



CS 152 Computer Architecture and Engineering

Lecture 6 - Memory

Krste Asanovic

Electrical Engineering and Computer Sciences
University of California at Berkeley

<http://www.eecs.berkeley.edu/~krste>
<http://inst.eecs.berkeley.edu/~cs152>



Last time in Lecture 5

- Control hazards (branches, interrupts) are most difficult to handle as they change which instruction should be executed next
- Speculation commonly used to reduce effect of control hazards (predict sequential fetch, predict no exceptions)
- Branch delay slots make control hazard visible to software
- Precise exceptions: stop cleanly on one instruction, all previous instructions completed, no following instructions have changed architectural state
- To implement precise exceptions in pipeline, shift faulting instructions down pipeline to “commit” point, where exceptions are handled in program order



Early Read-Only Memory Technologies

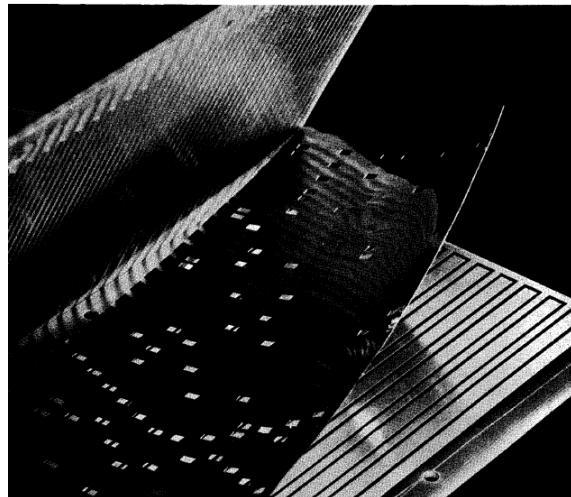
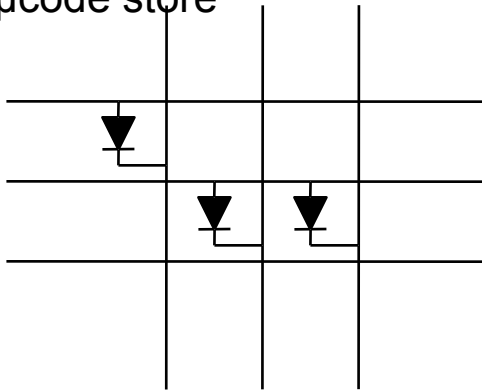


