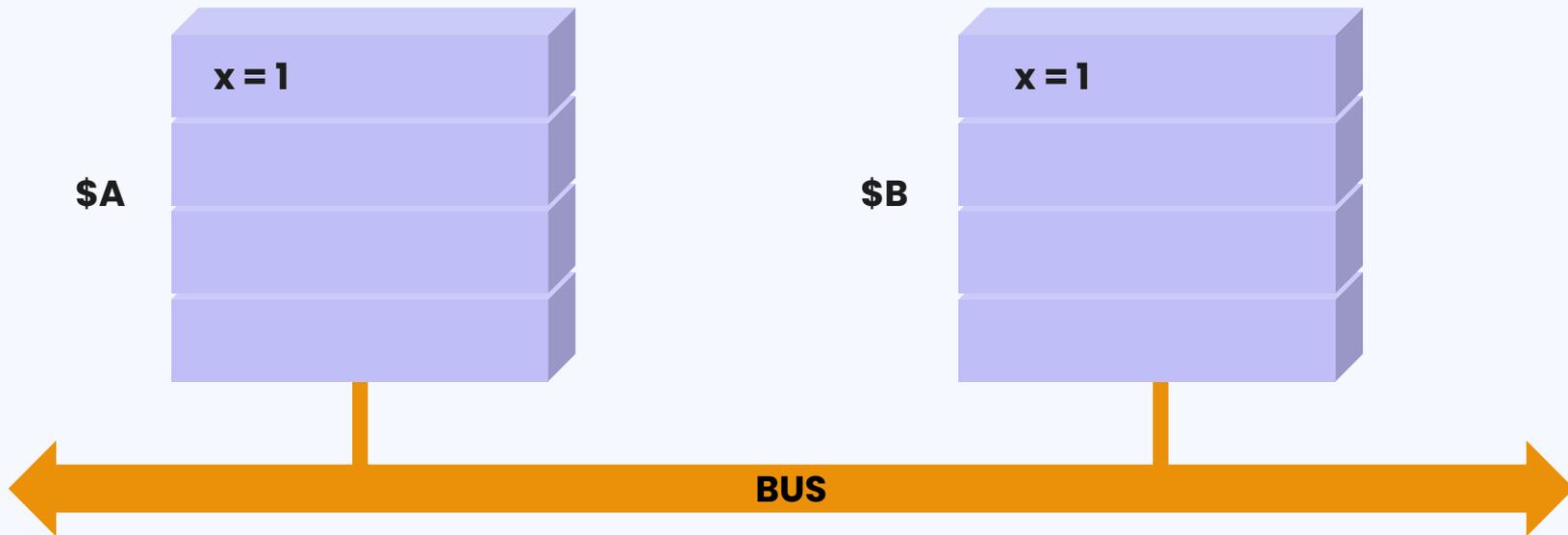
Snooping

3. A writes x



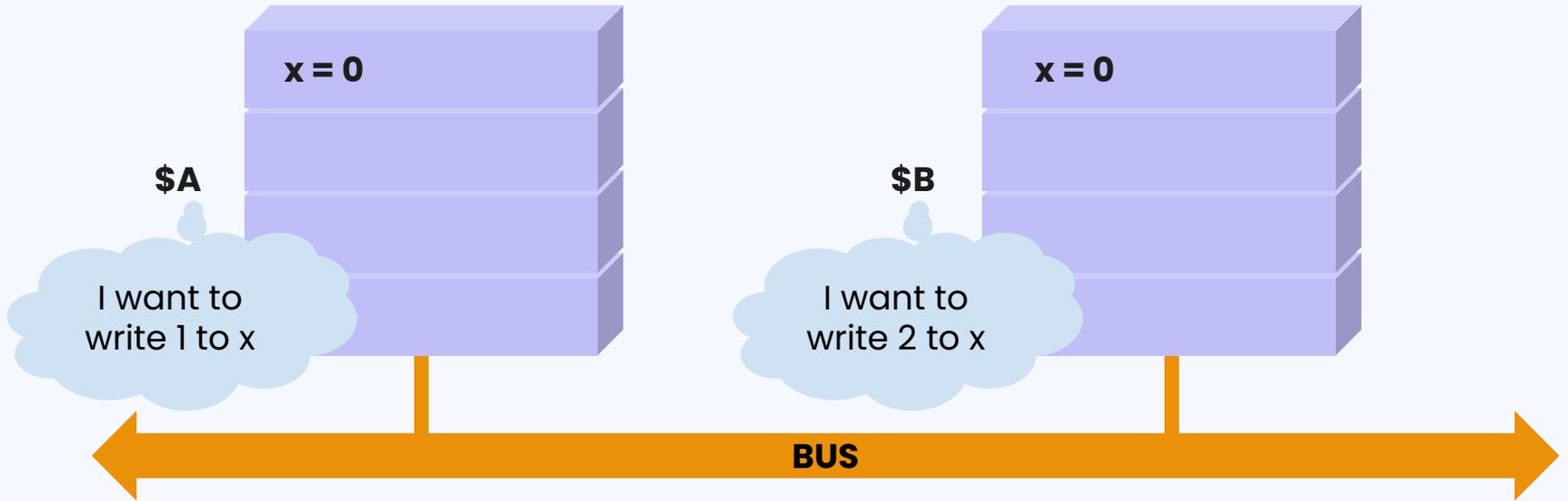
Snooping

4. B reads x



Cache needs to keep state for each block based on bus messages (common protocol: MESI (upgrade of MSI))

What should happen here?





What effects does block size have on
coherence protocols?



(Pivoting away from coherence)

