# Lecture-13 and 14 (Memory Hierarchy)
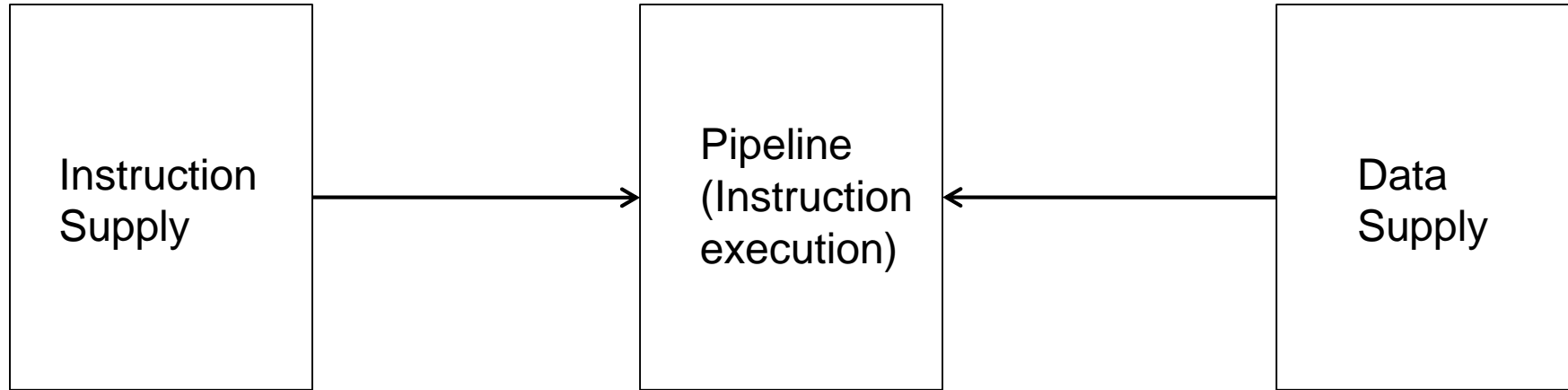## CS422-Spring 2020

Biswa@CSE-IITK

# The Ideal World

| Instruction Supply | Pipeline (Instruction execution) | Data Supply |
|---|---|---|

- Zero-cycle latency

- Infinite capacity

- Zero cost

- Perfect control flow

- No pipeline stalls

-Perfect data flow
 (reg/memory dependencies)

- Zero-cycle interconnect
 (operand communication)

- Enough functional units

- Zero latency compute

- Zero-cycle latency
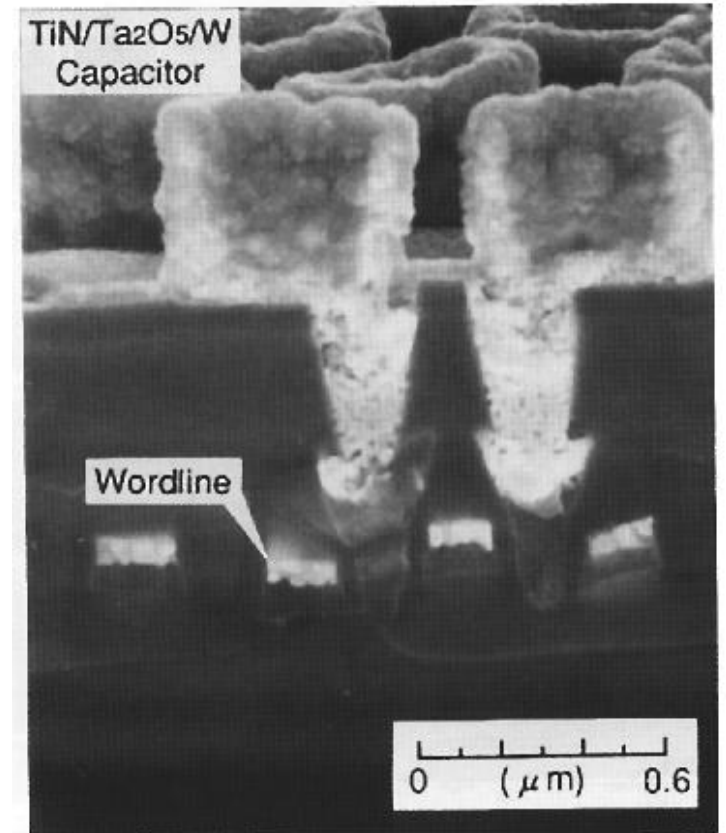
- Infinite capacity
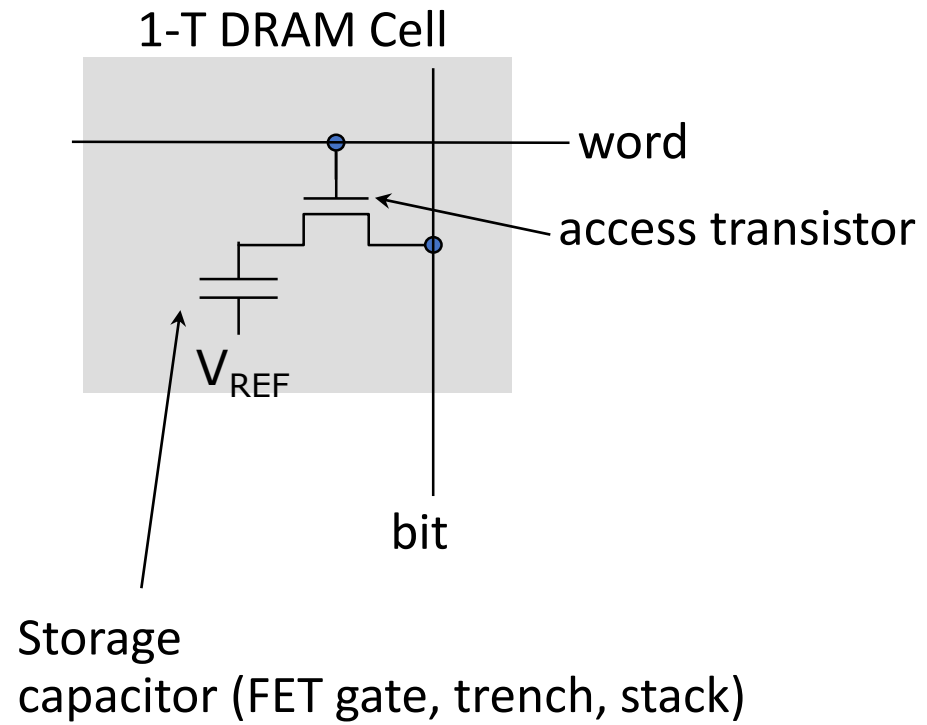
- Infinite bandwidth

