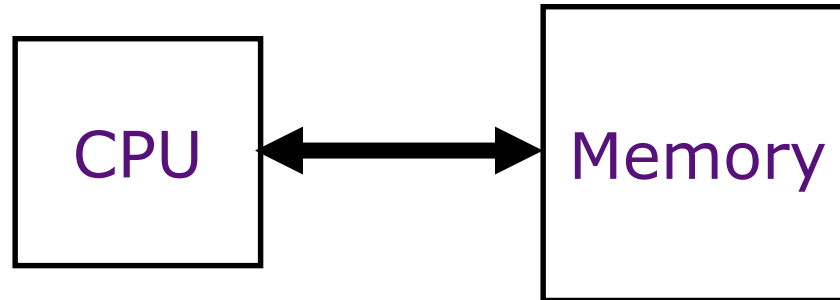


Cache Organization

Daniel Sanchez

Computer Science and Artificial Intelligence Laboratory
M.I.T.

CPU-Memory Bottleneck



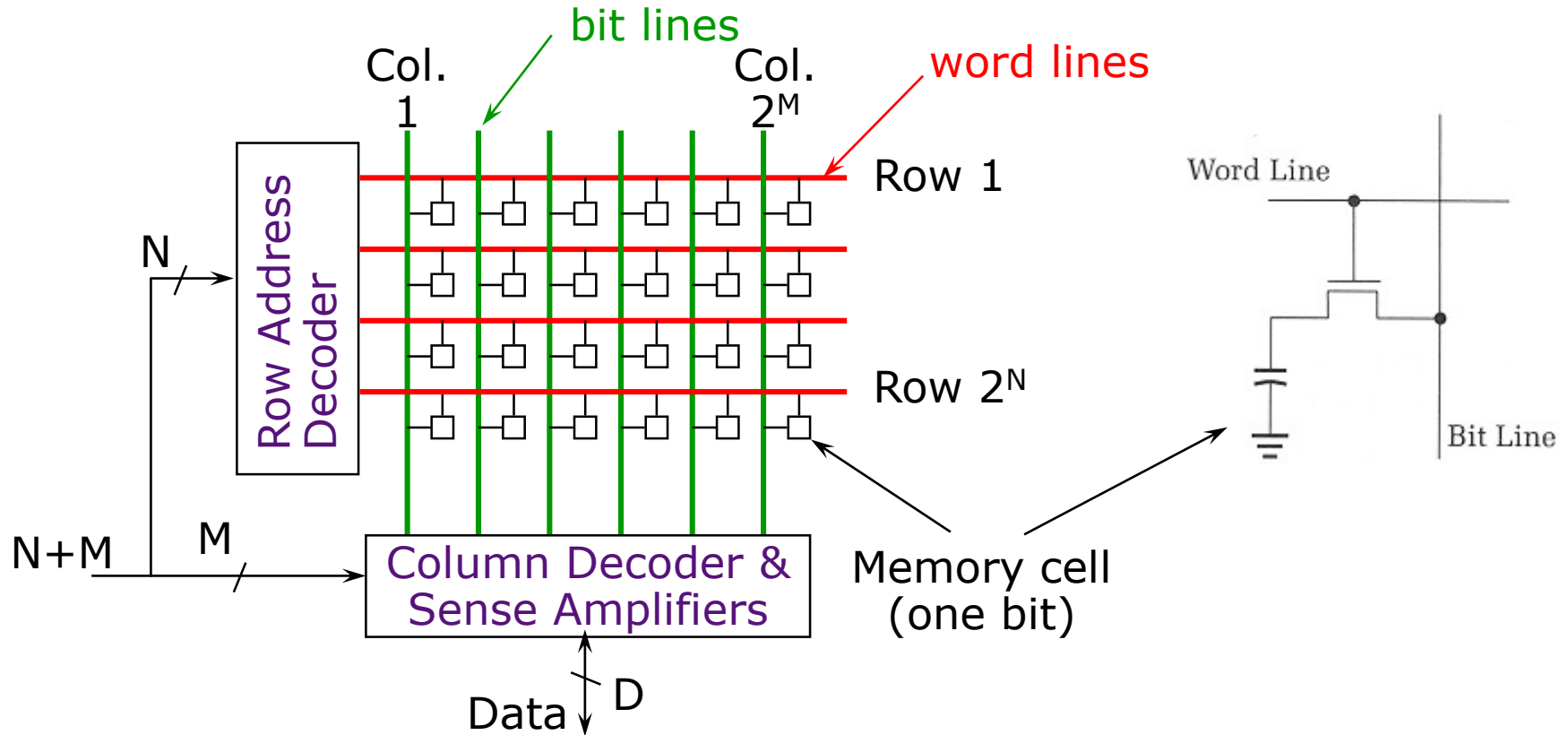
Performance of high-speed computers is usually limited by memory *bandwidth* & *latency*

- Latency (time for a single access)
Memory access time \gg Processor cycle time
- Bandwidth (number of accesses per unit time)
if fraction m of instructions access memory,
 $\Rightarrow 1+m$ memory references / instruction
 $\Rightarrow \text{CPI} = 1$ requires $1+m$ memory refs / cycle

Memory Technology

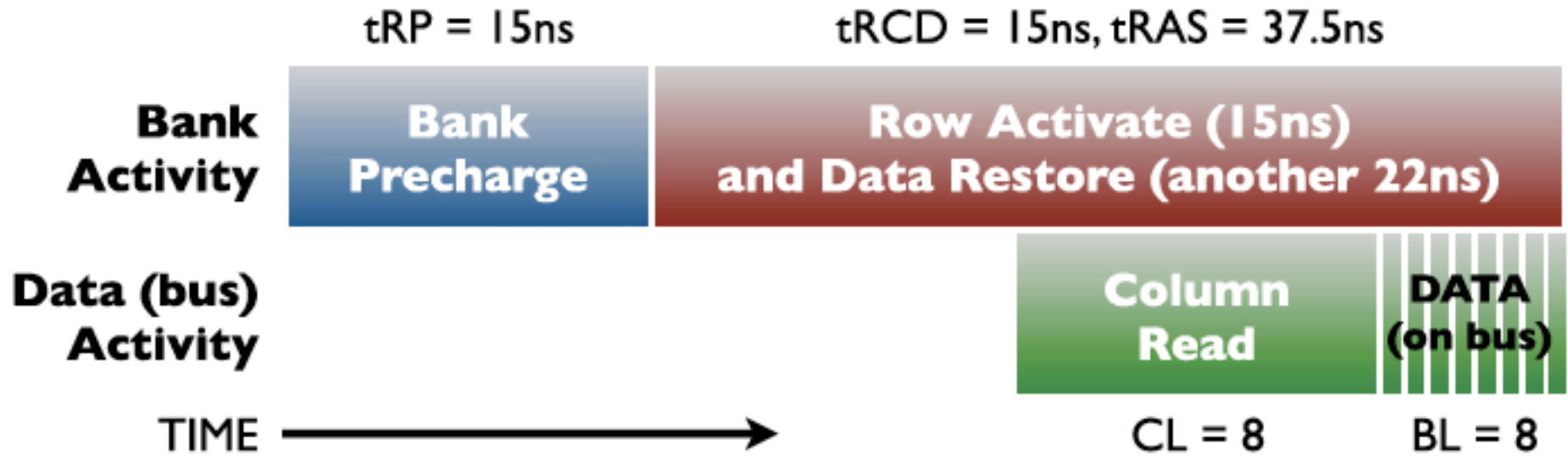
- Early machines used a variety of memory technologies
 - Manchester Mark I used CRT Memory Storage
 - EDVAC used a mercury delay line
- Core memory was first large scale reliable main memory
 - Invented by Forrester in late 40s at MIT for Whirlwind project
 - Bits stored as magnetization polarity on small ferrite cores threaded onto 2 dimensional grid of wires
- First commercial DRAM was Intel 1103
 - 1Kbit of storage on single chip
 - charge on a capacitor used to hold value
- Semiconductor memory quickly replaced core in 1970s
 - Intel formed to exploit market for semiconductor memory
- Flash memory
 - Slower, but denser than DRAM. Also non-volatile, but with wearout issues
- Phase change memory (PCM, 3D XPoint)
 - Slightly slower, but much denser than DRAM and non-volatile

DRAM Architecture



- Bits stored in 2-dimensional arrays on chip
- Modern chips have around 8 logical banks on each chip
 - Each logical bank physically implemented as many smaller arrays