Punched cards, From early 1700s through Jacquard Loom, Babbage, and then IBM

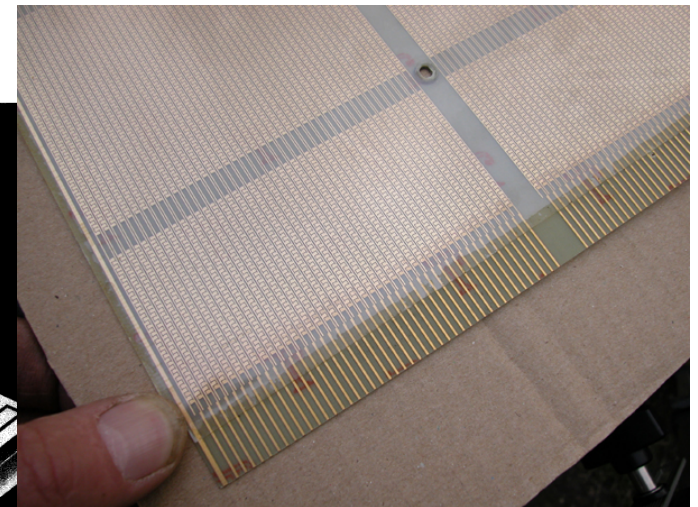


Punched paper tape, instruction stream in Harvard Mk 1

Diode Matrix, EDSAC-2 μ code store



IBM Card Capacitor ROS

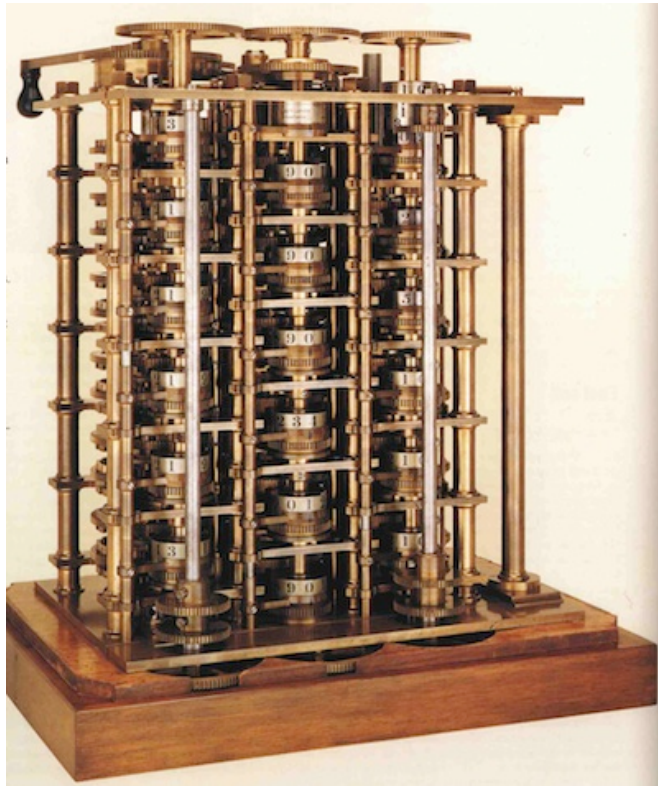


IBM Balanced Capacitor ROS

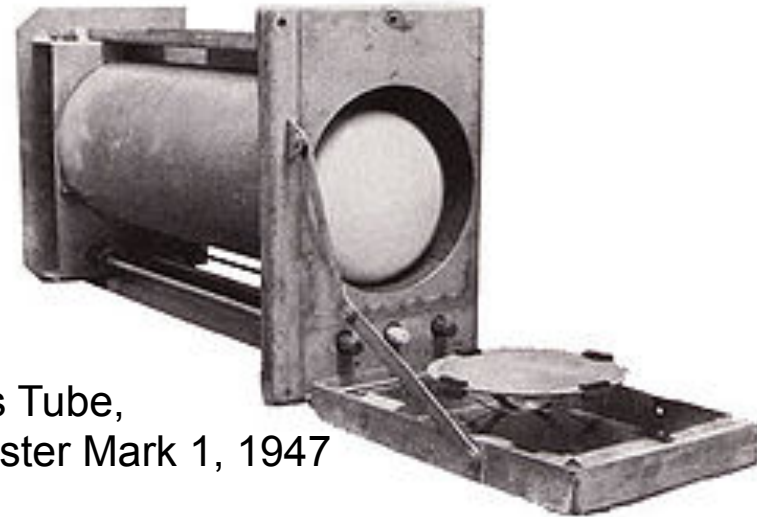


Early Read/Write Main Memory Technologies

Babbage, 1800s: Digits stored on mechanical wheels

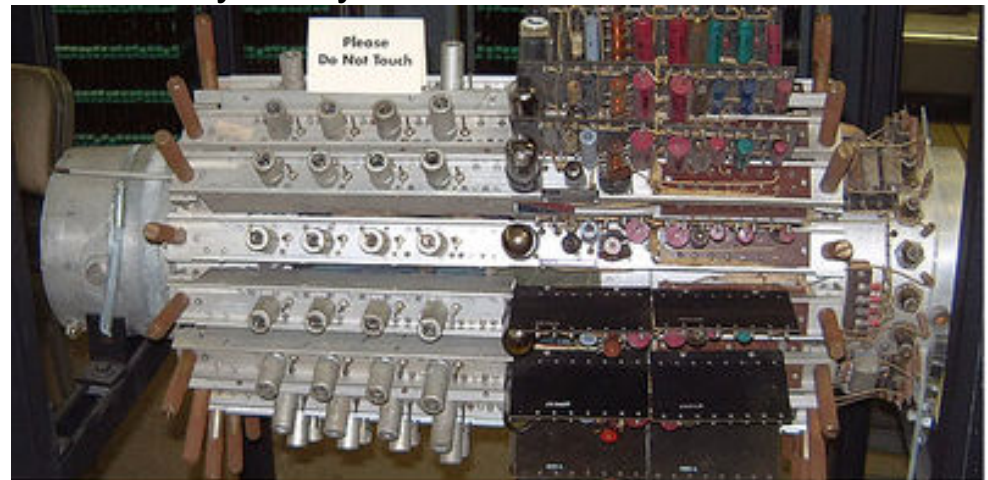


Also, regenerative capacitor memory on Atanasoff-Berry computer, and rotating magnetic drum memory on IBM 650



Williams Tube,
Manchester Mark 1, 1947

Mercury Delay Line, Univac 1, 1951

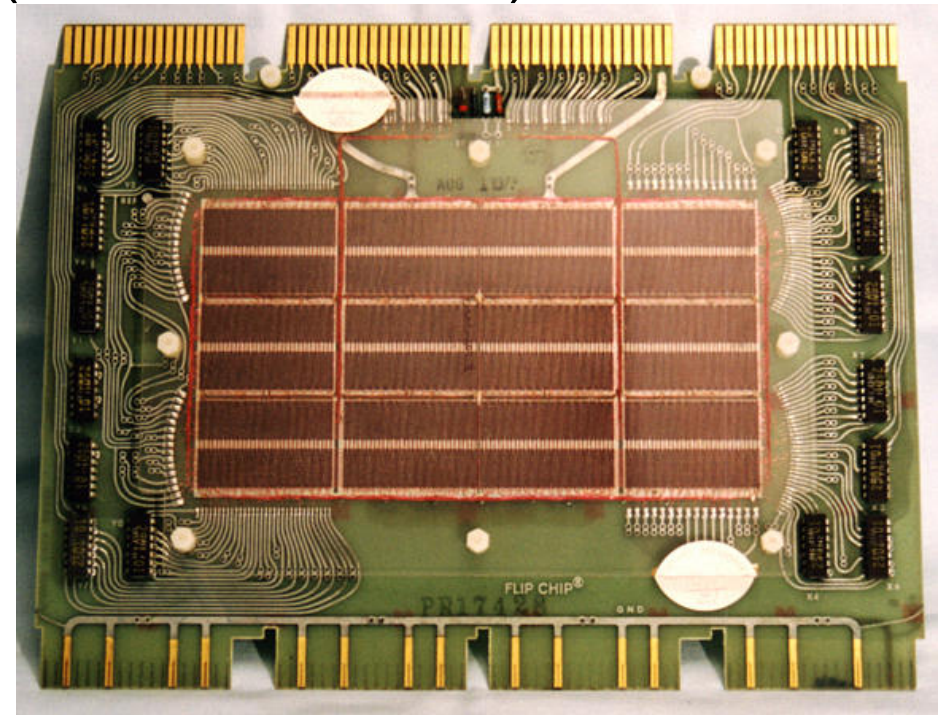




Core Memory

- Core memory was first large scale reliable main memory
 - invented by Forrester in late 40s/early 50s at MIT for Whirlwind project
- Bits stored as magnetization polarity on small ferrite cores threaded onto 2 dimensional grid of wires
- Coincident current pulses on X and Y wires would write cell and also sense original state (destructive reads)
- Robust, non-volatile storage
- Used on space shuttle computers until recently
- Cores threaded onto wires by hand (25 billion a year at peak production)
- Core access time $\sim 1\mu\text{s}$

DEC PDP-8/E Board,
4K words x 12 bits, (1968)