- Zero cost

# Semiconductor Memory

- Semiconductor memory began to be competitive in early 1970s
  - Intel formed to exploit market for semiconductor memory
  - Early semiconductor memory was Static RAM  (SRAM).  SRAM cell internals similar to a latch (cross-coupled inverters).

- First commercial Dynamic RAM (DRAM) was Intel 1103
  - 1Kbit of storage on single chip
  - charge on a capacitor used to hold value

*Semiconductor memory quickly replaced core in '70s*

# One-transistor DRAM

1-T DRAM Cell

word

access transistor

$V_{REF}$

bit

Storage
capacitor (FET gate, trench, stack)

TiN/Ta₂O₅/W
Capacitor

Wordline

0        ($\mu$m)        0.6

# DRAM vs SRAM

- DRAM
  - Slower access (capacitor)
  - Higher density (1T 1C cell)
  - Lower cost
  - Requires refresh (power, performance, circuitry)

- SRAM
  - Faster access (no capacitor)
  - Lower density (6T cell)
  - Higher cost
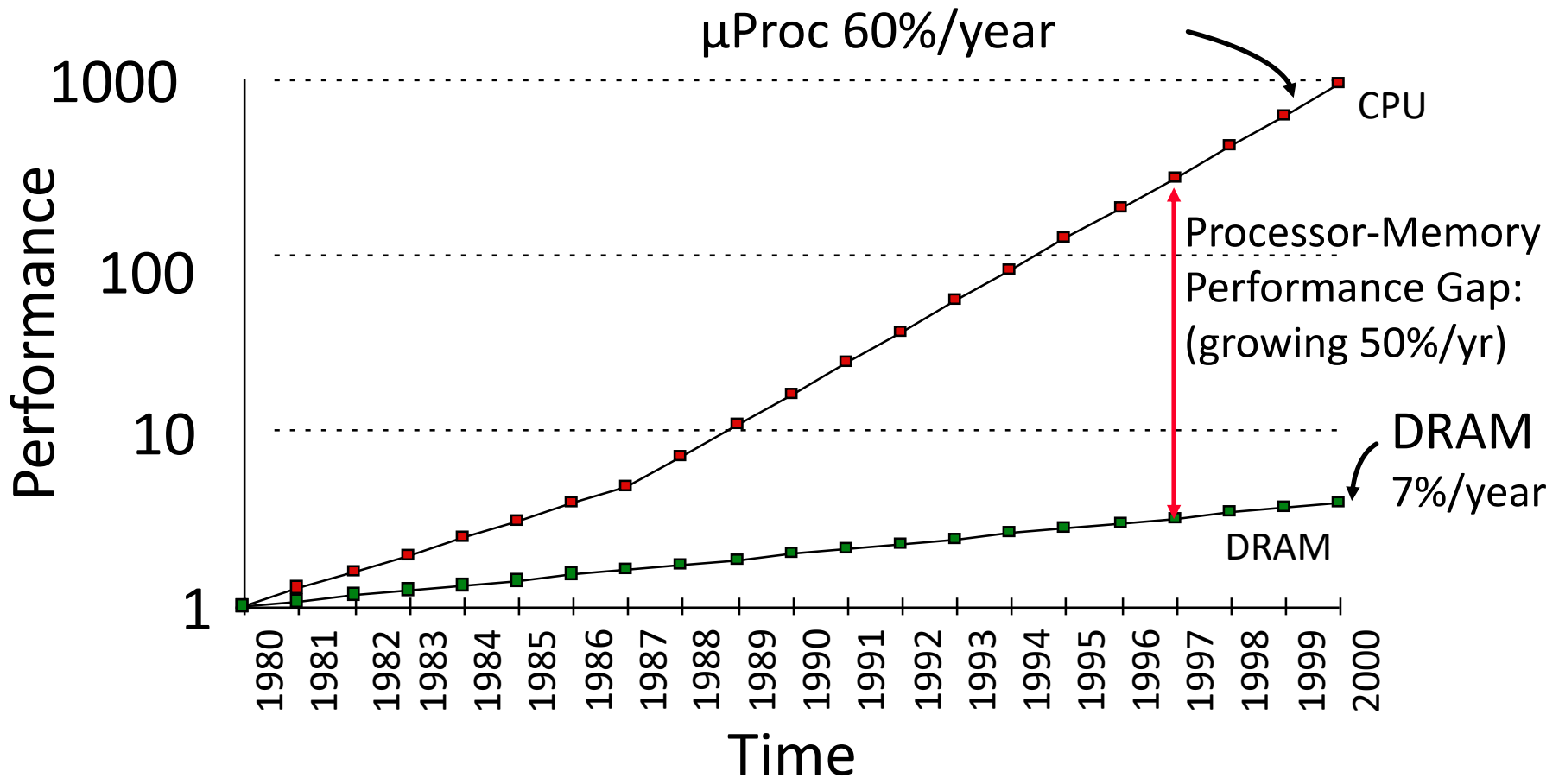  - No need for refresh