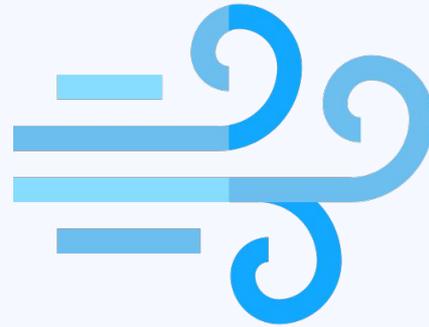
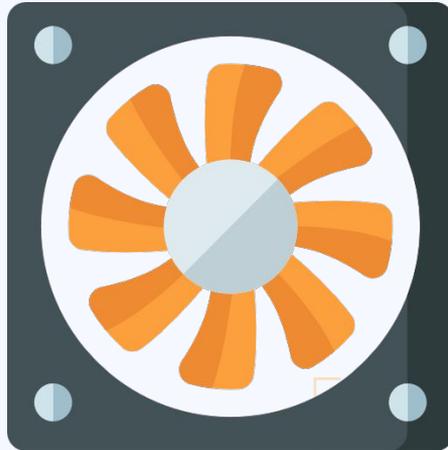
How can we *measurably* determine whether or not a piece of data is already present in a cache?

Security and side channels

Security: protection against threats from malicious actor

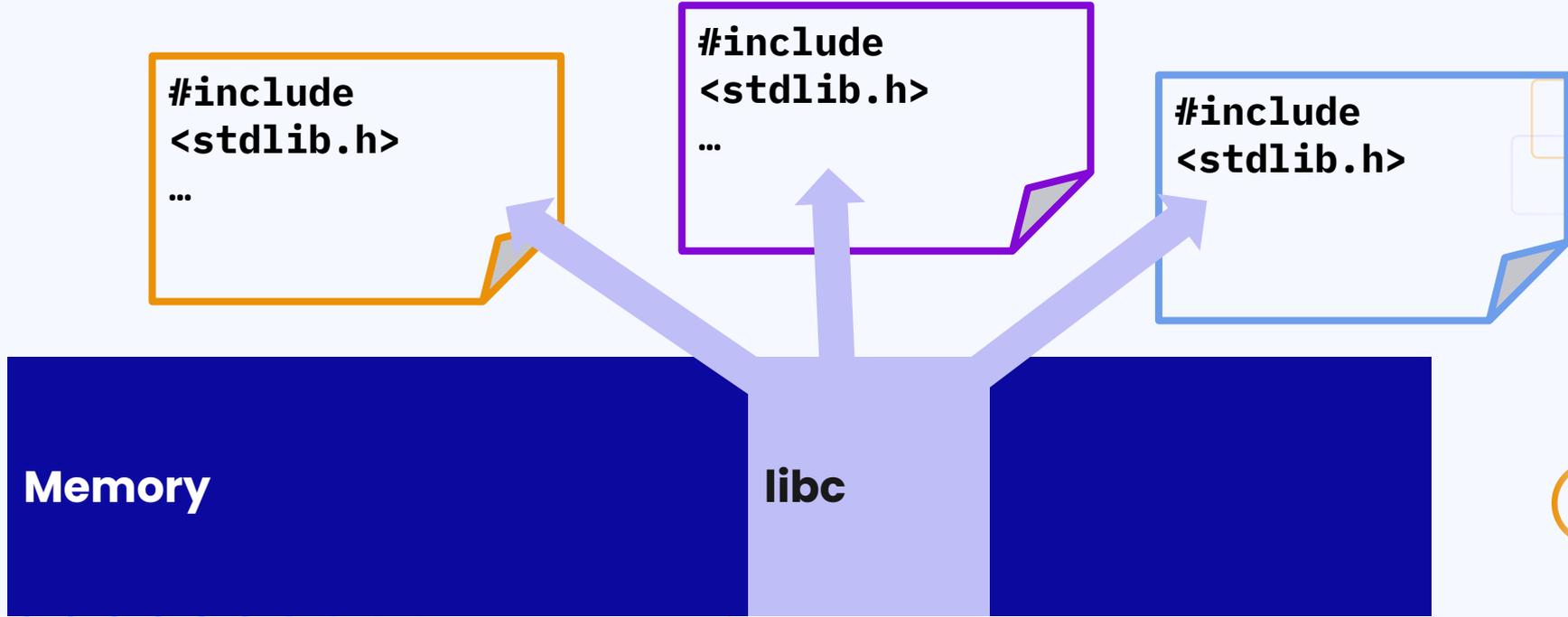
Obviously a **large** field, this is just a taste