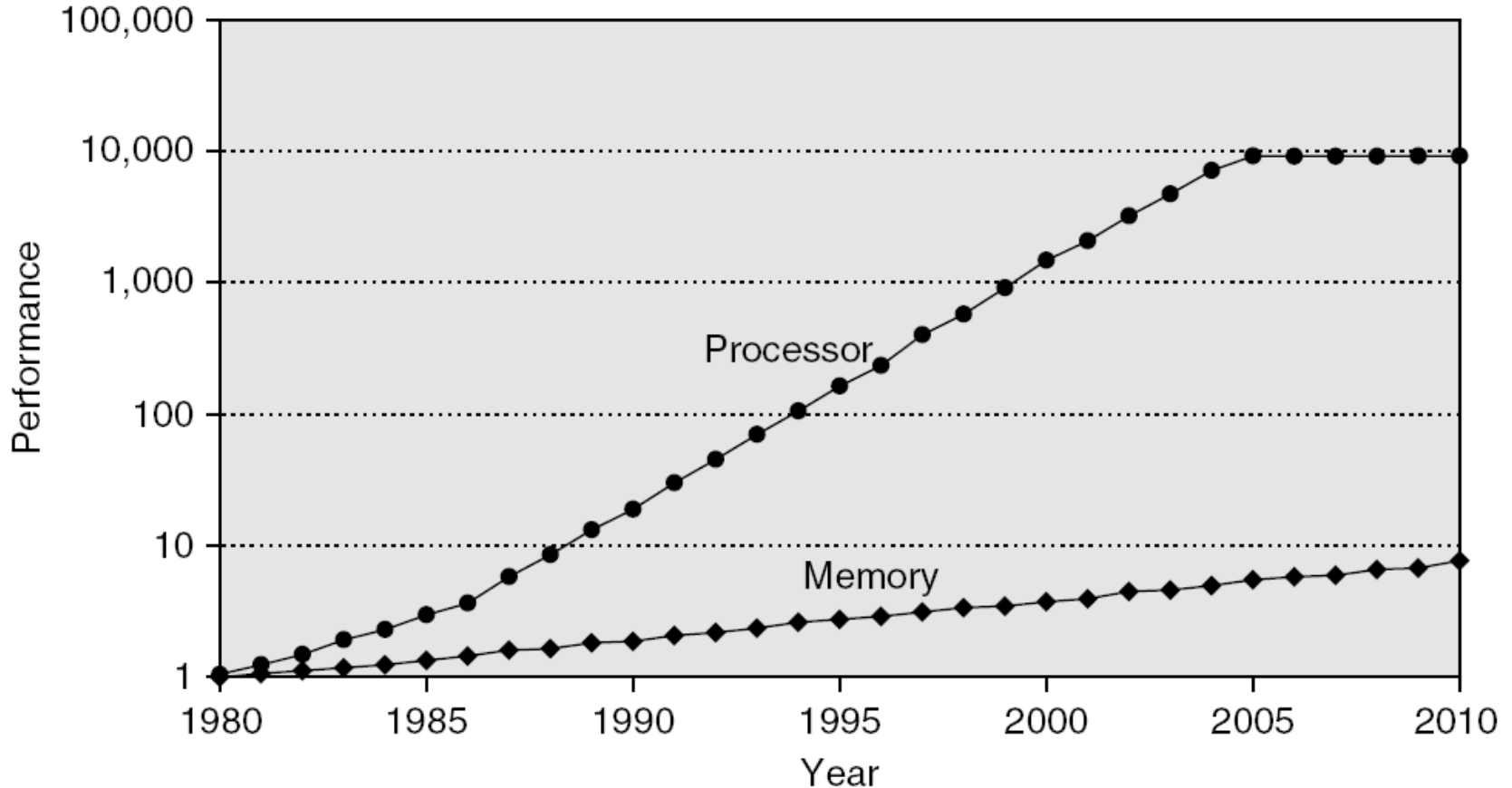
DRAM timing



DRAM Spec:

CL, tRCD, tRP, tRAS, e.g., 9-9-9-24

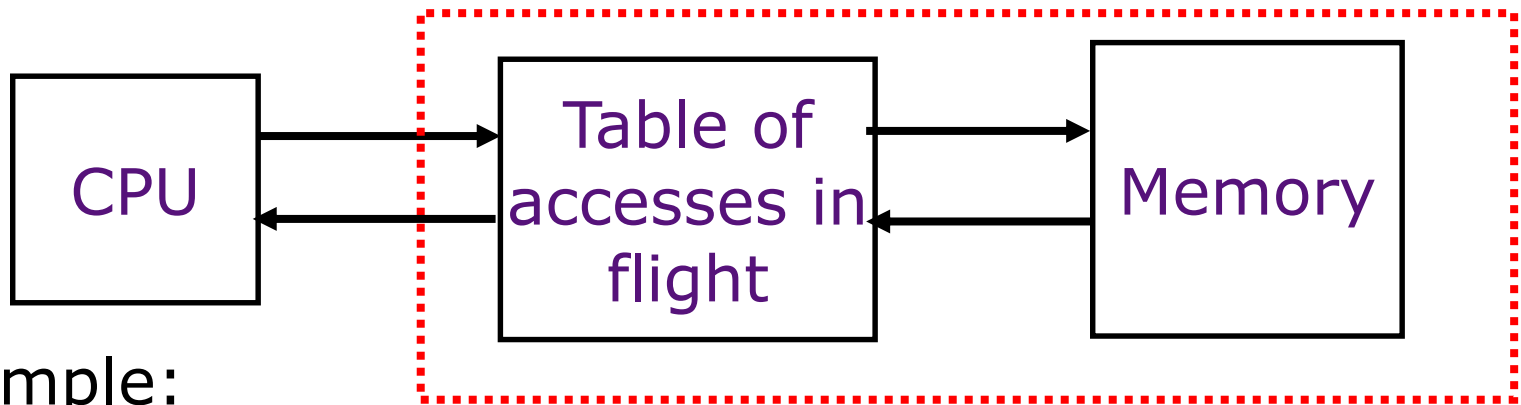
Processor-DRAM Gap (latency)



Four-issue 2GHz superscalar accessing 100ns DRAM could execute 800 instructions during time for one memory access!

Little's Law

Throughput (T) = Number in Flight (N) / Latency (L)



Example:

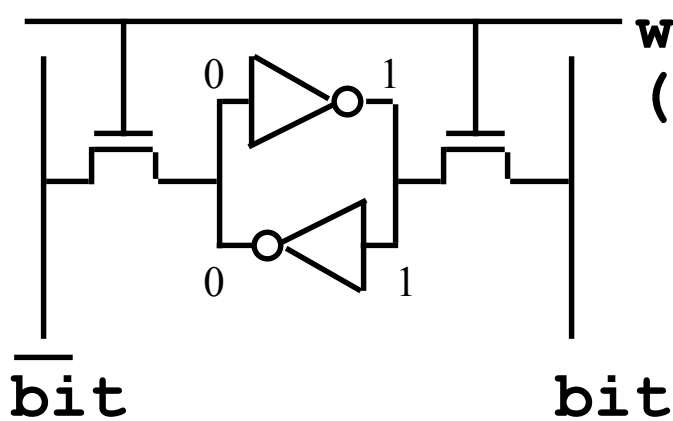
- Assume infinite-bandwidth memory*
- 100 cycles / memory reference*
- 1 + 0.2 memory references / instruction*

*⇒ Table size = 1.2 * 100 = 120 entries*

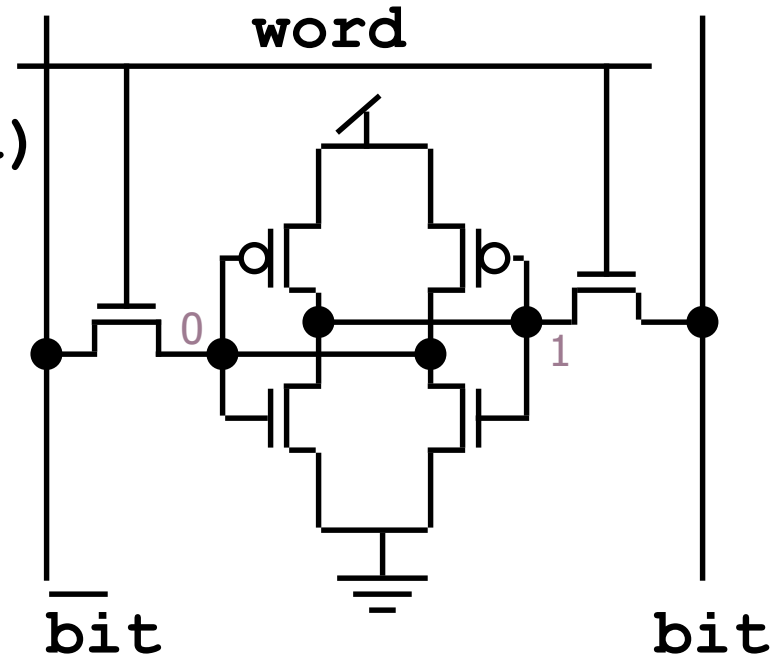
120 independent memory operations in flight!