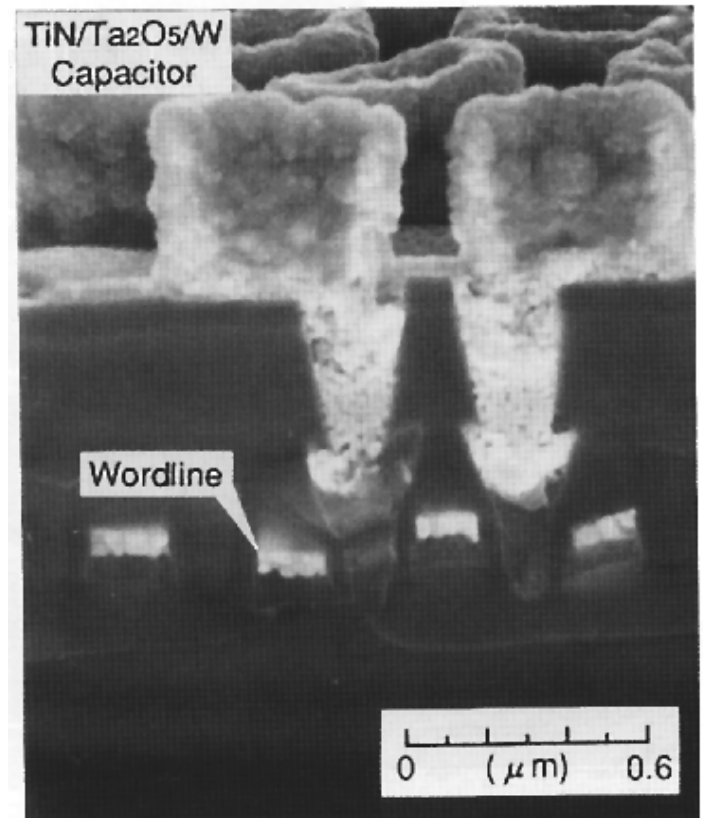
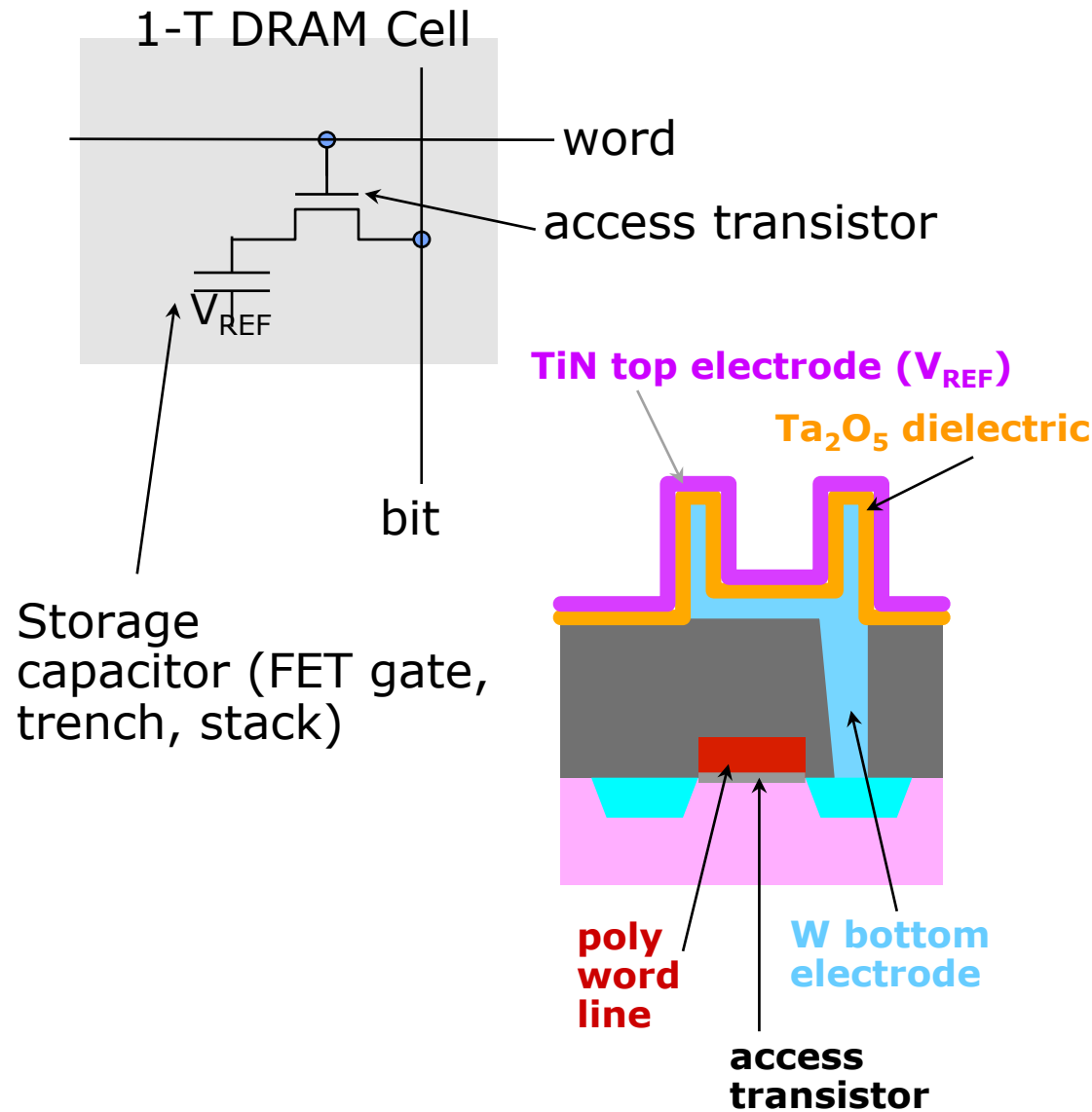
Semiconductor Memory