# The Problem?

- Bigger is slower
  - SRAM, 512 Bytes, sub-nanosec
  - SRAM,  KByte~MByte, ~nanosec
  - DRAM, Gigabyte, ~50 nanosec
  - Hard Disk, Terabyte, ~10 millisec

- Faster is more expensive (dollars and chip area)
  - SRAM, < 10$ per Megabyte
  - DRAM, < 1$ per Megabyte
  - Hard Disk < 1$ per Gigabyte
  - These sample values scale with time

- Other technologies have their place as well
  - Flash memory, PC-RAM, MRAM, RRAM (not mature yet)
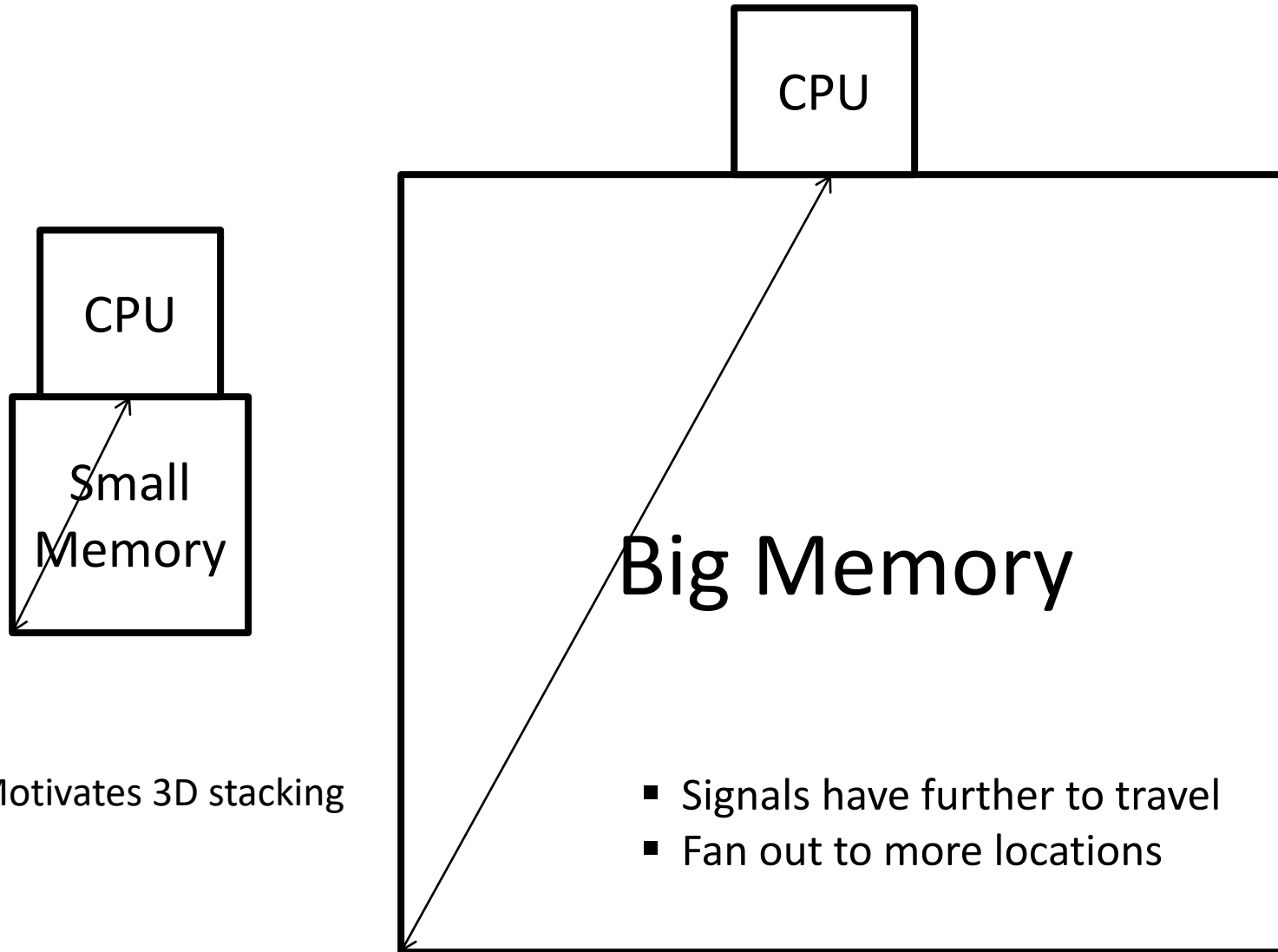
# Why Memory Hierarchy?