Basic Static RAM Cell

6-Transistor SRAM Cell



- Write:
 1. Drive bit lines ($\text{bit}=1, \overline{\text{bit}}=0$)
 2. Select word line
- Read:
 1. Precharge bit and $\overline{\text{bit}}$ to Vdd
 2. Select word line
 3. Cell pulls one bit line low
 4. Column sense amp detects difference between bit & $\overline{\text{bit}}$



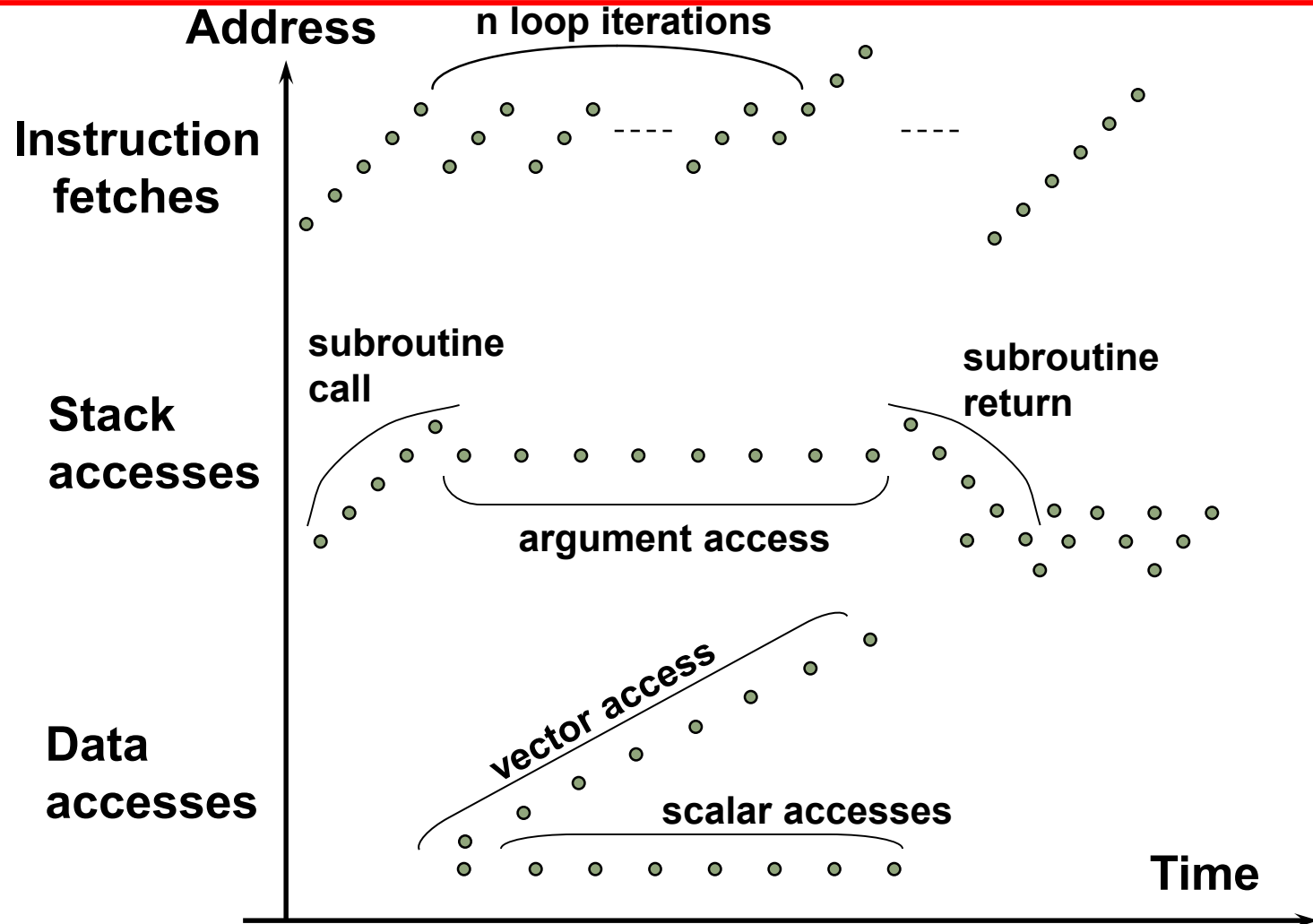
Multilevel Memory

Strategy: Reduce average latency using small, fast memories called caches.

Caches are a mechanism to reduce memory latency based on the empirical observation that the patterns of memory references made by a processor are often highly predictable:

	<u>PC</u>
...	96
<i>Loop:</i> <i>add r2, r1, r1</i>	100
<i>subi r3, r3, #1</i>	104
<i>bnez r3, loop</i>	108
...	112

Typical Memory Reference Patterns

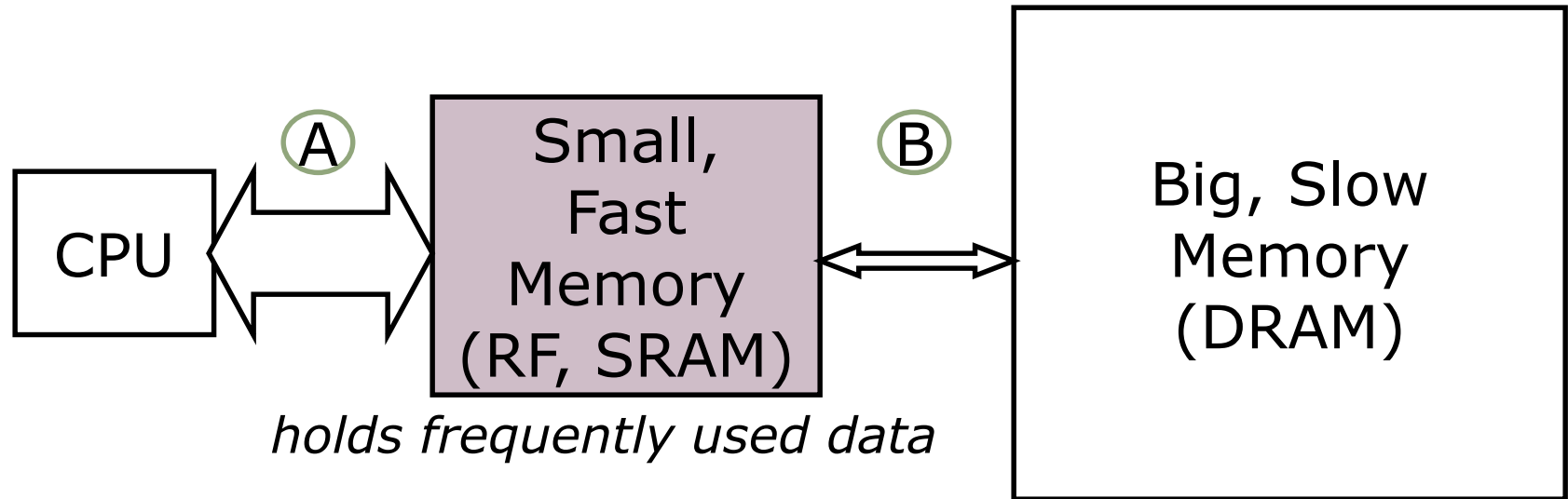


Common Predictable Patterns

Two predictable properties of memory references:

- *Temporal Locality*: If a location is referenced, it is likely to be referenced again in the near future
- *Spatial Locality*: If a location is referenced, it is likely that locations near it will be referenced in the near future

Memory Hierarchy



- *size:* Register \ll SRAM \ll DRAM *why?*
- *latency:* Register \ll SRAM \ll DRAM *why?*
- *bandwidth:* on-chip \gg off-chip *why?*

On a data access:

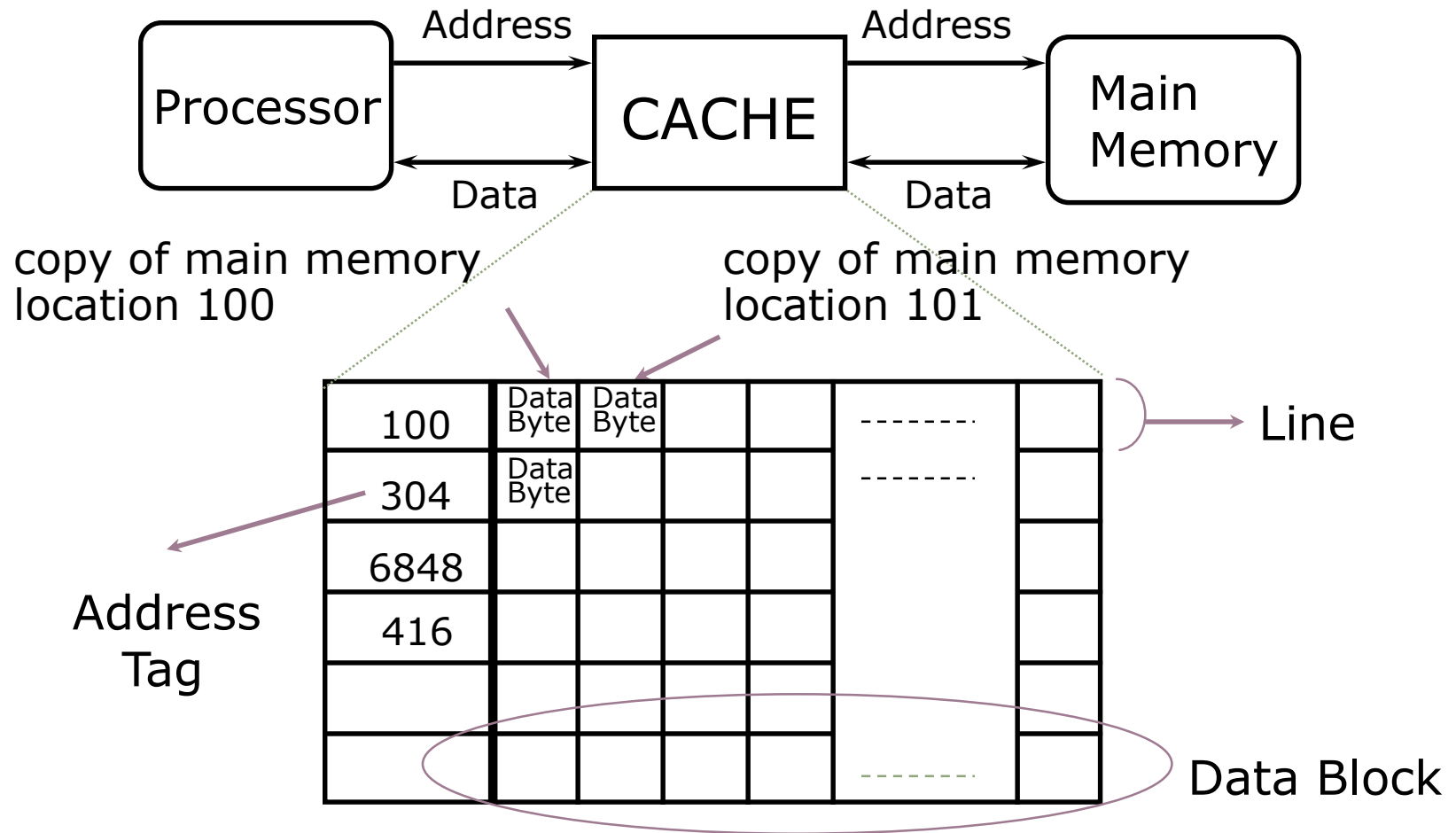
hit (data \in fast memory) \Rightarrow low latency access

miss (data \notin fast memory) \Rightarrow long latency access (*DRAM*)