- 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



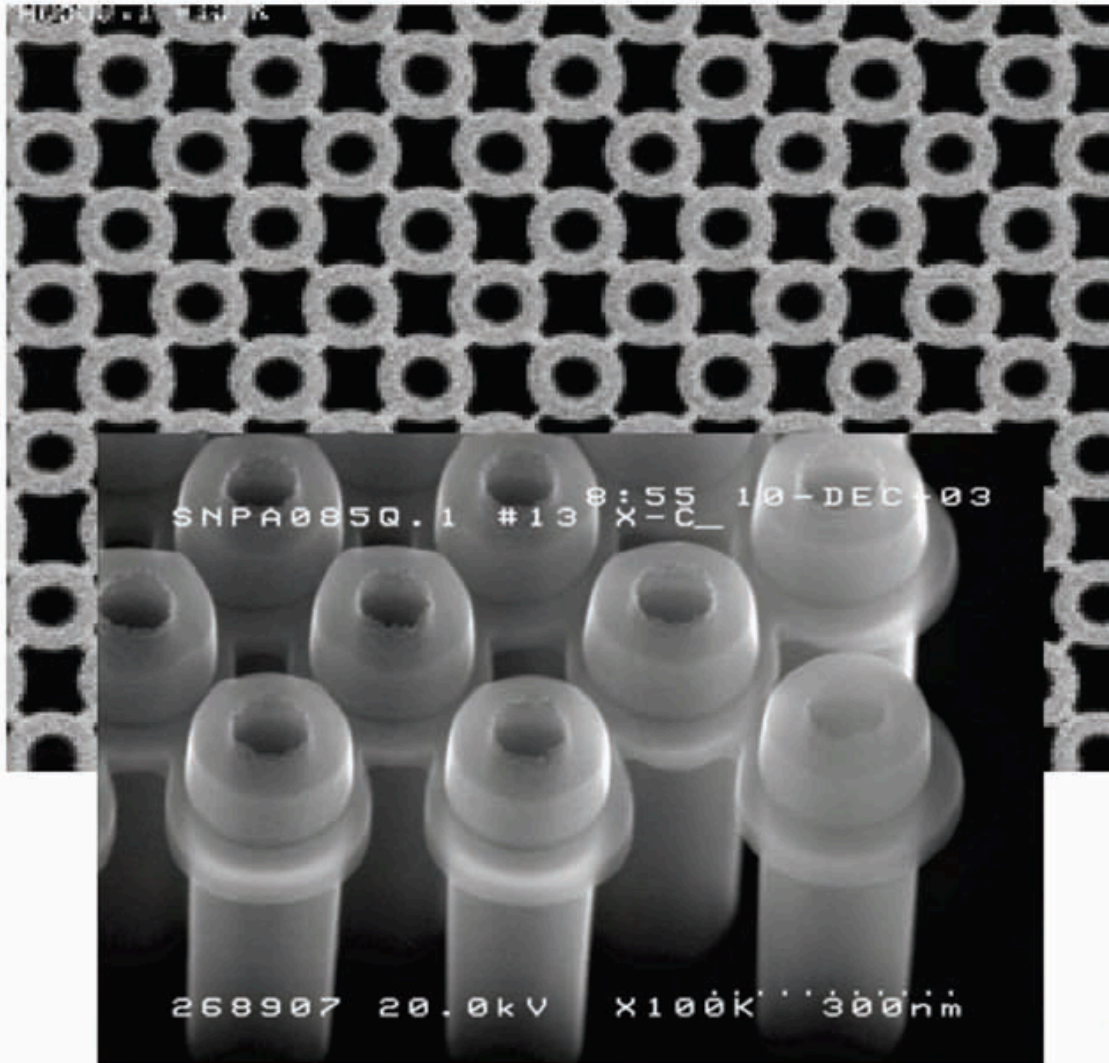
One Transistor Dynamic RAM

[Dennard, IBM]





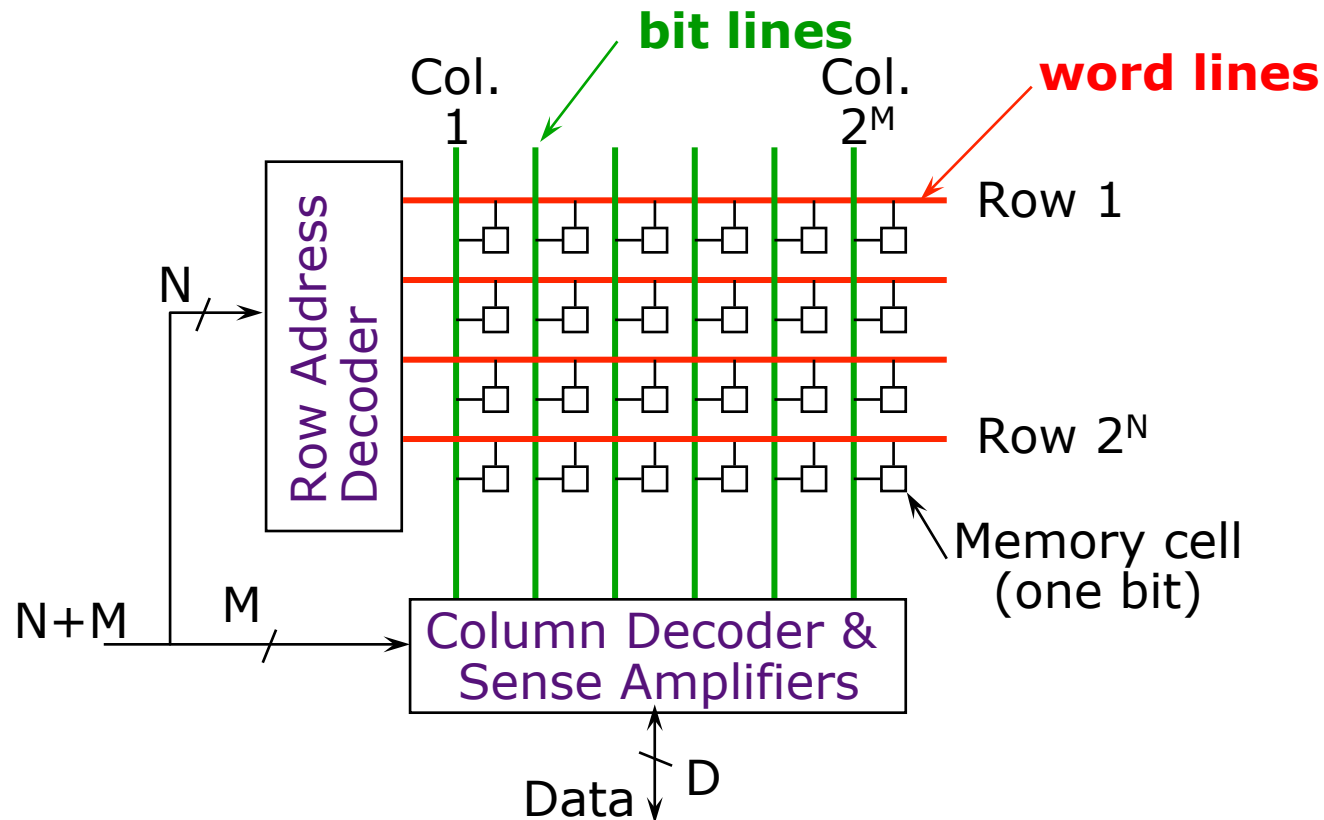
Modern DRAM Structure



[Samsung, sub-70nm DRAM, 2004]



DRAM Architecture



- Bits stored in 2-dimensional arrays on chip
- Modern chips have around 4 logical banks on each chip
 - each logical bank physically implemented as many smaller arrays