- We want both fast and large

- But we cannot achieve both with a single level of memory

- Idea: Have multiple levels of storage (progressively bigger and slower as the levels are farther from the processor) and ensure most of the data the processor needs is kept in the fast(er) level(s)
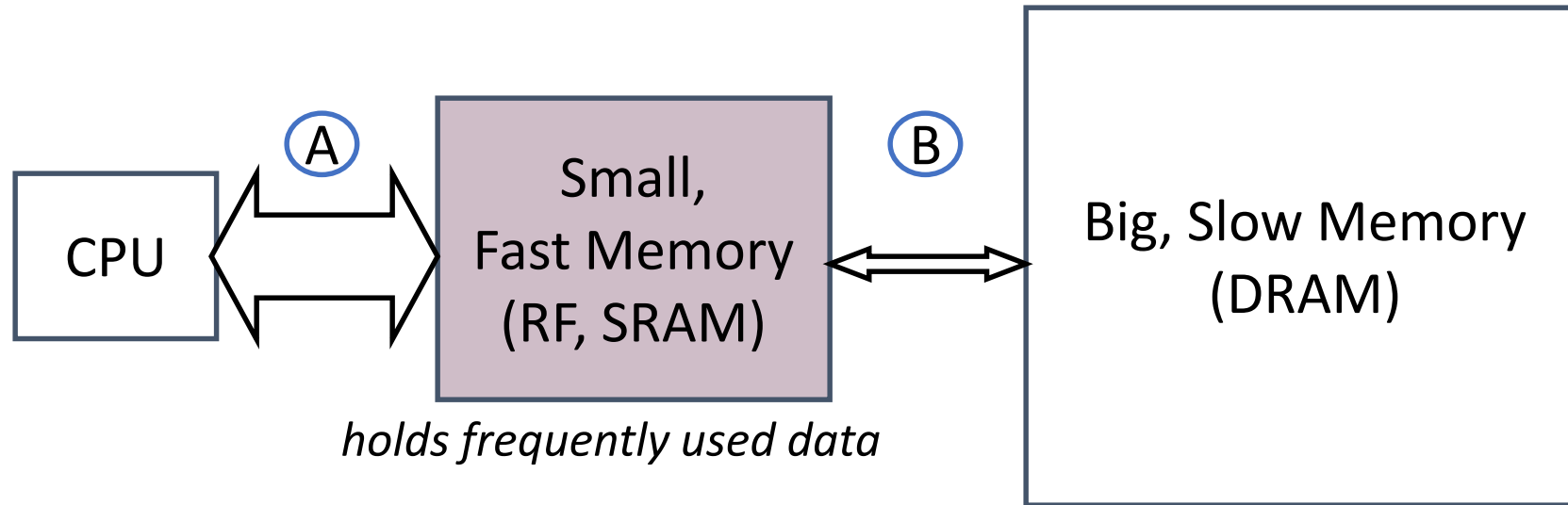
# Memory Wall Problem

# Size Affects Latency

CPU

CPU

Small Memory

Big Memory

Motivates 3D stacking

- Signals have further to travel
- Fan out to more locations

# Memory Hierarchy

CPU ⟷ (A) ⟷ Small, Fast Memory (RF, SRAM) ⟷ (B) ⟷ Big, Slow Memory (DRAM)

*holds frequently used data*

- *capacity*:  Register << SRAM << DRAM
- *latency*:   Register << SRAM << DRAM
- *bandwidth:* on-chip >> off-chip

On a data access:

  *if* data ∈ fast memory ⟹ low latency access *(SRAM)*
  *if* data ∉ fast memory ⟹ high latency access *(DRAM)*