Management of Memory Hierarchy

- Small/fast storage, e.g., registers
 - Address usually specified in instruction
 - Generally implemented directly as a register file
 - but hardware might do things behind software's back, e.g., stack management, register renaming
- Large/slower storage, e.g., memory
 - Address usually computed from values in register
 - Generally implemented as a cache hierarchy
 - hardware decides what is kept in fast memory
 - but software may provide "hints", e.g., don't cache or prefetch

Inside a Cache



Q: How many bits needed in tag? _____

Cache Algorithm (Read)

Look at Processor Address, search cache tags to find match.
Then either

Found in cache
a.k.a. HIT

Return copy
of data from
cache

Not in cache
a.k.a. MISS

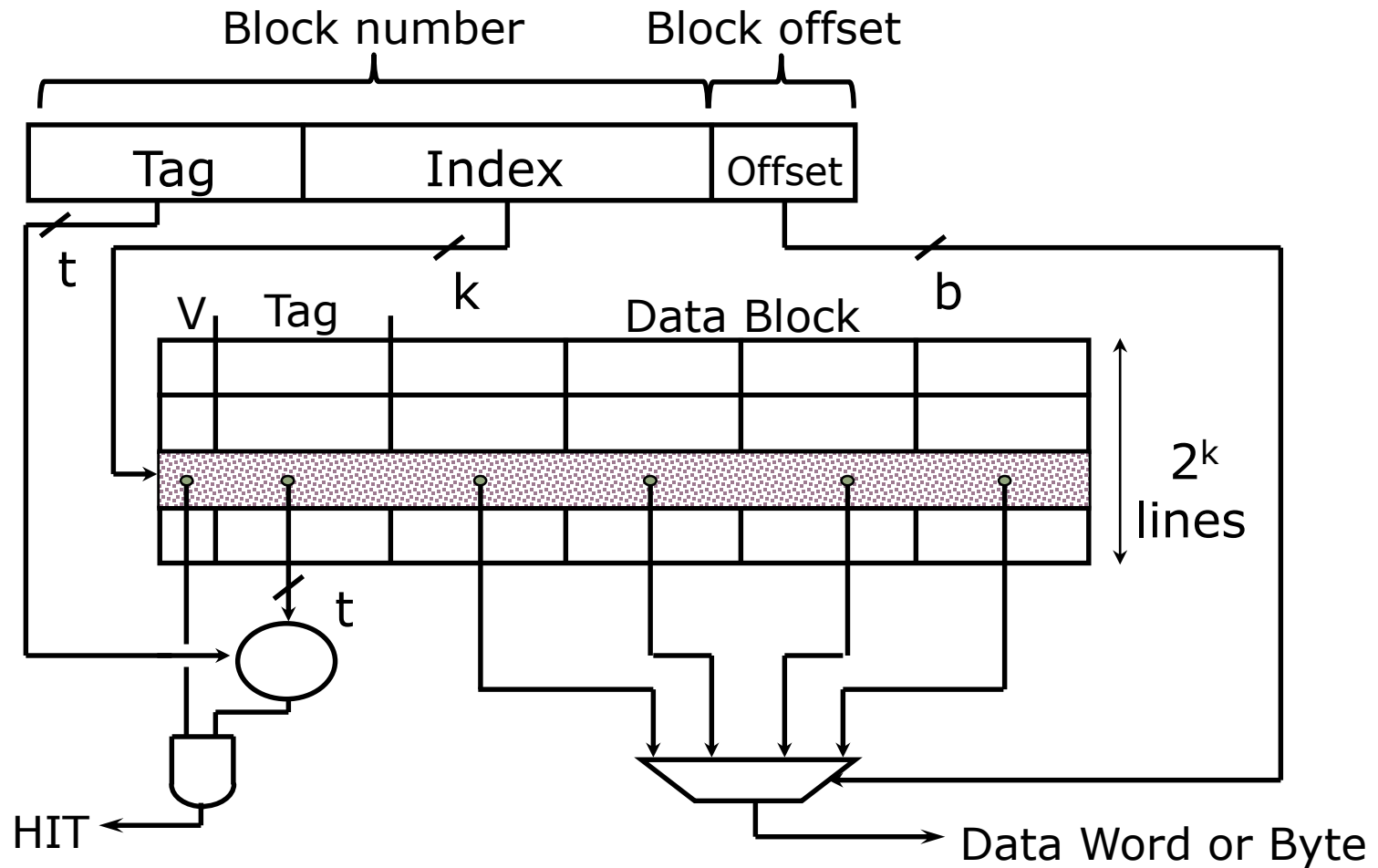
Read block of data from
Main Memory

Wait ...

Return data to processor
and update cache

Which line do we replace?