DRAM Operation

Three steps in read/write access to a given bank

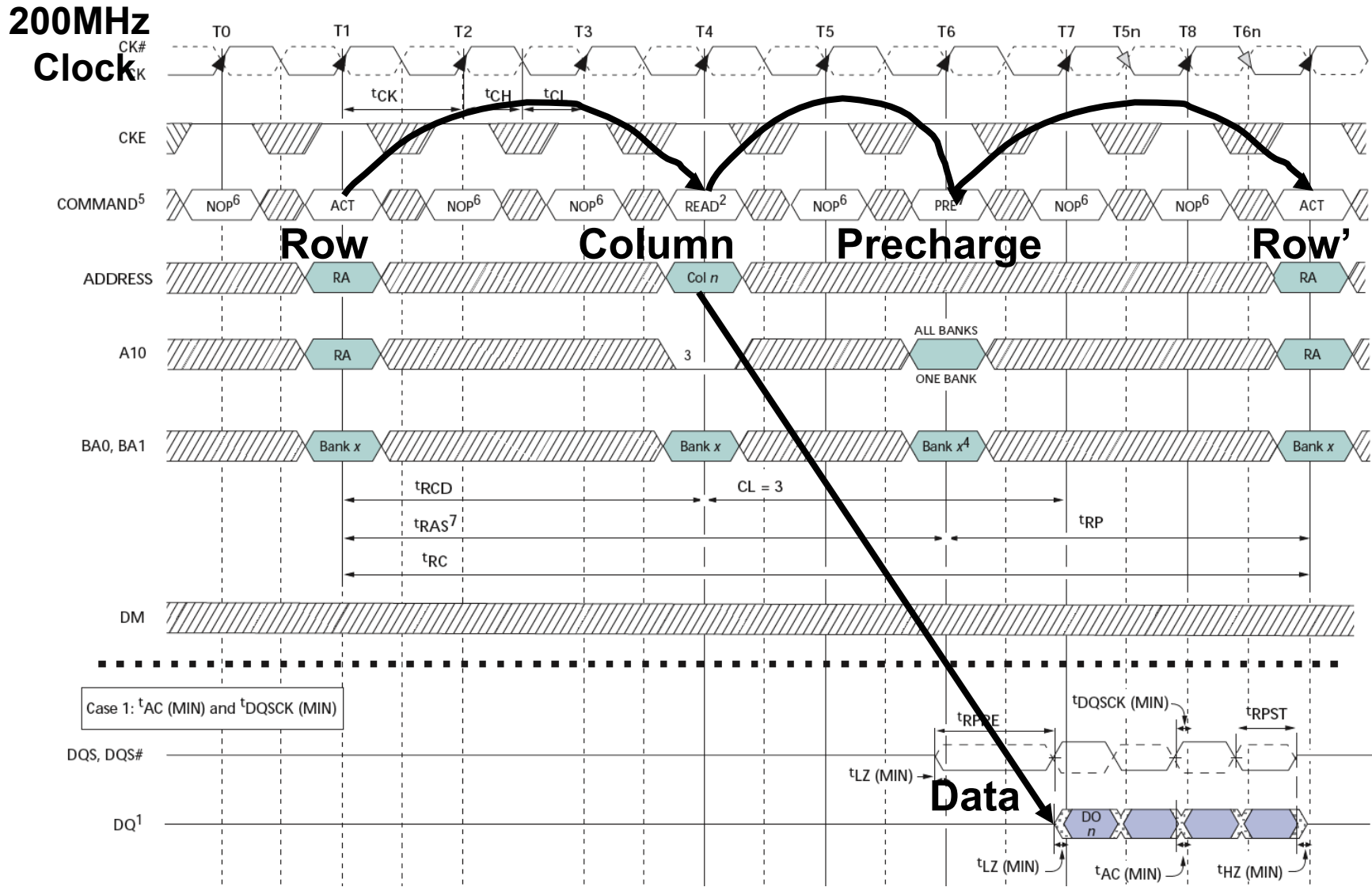
- Row access (RAS)
 - decode row address, enable addressed row (often multiple Kb in row)
 - bitlines share charge with storage cell
 - small change in voltage detected by sense amplifiers which latch whole row of bits
 - sense amplifiers drive bitlines full rail to recharge storage cells
- Column access (CAS)
 - decode column address to select small number of sense amplifier latches (4, 8, 16, or 32 bits depending on DRAM package)
 - on read, send latched bits out to chip pins
 - on write, change sense amplifier latches which then charge storage cells to required value
 - can perform multiple column accesses on same row without another row access (burst mode)
- Precharge
 - charges bit lines to known value, required before next row access

Each step has a latency of around 15-20ns in modern DRAMs

Various DRAM standards (DDR, RDRAM) have different ways of encoding the signals for transmission to the DRAM, but all share same core architecture



Double-Data Rate (DDR2) DRAM

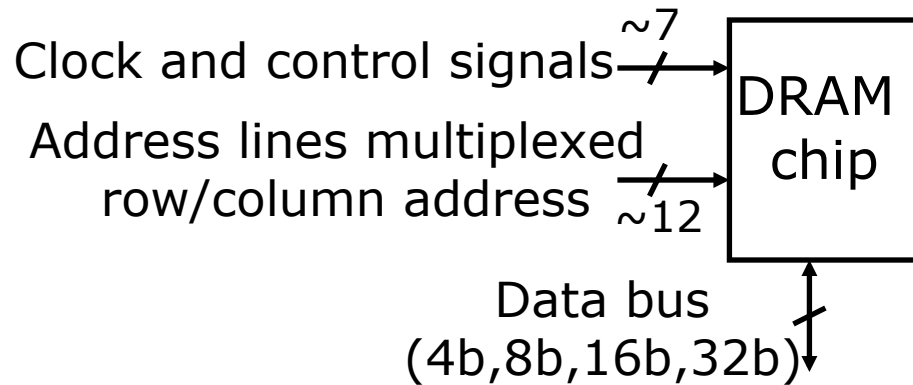


**400Mb/s
Data Rate**

[Micron, 256Mb DDR2 SDRAM datasheet]



DRAM Packaging



- DIMM (Dual Inline Memory Module) contains multiple chips with clock/control/address signals connected in parallel (sometimes need buffers to drive signals to all chips)
- Data pins work together to return wide word (e.g., 64-bit data bus using 16x4-bit parts)