# Access Patterns



Donald J. Hatfield, Jeanette Gerald: Program Restructuring for Virtual Memory. IBM Systems Journal 10(3): 168-192 (1971)
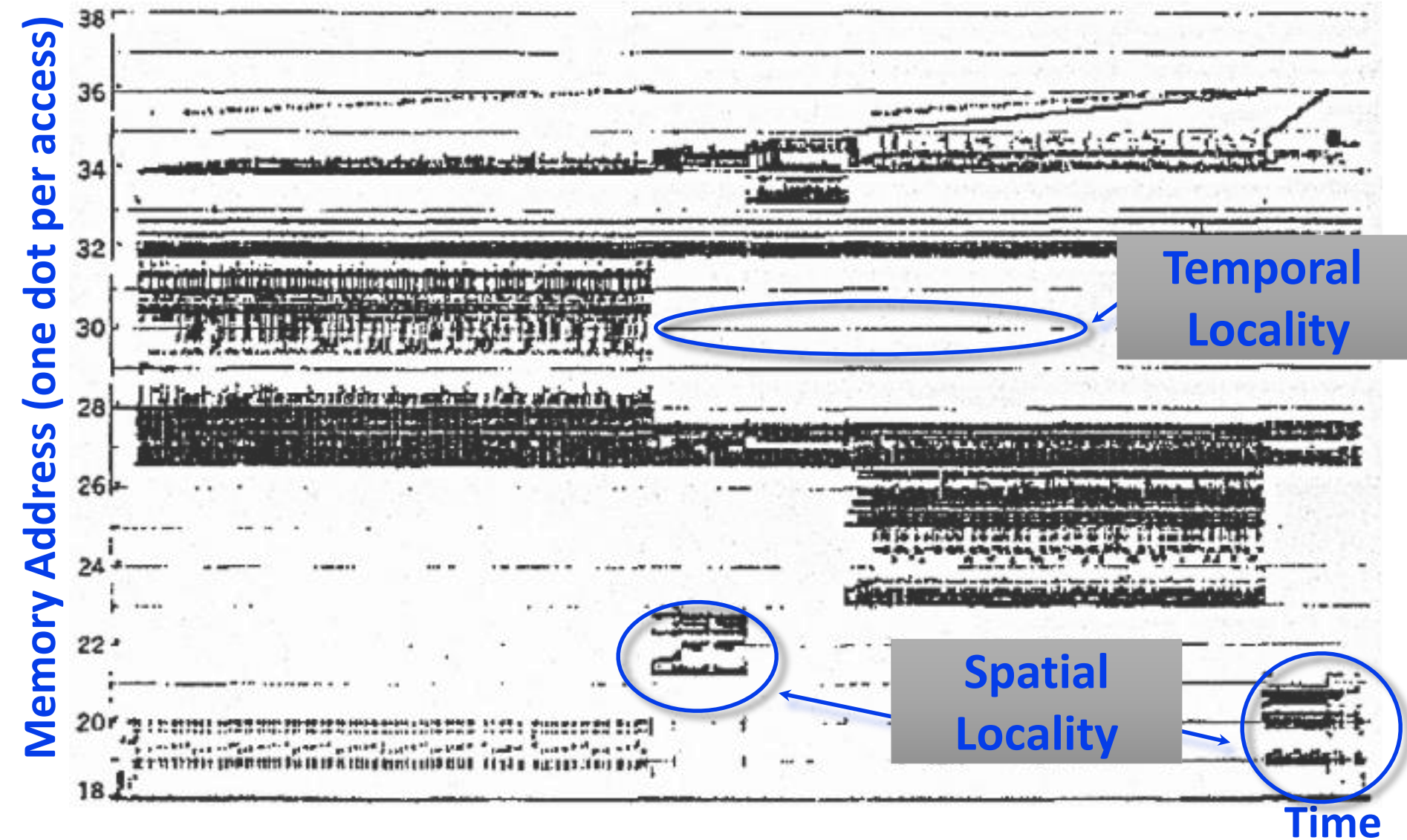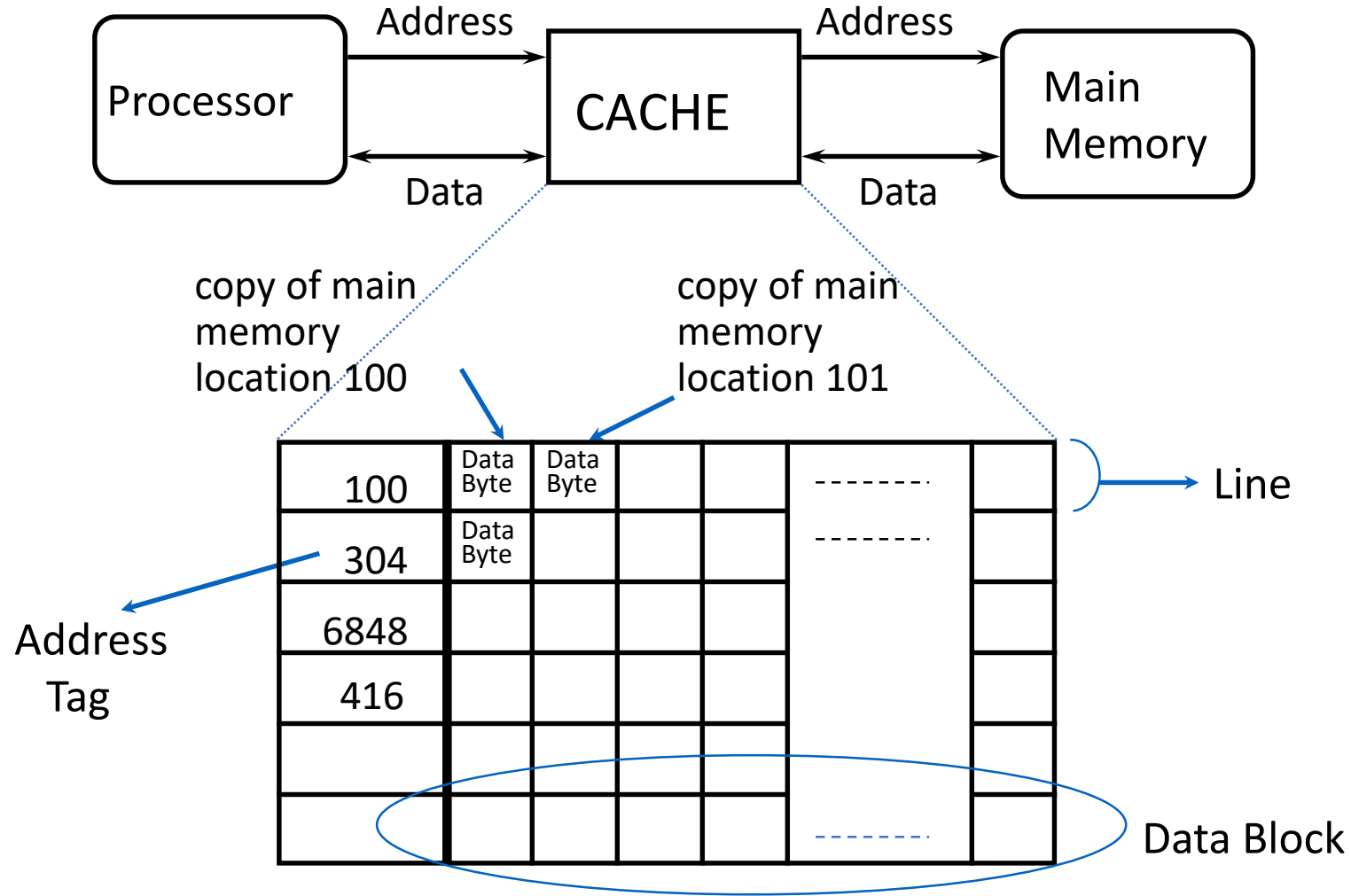
# Examples

# Locality of Reference

- **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.

# Again

# Inside a Cache

# Cache Events

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

Found in cache
a.k.a.  HIT

Not in cache
a.k.a. MISS

Return copy
of data from
cache

Read block of data from
Main Memory

Wait ...

Return data to processor
and update cache

Q: Which line do we replace?

# Placement Policy

Block Number

1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 3 3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1