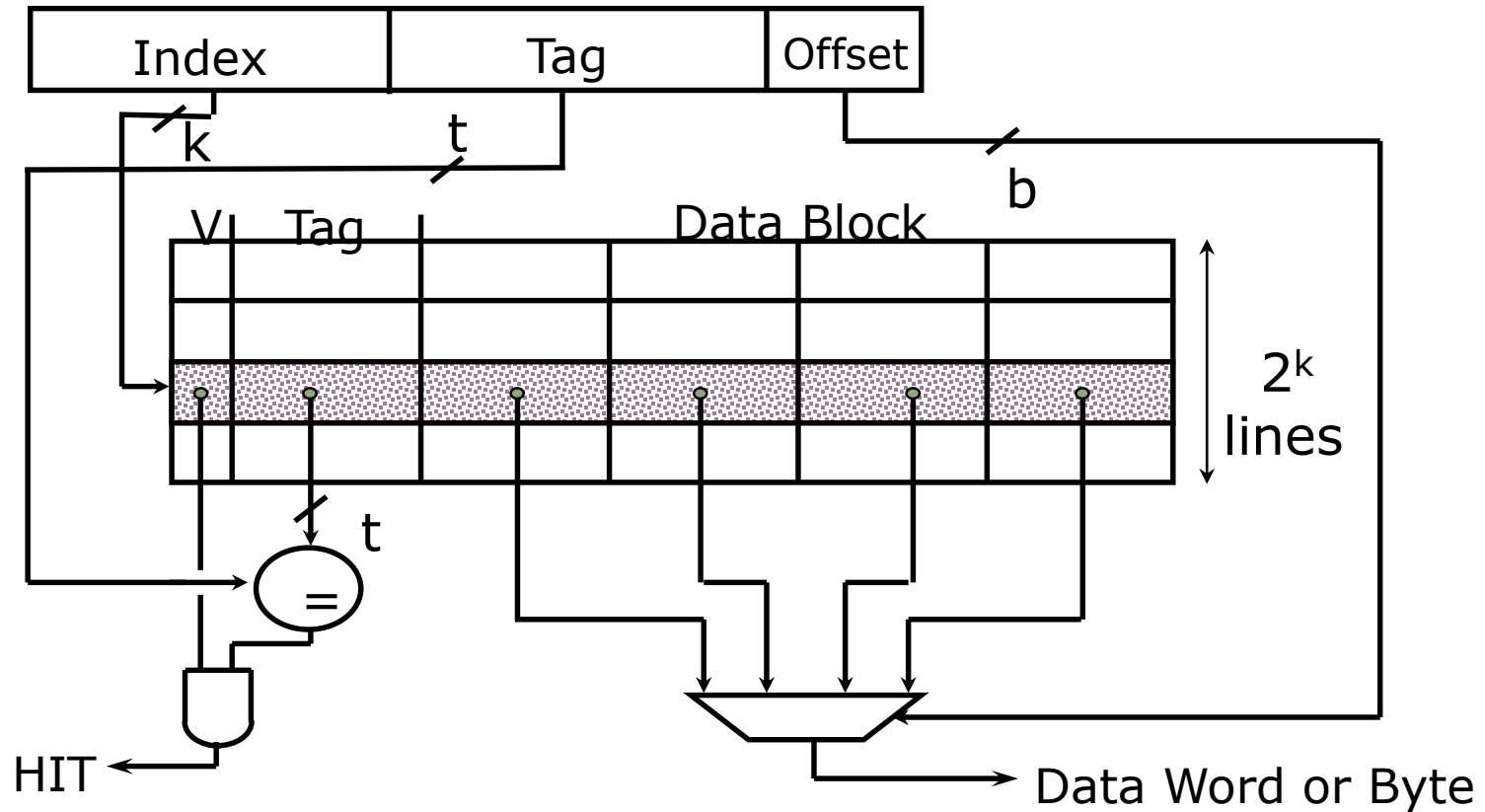
Direct-Mapped Cache



Q: What is a bad reference pattern? _____

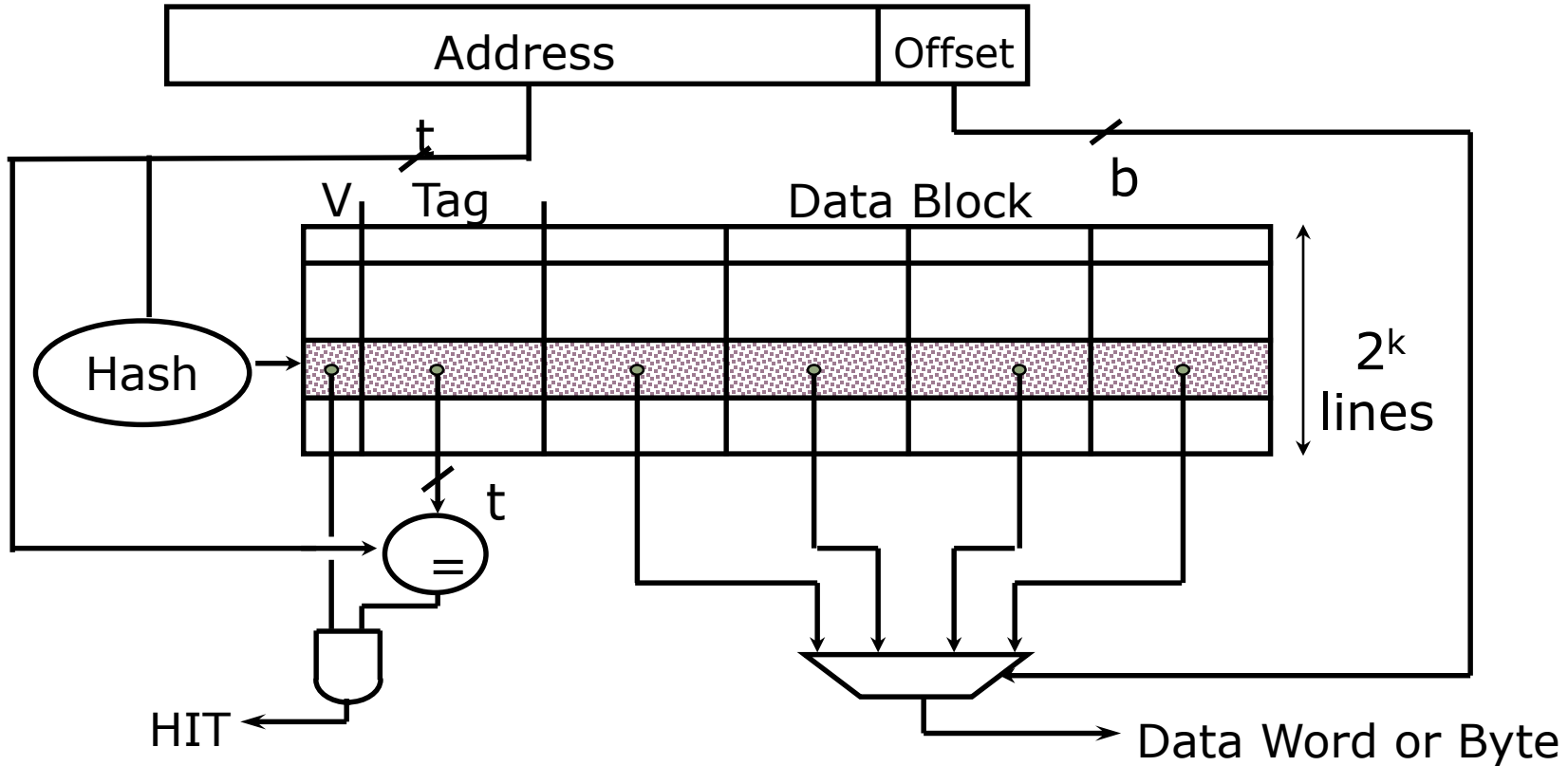
Direct Map Address Selection

higher-order vs. lower-order address bits



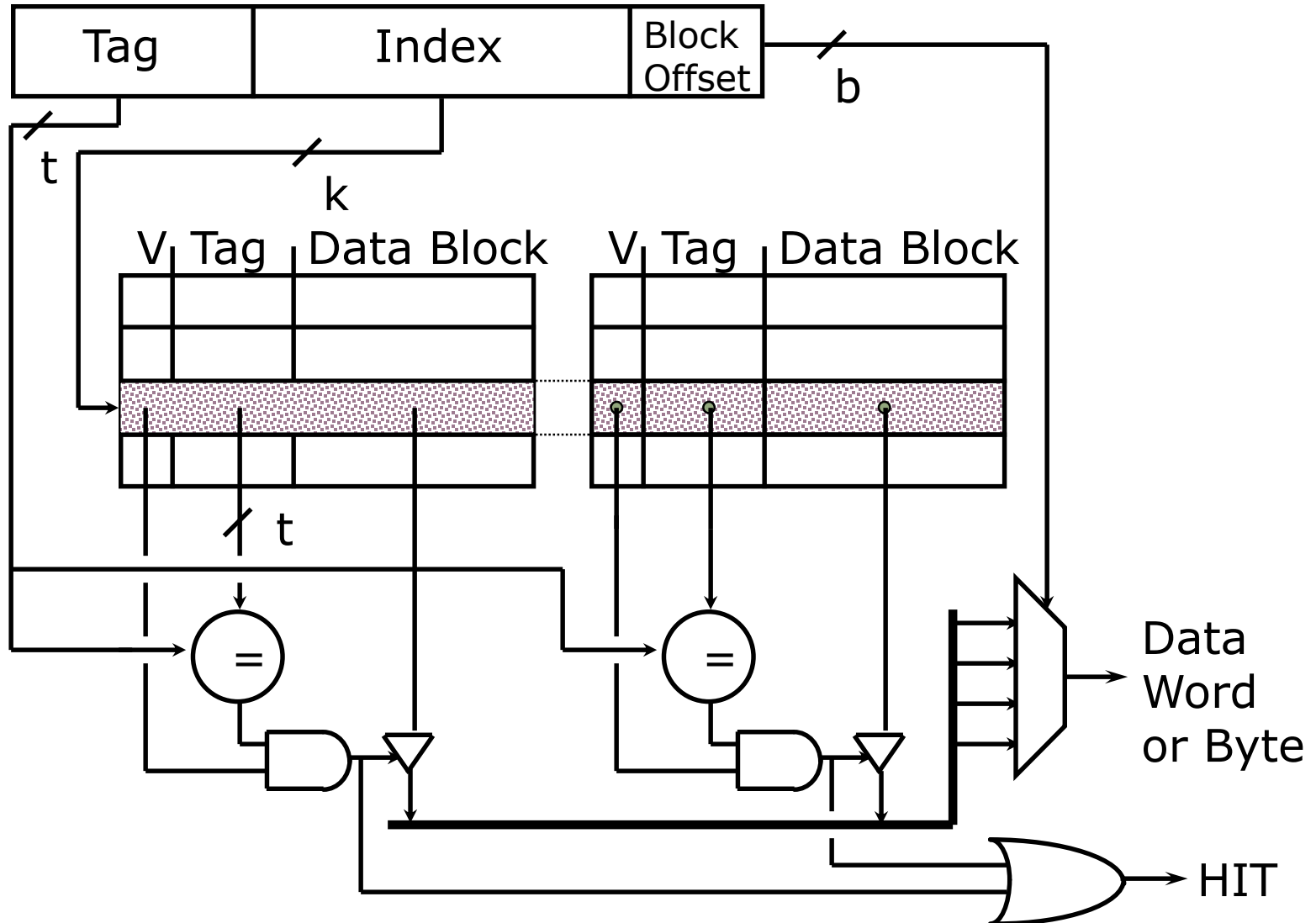
Q: Why might this be undesirable? _____

Hashed Address Mapping

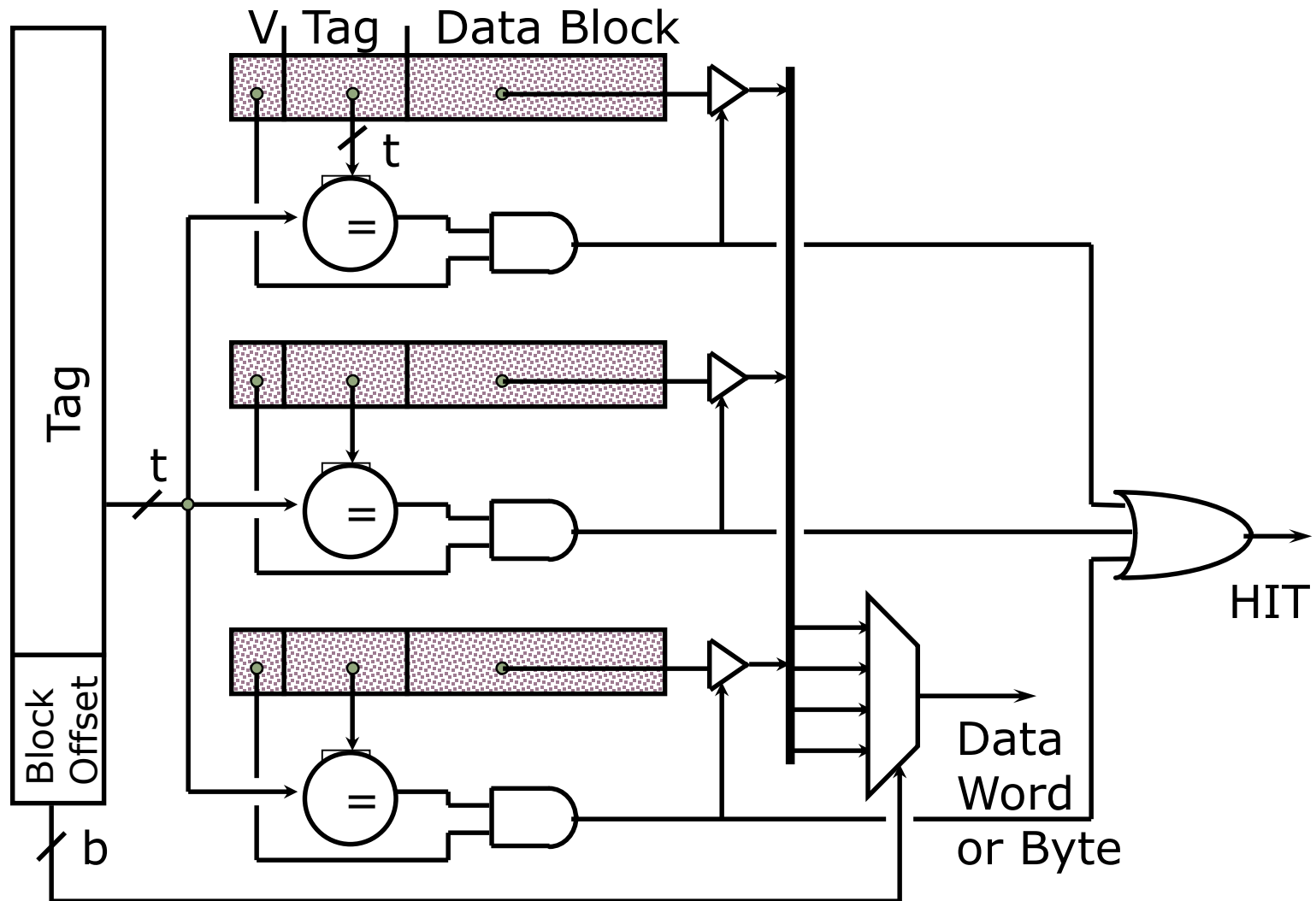


Q: What are the tradeoffs of hashing?

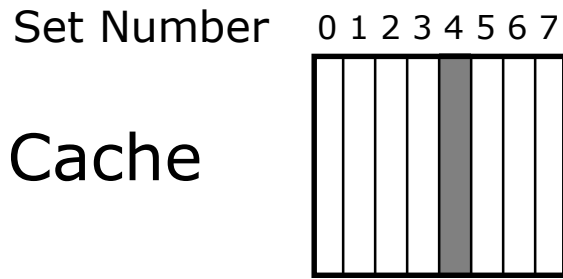
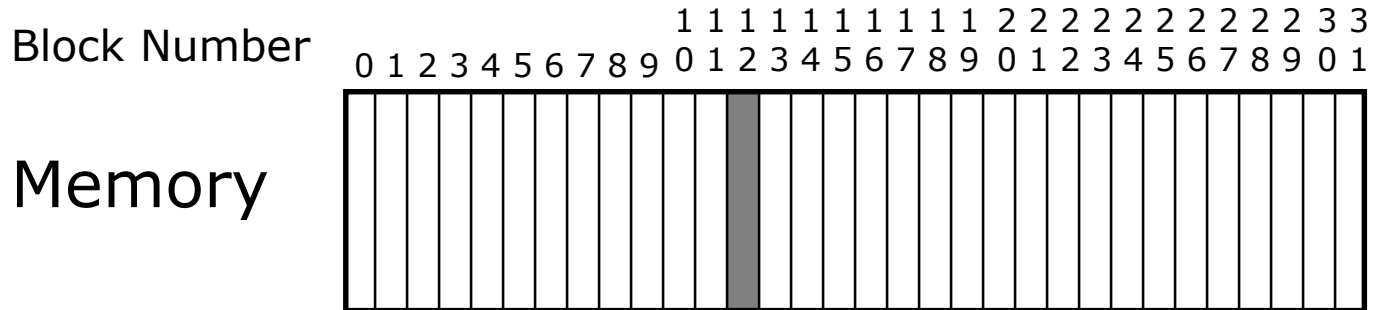
2-Way Set-Associative Cache



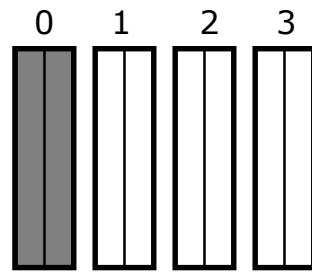
Fully Associative Cache



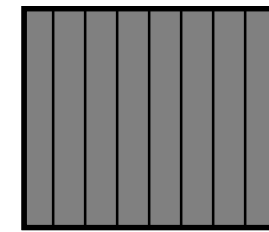
Placement Policy



Direct
Mapped
only into
block 4
 $(12 \bmod 8)$



(2-way) Set
Associative
anywhere in
set 0
 $(12 \bmod 4)$



Fully
Associative
anywhere

block 12
can be placed

Improving Cache Performance

Average memory access time =
Hit time + Miss rate x Miss penalty

To improve performance:

- reduce the hit time
- reduce the miss rate (e.g., larger, better policy)
- reduce the miss penalty (e.g., L2 cache)

What is the simplest design strategy?

Causes for Cache Misses

- *Compulsory:*
First reference to a block *a.k.a.* cold start misses
 - misses that would occur even with infinite cache
- *Capacity:*
cache is too small to hold all data the program needs
 - misses that would occur even under perfect placement & replacement policy