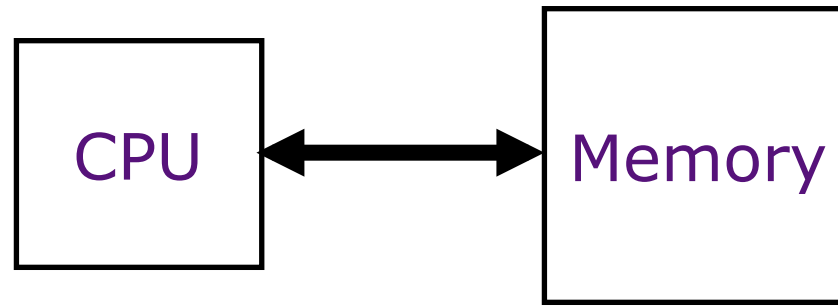
72-pin SO DIMM



168-pin DIMM



CPU-Memory Bottleneck

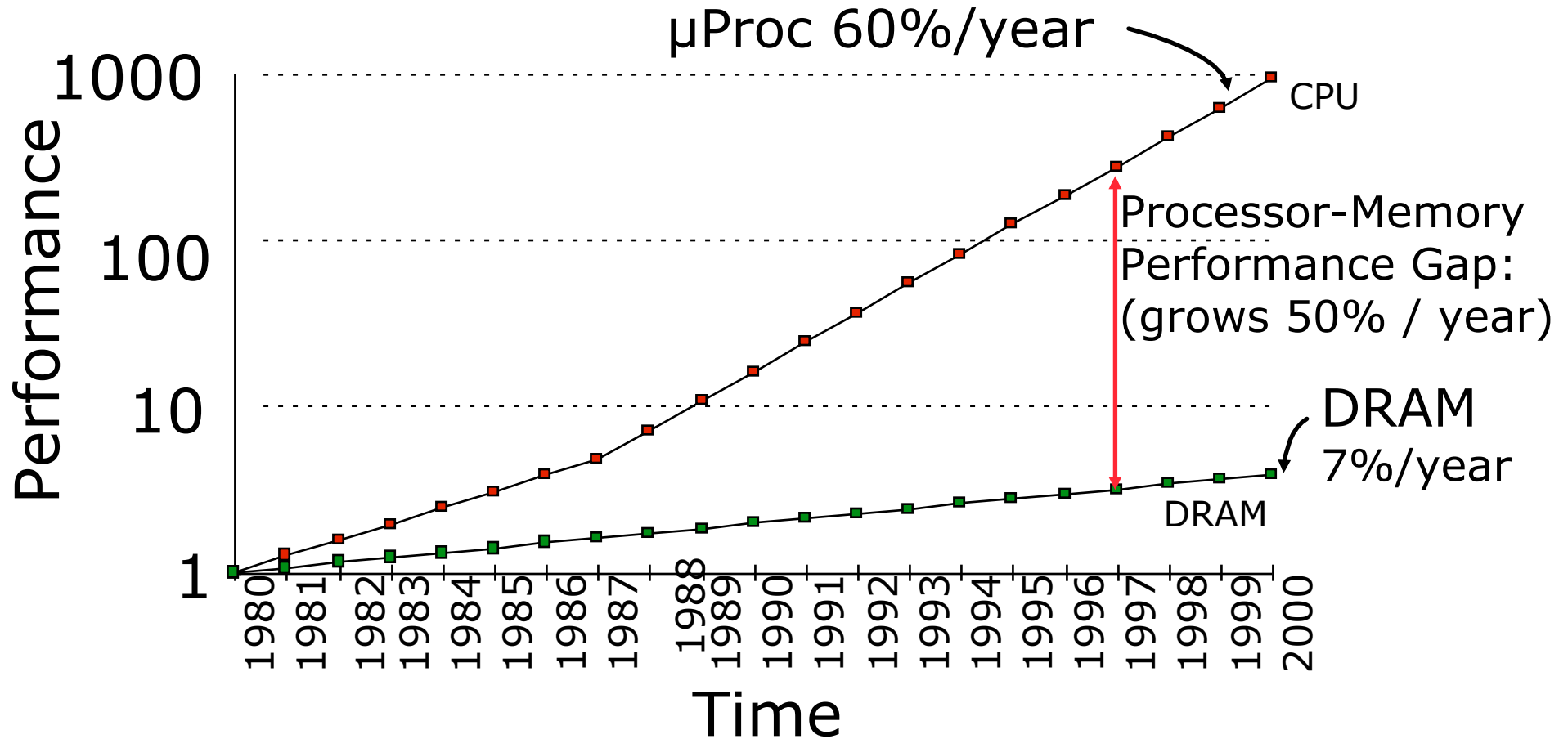


Performance of high-speed computers is usually limited by memory *bandwidth & latency*

- 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
(assuming MIPS RISC ISA)



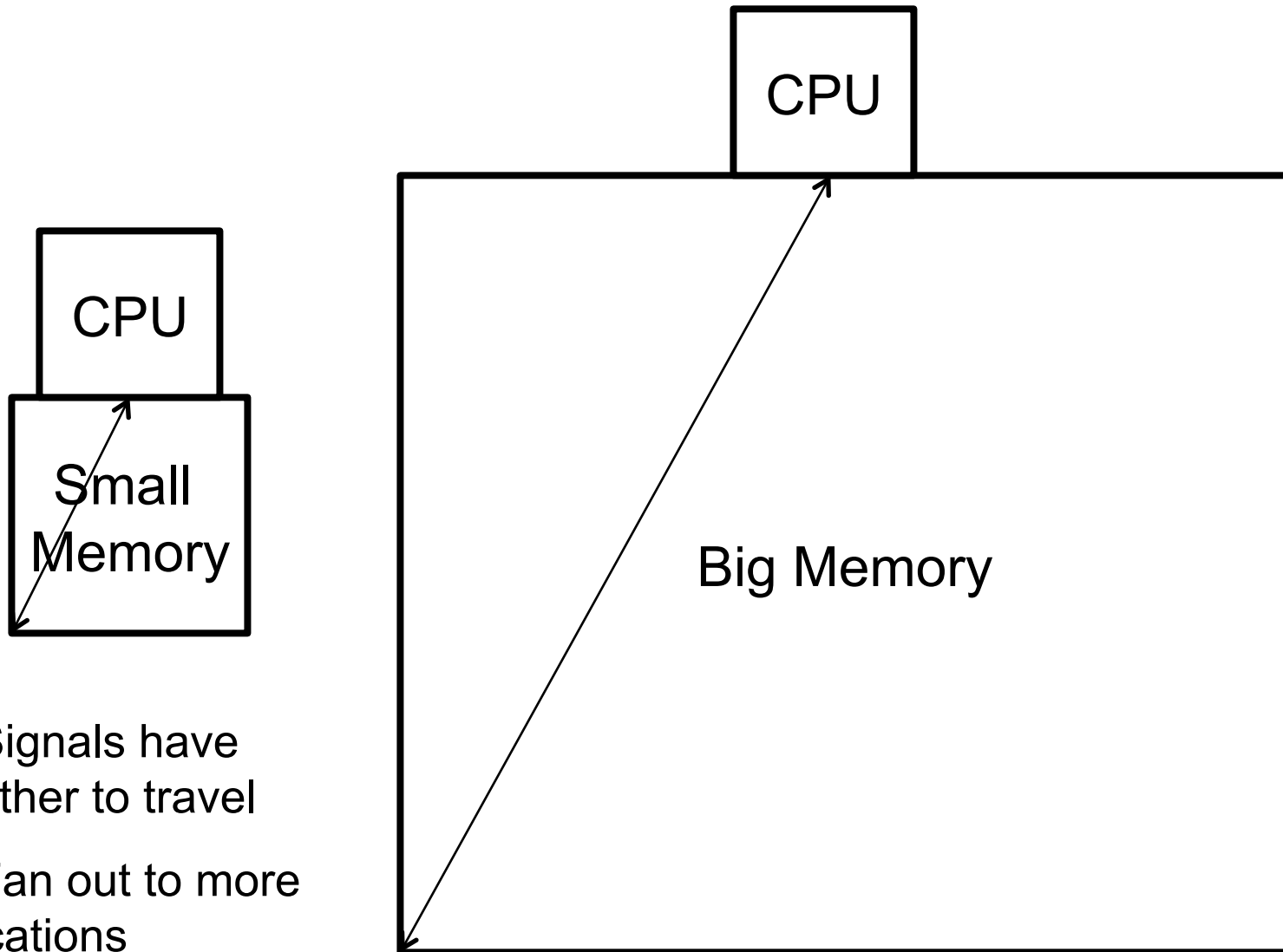
Processor-DRAM Gap (latency)