Memory

Set Number

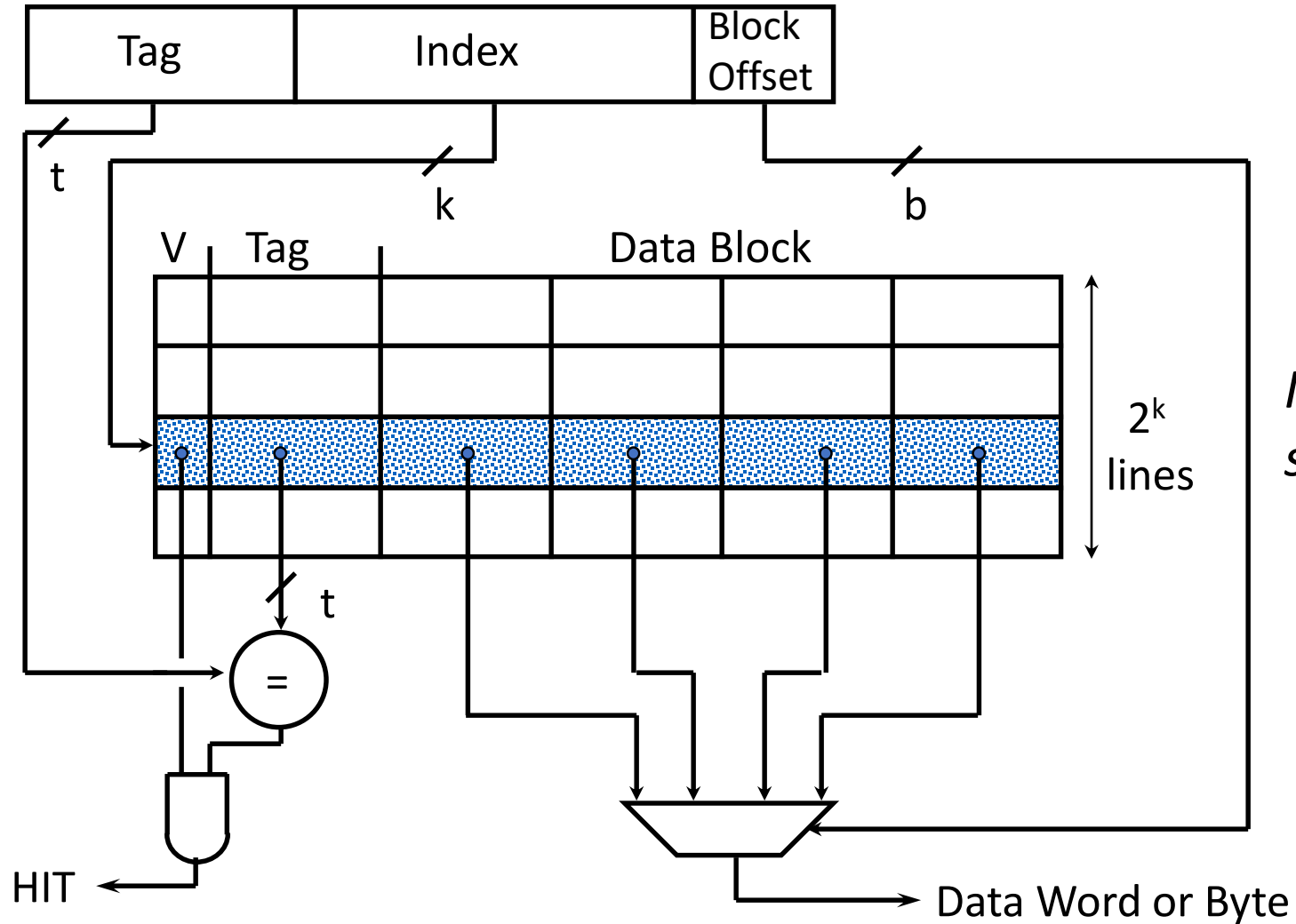    0   1   2   3       0 1 2 3 4 5 6 7

Cache

Fully
Associative
anywhere

(2-way) Set
Associative
anywhere in
set 0
*(12 mod 4)*

Direct
Mapped
only into
block 4
*(12 mod 8)*

block 12
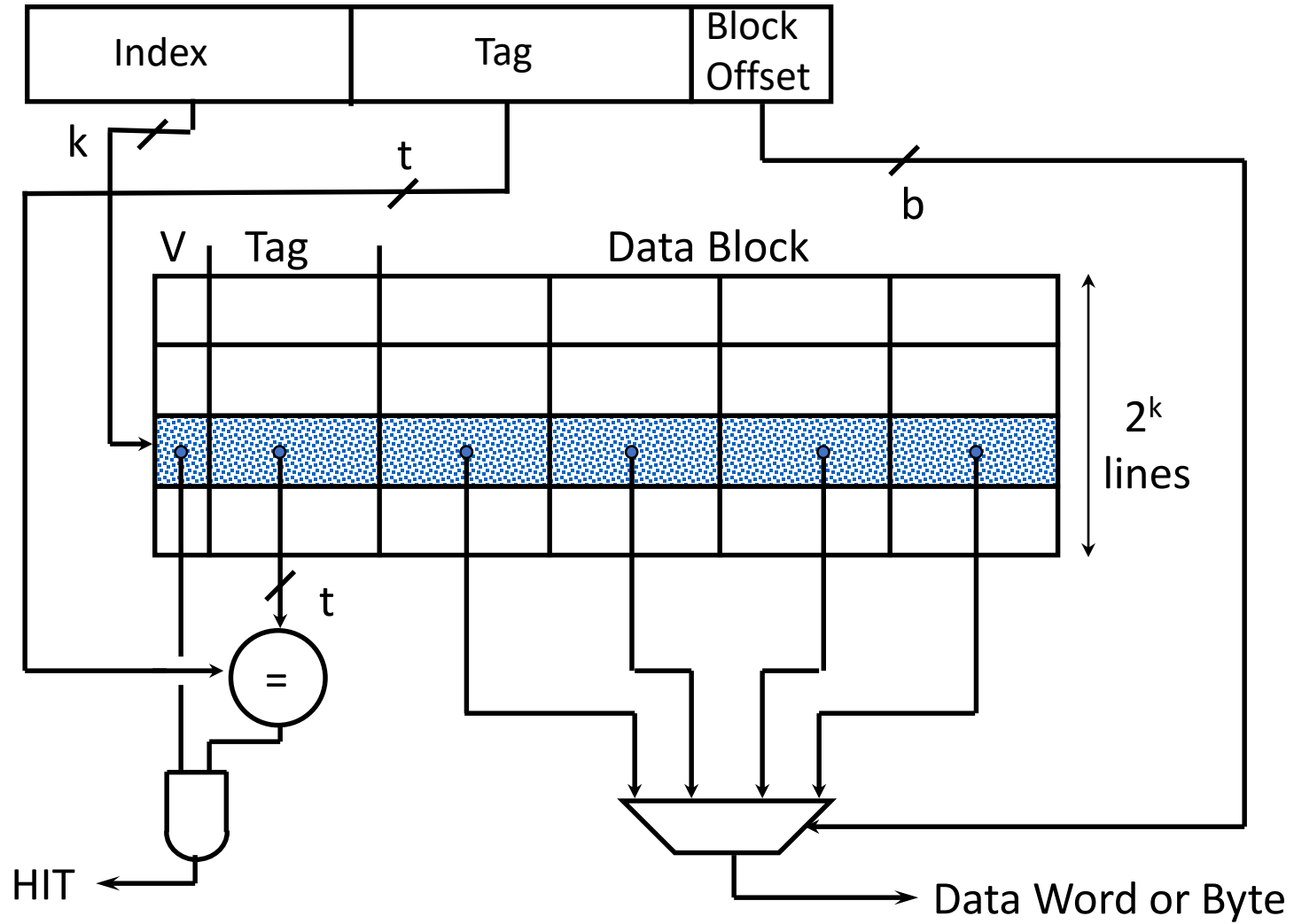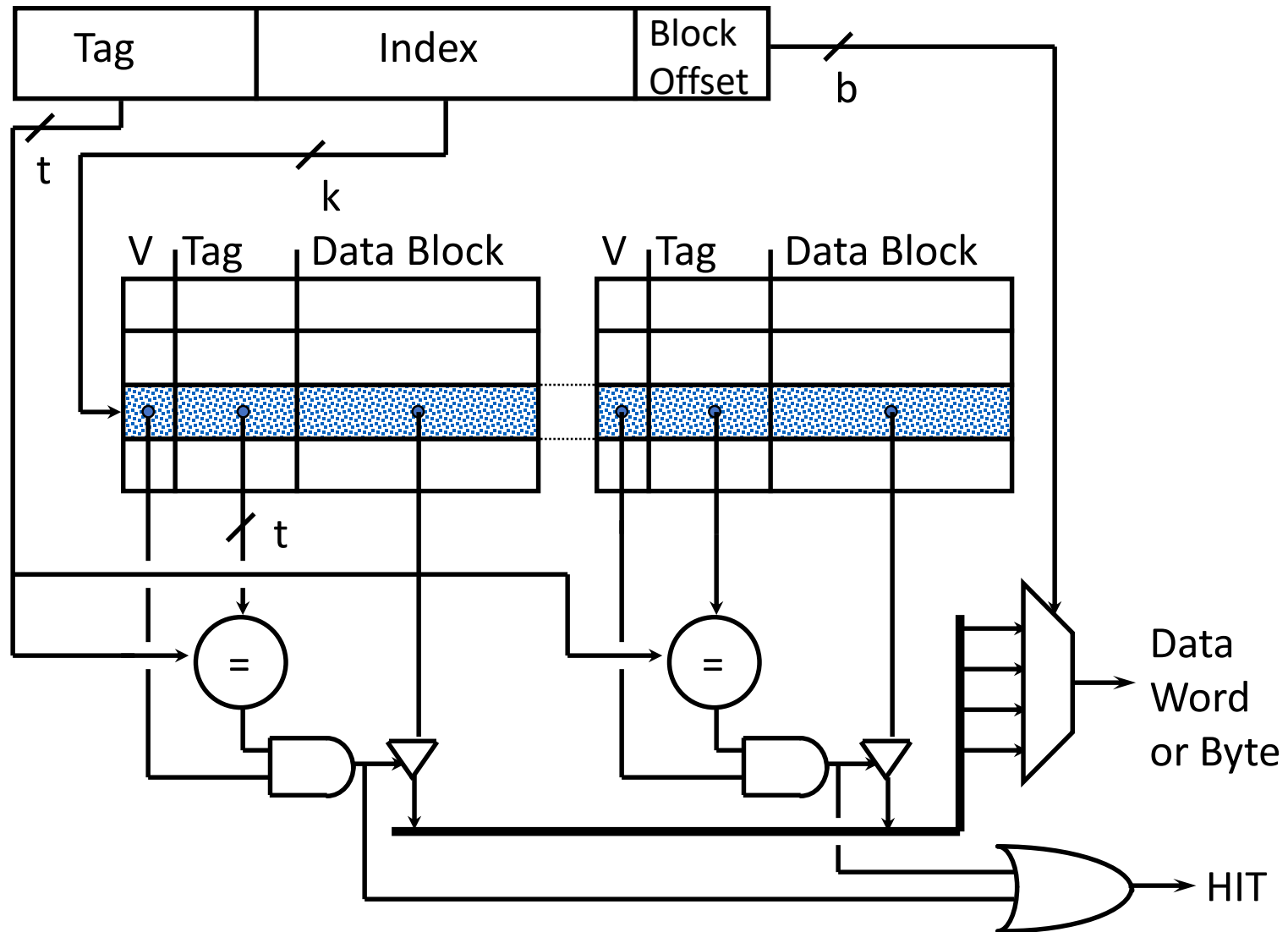can be placed

# Direct Mapped

# High bits or Low bits

# Set-Associative

# Fully-associative

# Block (line) Size ?

| Tag |
|-----|

| Word0 | Word1 | Word2 | Word3 |
|-------|-------|-------|-------|

4 word block, b=2

Split CPU address

| block address | offset$_b$ |
|---------------|------------|

32-b bits                b bits

$2^b$ = block size *a.k.a* line size (in bytes)