Side channels: Incidental information leakage inferred from observing normal execution

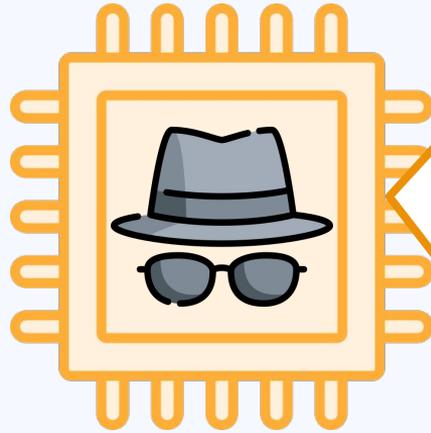
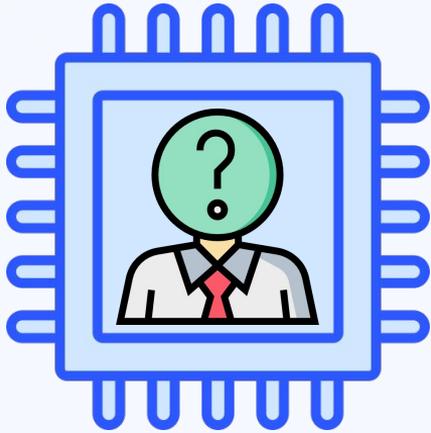


Shared cache side channels

What precisely can CPU B learn about CPU A when data is shared?



Flush & Reload attack



```
// flush the block  
asm(clflush 0xLIBCADDR);
```

```
// wait some time  
t1 = time.now();  
while(time.now() - t1 <  
100ns);
```

```
// access block  
t2 = time.now();  
x = * 0xLIBCADDR;  
accessTime = time.now()  
- t2;
```

```
// if slow: unused  
// if fast: used!
```

Shared \$



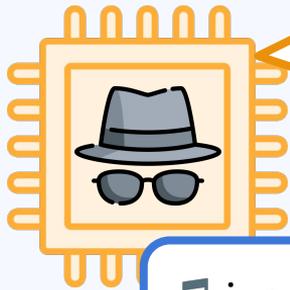
0xLIBCADDR



What does flush & reload depend on?

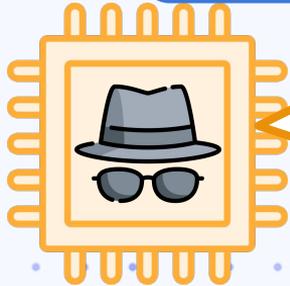
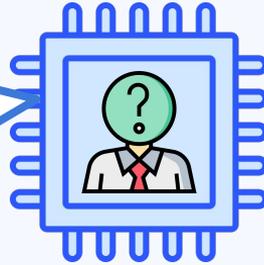
Prime & Probe

Doesn't require precise timing or access to victim's memory; measures cache contention (less granular)

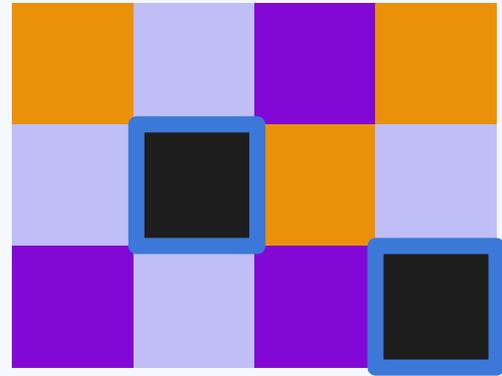


Fill up the cache!

♪ just executin' some code ♪



Read all addresses
Which ones were misses?



Caveats

It takes a *lot* to instrument a side-channel attack

Often can't learn *everything*, but narrow down search space

What can be done to guard against attacks?

- Oblivious RAM (largely theoretical)

- Cache partitioning/slicing

- Use an abacus