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



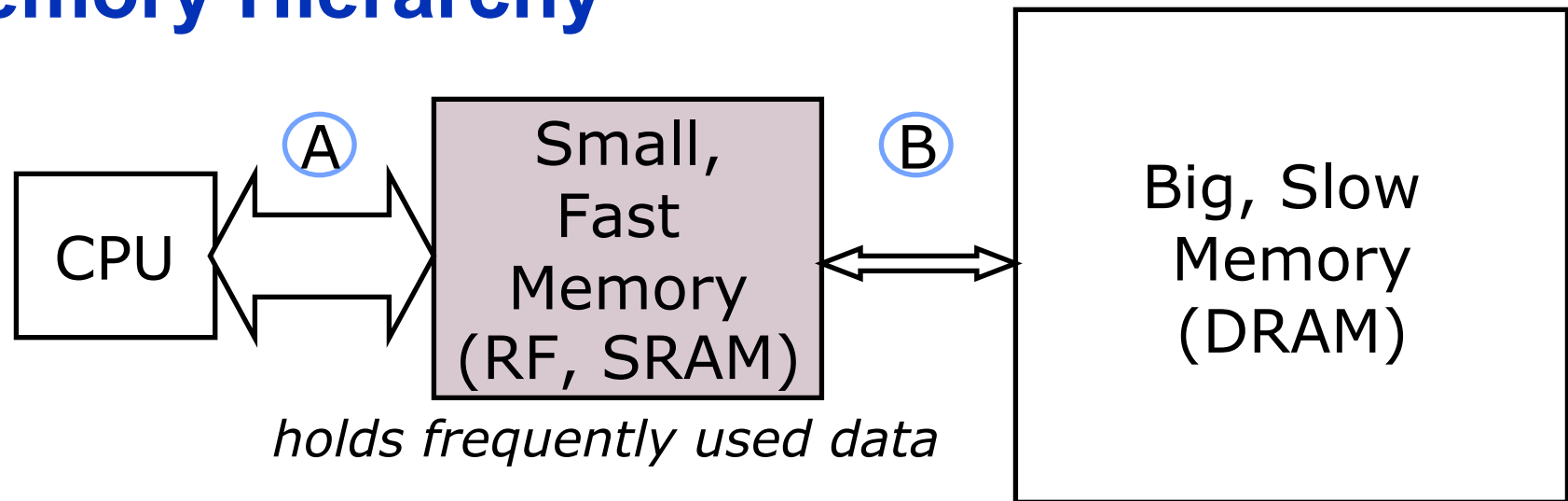
Physical Size Affects Latency



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



Memory Hierarchy



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

On a data access:

if data \in fast memory \Rightarrow low latency access (*SRAM*)

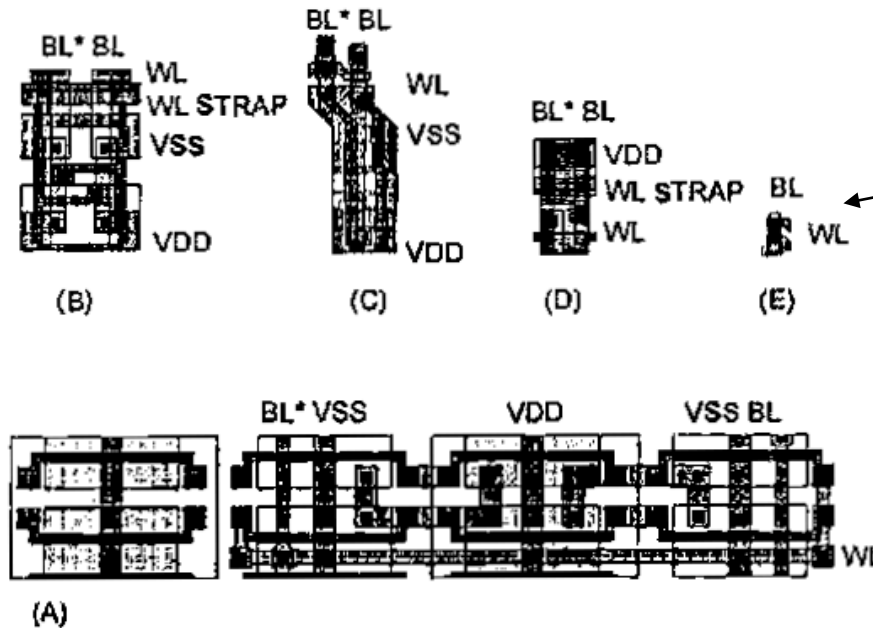
If data \notin fast memory \Rightarrow long latency access (*DRAM*)



Relative Memory Cell Sizes

On-Chip SRAM in logic chip

DRAM on memory chip



1 Memory cell in $0.5\mu\text{m}$ processes

- a) Gate Array SRAM
- b) Embedded SRAM
- c) Standard SRAM (6T cell with local interconnect)
- d) ASIC DRAM
- e) Standard DRAM (stacked cell)

[Foss, "Implementing Application-Specific Memory", ISSCC 1996]

Memory	Process	Cell size (μm^2)	Cell efficiency	Bits in $100\text{mm}^2(10^3)$	Gate size (μm^2)	Gate utilization	Gates in $100\text{mm}^2(10^3)$
Gate array SRAM	3-metal ASIC	370	80%	216	185	70%	378
Embedded SRAM	3-metal ASIC	67	70%	1045	185	70%	378
Standard SRAM	2-metal 6T local int.	43	65%	1512	245	40%	163
Embedded ASIC-DRAM	3-metal ASIC	23	60%	2609	185	70%	378
Standard DRAM	2-metal stacked cell	3.2	50%	15625	411	40%	97

Table 1: Memory and logic density for a variety of $0.5\mu\text{m}$ implementations.