Larger block size has distinct hardware advantages
- less tag overhead
- exploit fast burst transfers from DRAM
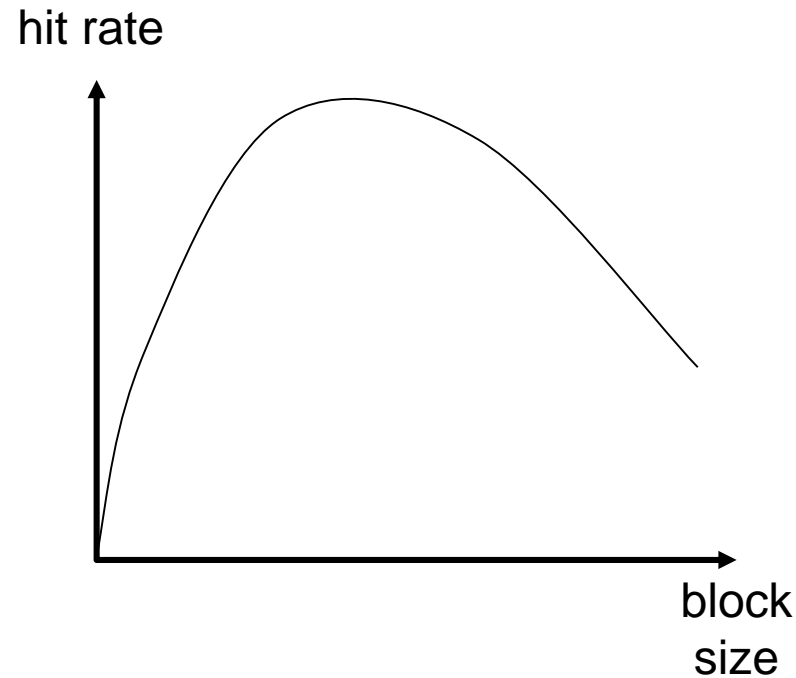- exploit fast burst transfers over wide busses

*What are the disadvantages of increasing block size?*

*Fewer blocks => more conflicts.  Can waste bandwidth.*
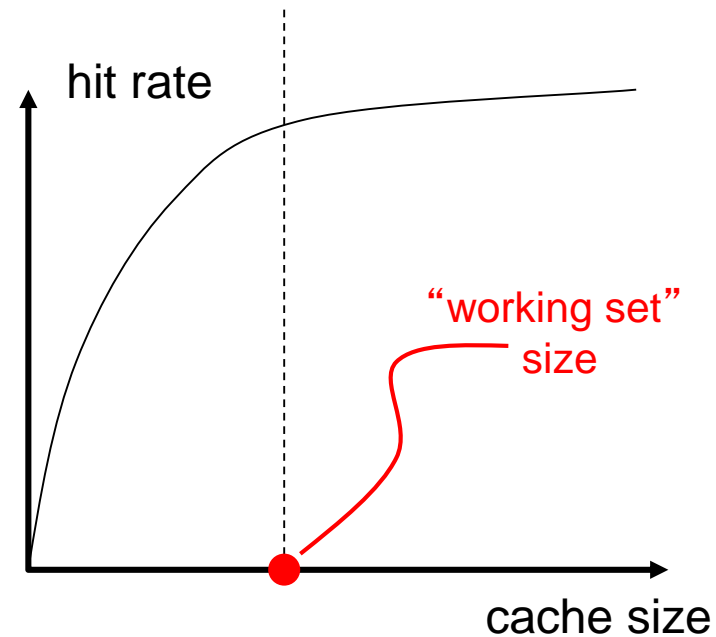
# Block Size?

- Block size is the data that is associated with an address tag
  - not necessarily the unit of transfer between hierarchies
    - Sub-blocking: A block divided into multiple pieces (each with V bit)
      - Can improve "write" performance

- Too small blocks
  - don't exploit spatial locality well
  - have larger tag overhead

- Too large blocks
  - too few total # of blocks
    - likely-useless data transferred
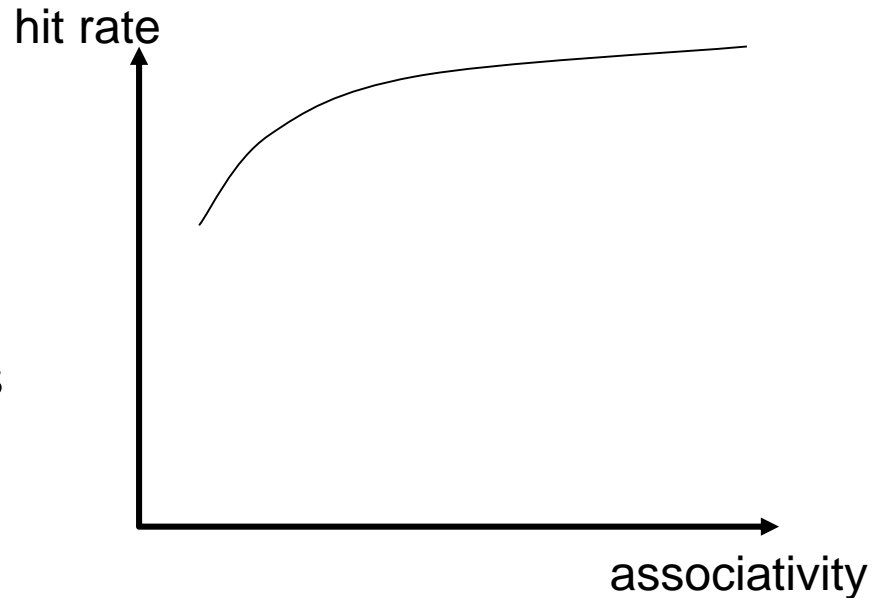    - Extra bandwidth/energy consumed

# Cache Size

- Cache size: total data (not including tag) capacity
  - ❑ bigger can exploit temporal locality better
  - ❑ not ALWAYS better
- Too large a cache adversely affects hit and miss latency
  - ❑ smaller is faster => bigger is slower
  - ❑ access time may degrade critical path
- Too small a cache
  - ❑ doesn't exploit temporal locality well
  - ❑ useful data replaced often

- Working set: the whole set of data the executing application references
  - ❑ Within a time interval

hit rate

"working set" size

cache size

# Associativity

- How many blocks can map to the same index (or set)?

- Larger associativity
  - lower miss rate, less variation among programs
  - diminishing returns, higher hit latency

- Smaller associativity
  - lower cost
  - lower hit latency
    - Especially important for L1 caches

- Power of 2 associativity?

# CPU – Cache Interaction