- *Conflict:*
misses from collisions due to block-placement strategy
 - misses that would not occur with full associativity

Effect of Cache Parameters on Performance

	Larger capacity cache	Higher associativity cache	Larger block size cache *
Compulsory misses			
Capacity misses			
Conflict misses			
Hit latency			
Miss latency			

* Assume substantial spatial locality

Block-level Optimizations

- Tags are too large, i.e., too much overhead
 - Simple solution: Larger blocks, but miss penalty could be large.
- Sub-block placement (aka sector cache)
 - A valid bit added to units smaller than the full block, called sub-blocks
 - Only read a sub-block on a miss
 - *If a tag matches, is the sub-block in the cache?*

100
300
204

1		1		1		1	
1		1		0		0	
0		1		0		1	

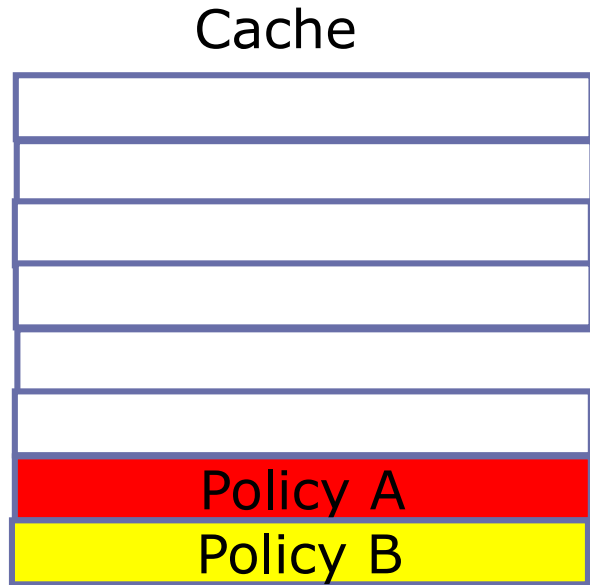
Replacement Policy

Which block from a set should be evicted?

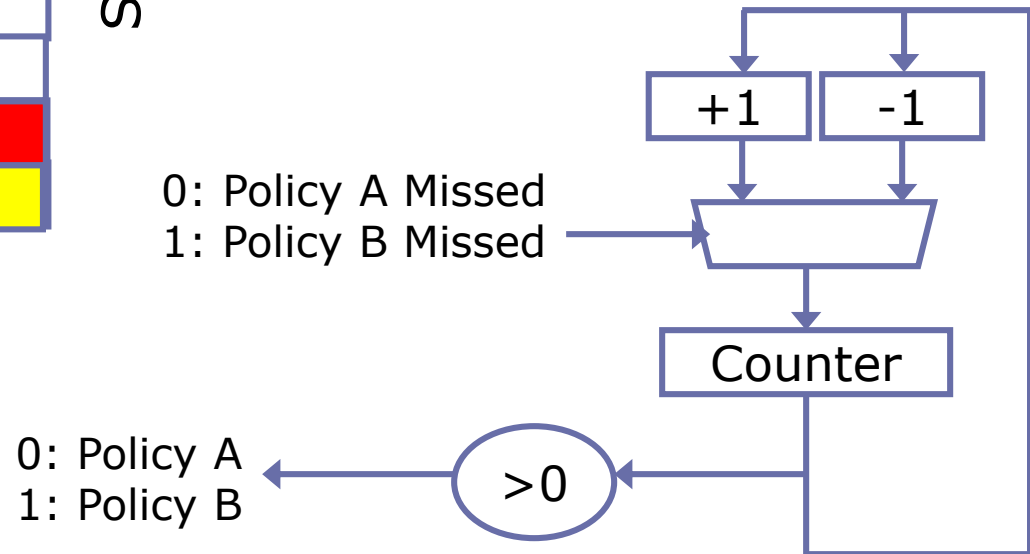
- Random
- Least Recently Used (LRU)
 - LRU cache state must be updated on every access
 - true implementation only feasible for small sets (2-way)
 - pseudo-LRU binary tree was often used for 4-8 way
- First In, First Out (FIFO) a.k.a. Round-Robin
 - used in highly associative caches
- Not Least Recently Used (NLRU)
 - FIFO with exception for most recently used block or blocks
- One-bit LRU
 - Each way represented by a bit. Set on use, replace first unused.

Multiple replacement policies

Use the best replacement policy for a program

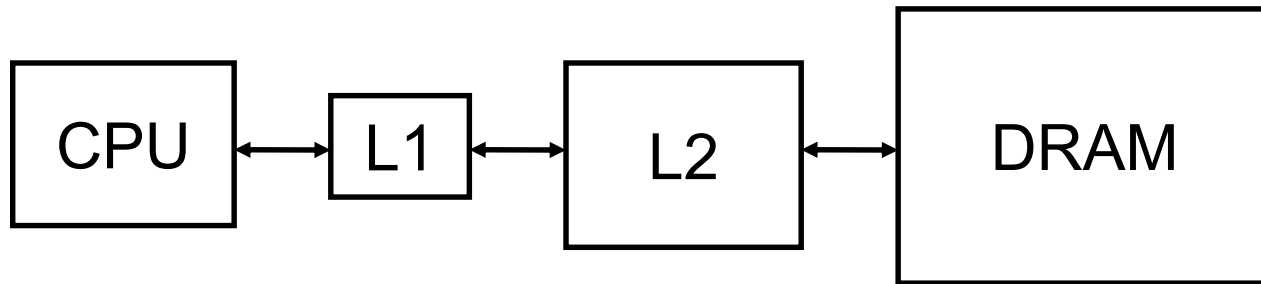


How do we decide which policy to use?



Multilevel Caches

- A memory cannot be large and fast
- Add level of cache to reduce miss penalty
 - Each level can have longer latency than level above
 - So, increase sizes of cache at each level



Metrics:

Local miss rate = misses in cache / accesses to cache

Global miss rate = misses in cache / CPU memory accesses

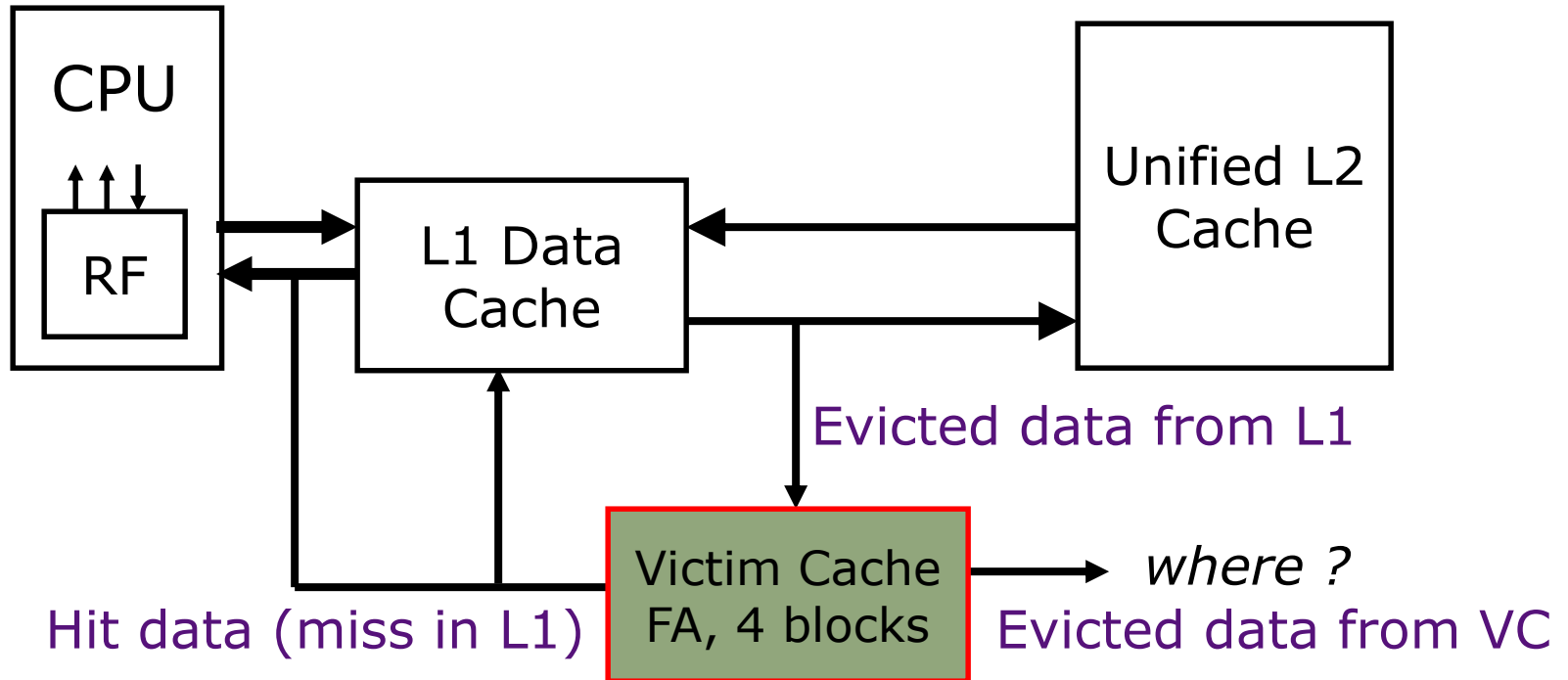
Misses per instruction (MPI) = misses in cache / number of instructions

Inclusion Policy

- **Inclusive multilevel cache:**
 - Inner cache holds copies of data in outer cache
 - On miss, line inserted in inner and outer cache; replacement in outer cache invalidates line in inner cache
 - External accesses need only check outer cache
 - Commonly used (e.g., Intel CPUs up to Broadwell)
- **Non-inclusive multilevel caches:**
 - Inner cache may hold data not in outer cache
 - Replacement in outer cache doesn't invalidate line in inner cache
 - Used in Intel Skylake, ARM
- **Exclusive multilevel caches:**
 - Inner cache and outer cache hold different data
 - Swap lines between inner/outer caches on miss
 - Used in AMD processors

Why choose one type or the other?

Victim Caches (HP 7200)

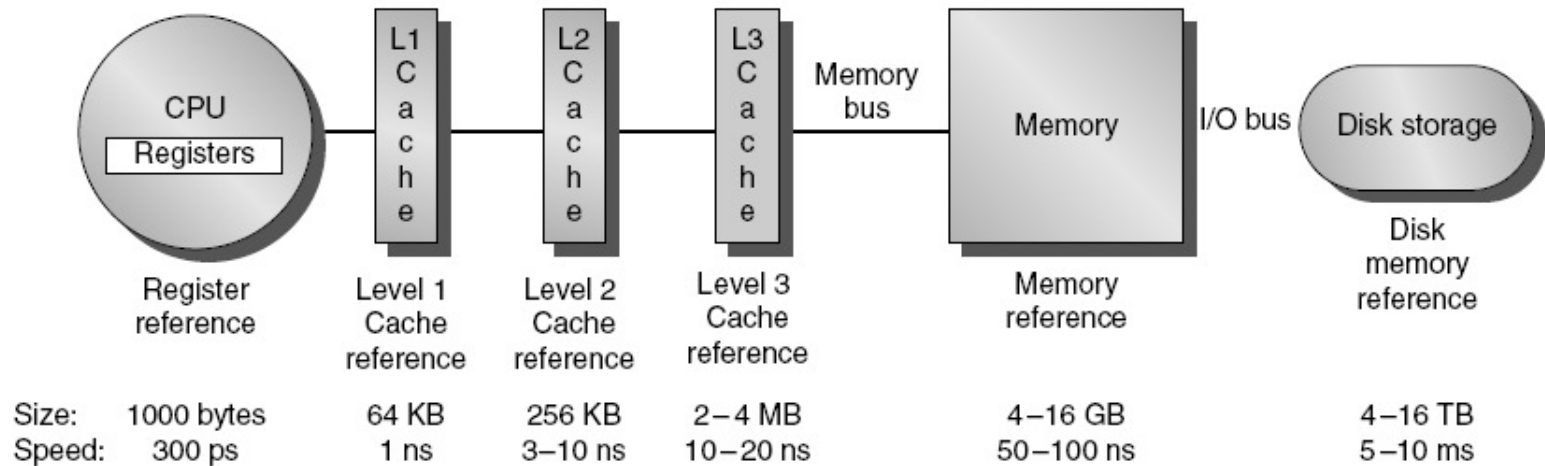


Victim cache is a small associative back up cache, added to a direct mapped cache, which holds recently evicted lines

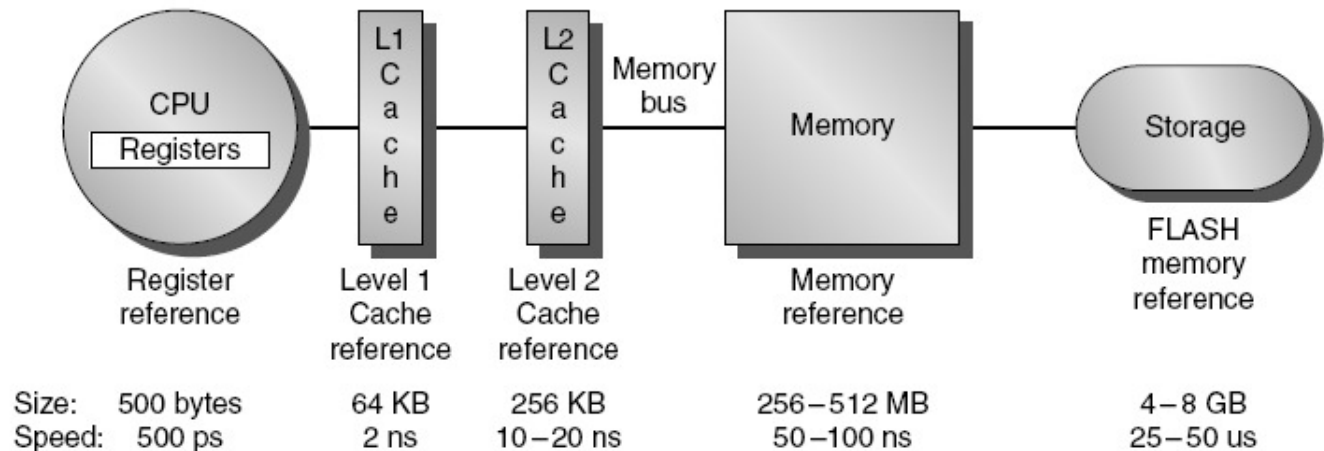
- First look up in direct mapped cache
- If miss, look in victim cache
- If hit in victim cache, swap hit line with line now evicted from L1
- If miss in victim cache, L1 victim -> VC, VC victim->?

Fast hit time of direct mapped but with reduced conflict misses

Typical memory hierarchies

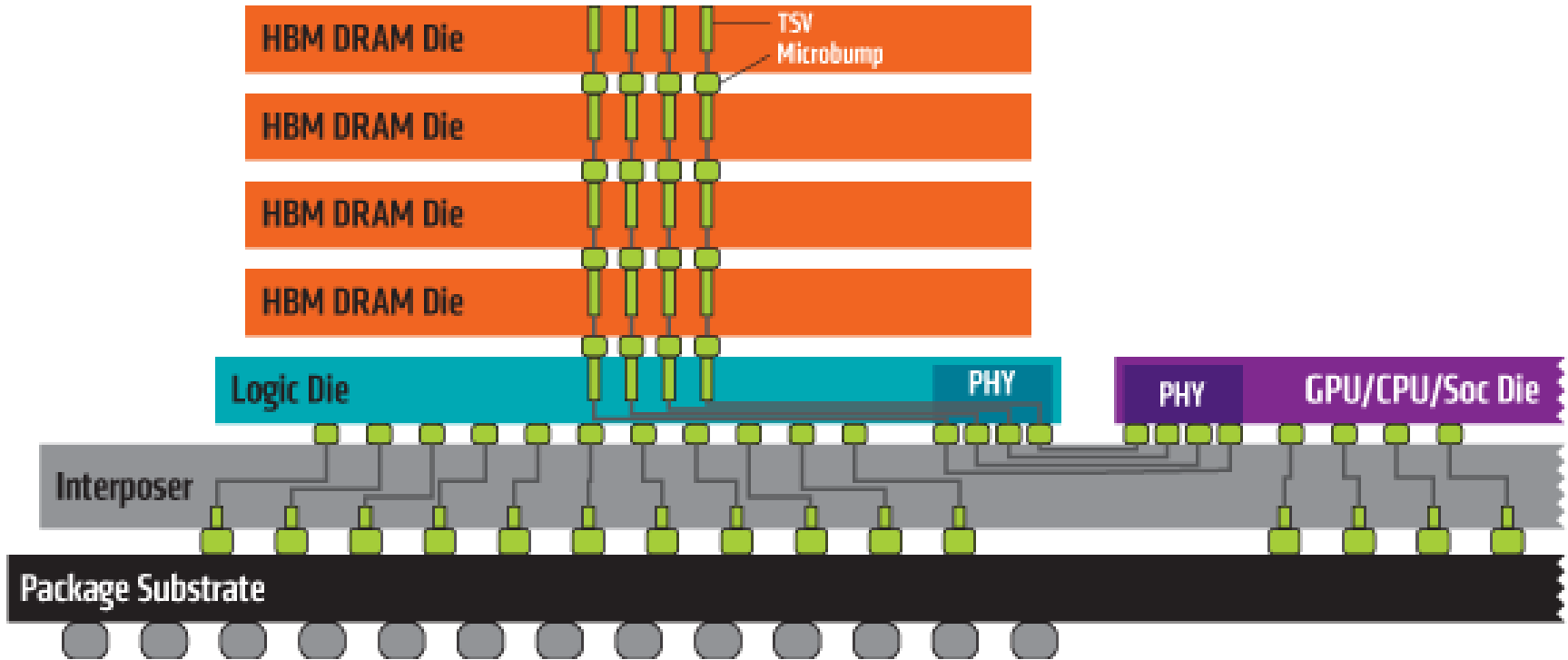


(a) Memory hierarchy for server



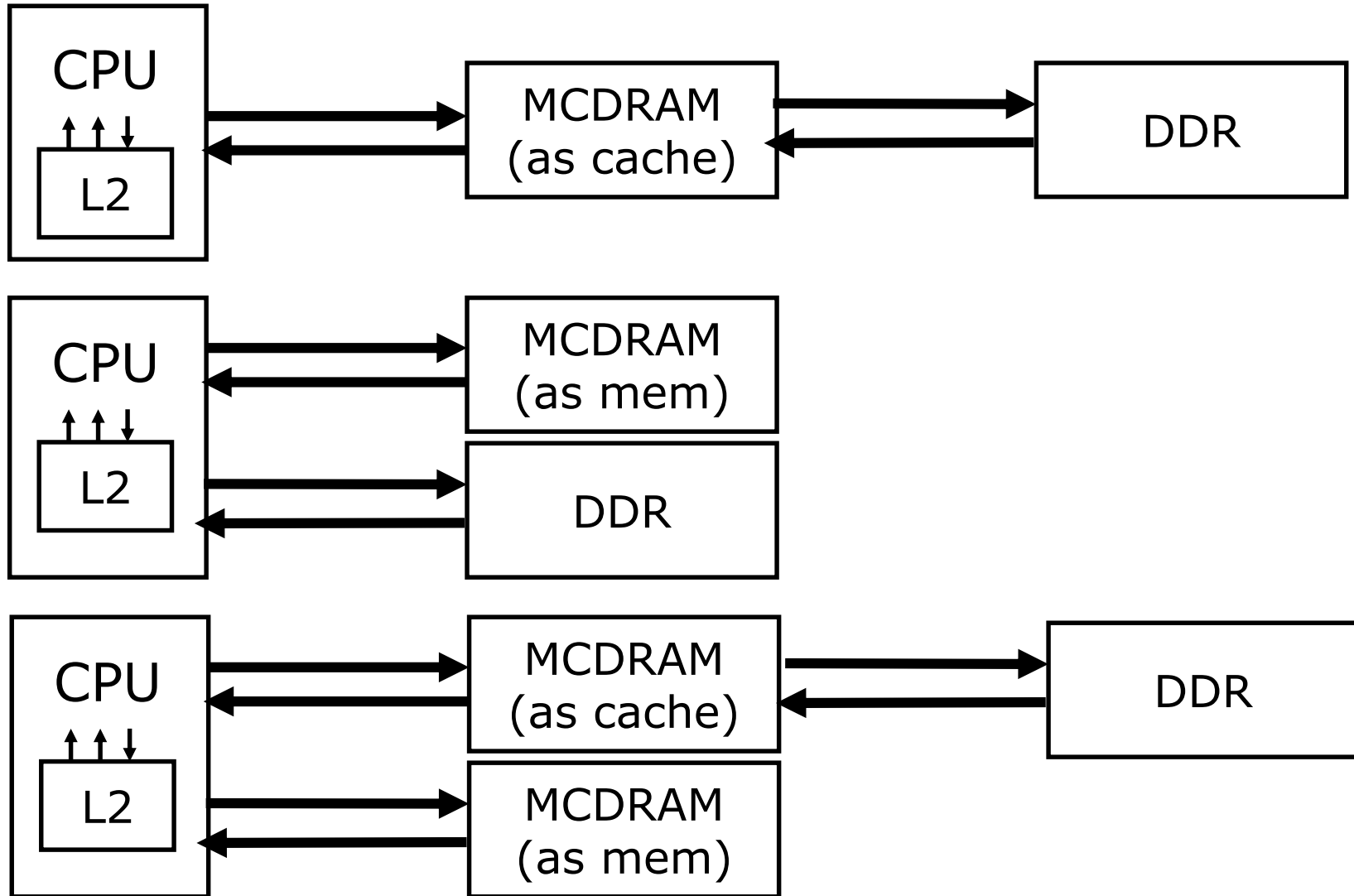
(b) Memory hierarchy for a personal mobile device

HBM DRAM or MCDRAM



Source: AMD

Mixed technology caching (Intel Knights Landing)



Thank you!

*Next lecture:
Virtual memory*