CS152 Administrivia

- Class accounts available today
- Handed out in Section at 2pm

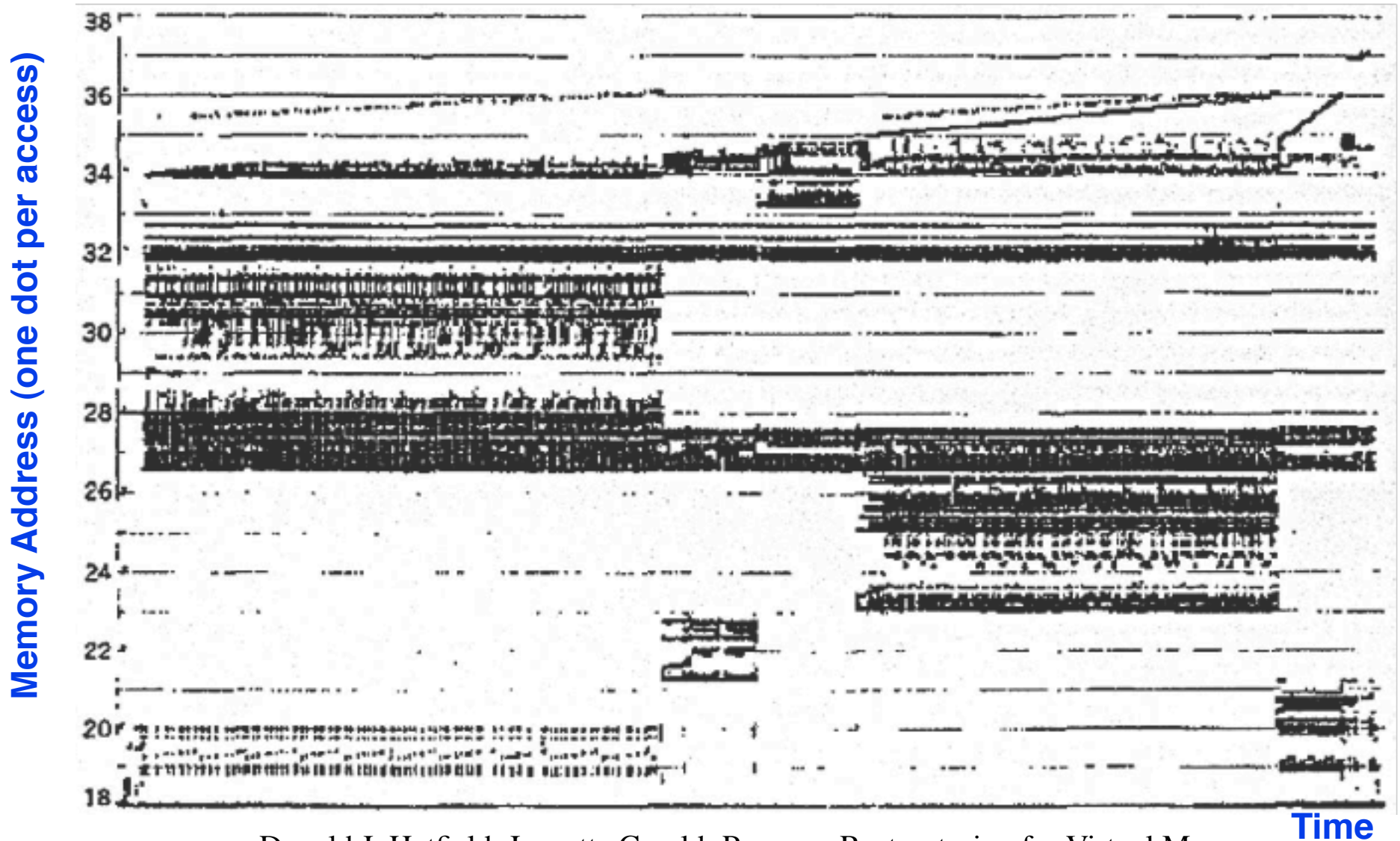


Management of Memory Hierarchy

- 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
- Larger/slower storage, e.g., main memory
 - Address usually computed from values in register
 - Generally implemented as a hardware-managed cache hierarchy
 - » hardware decides what is kept in fast memory
 - » but software may provide "hints", e.g., don't cache or prefetch



Real Memory Reference Patterns



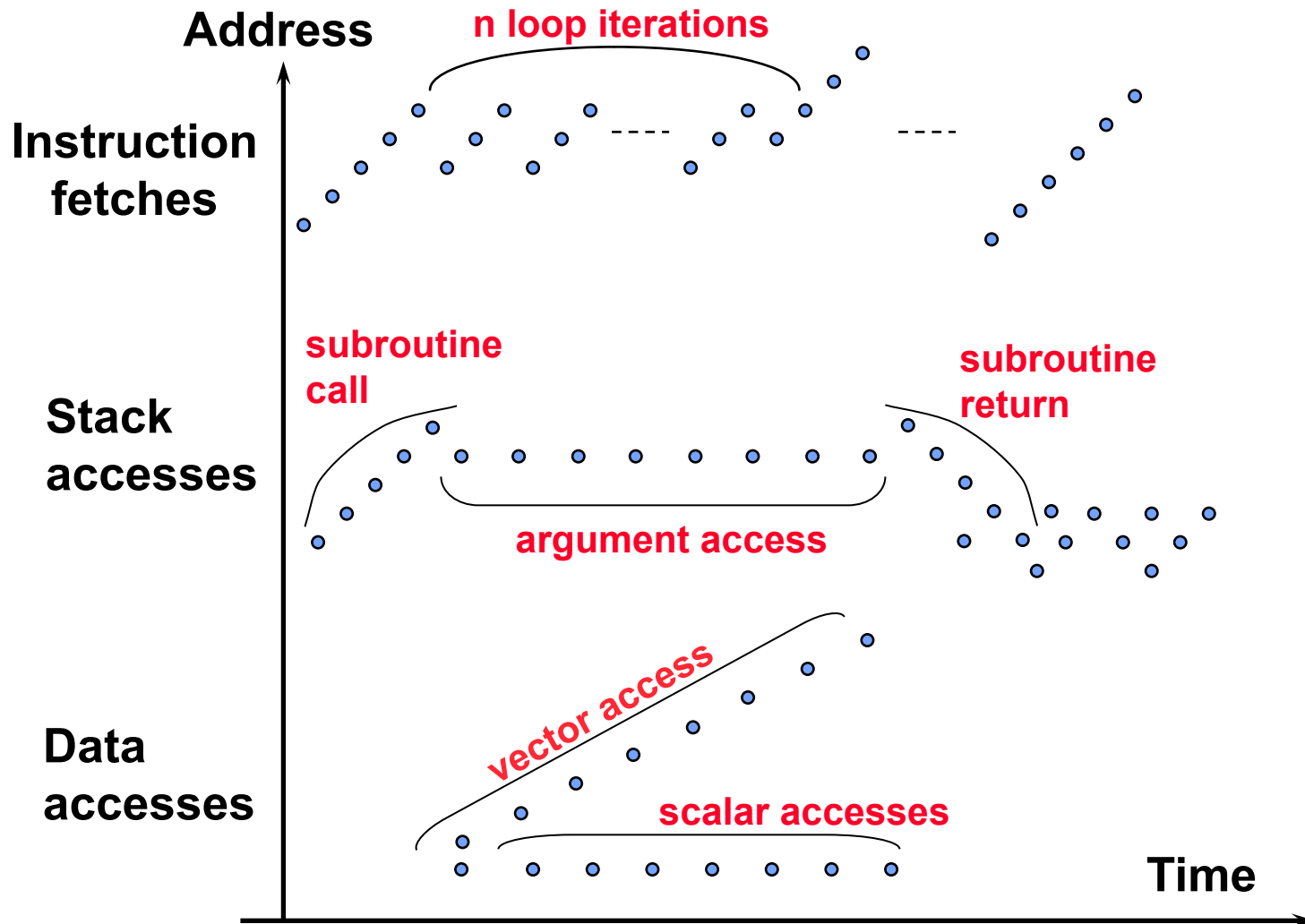
Donald J. Hatfield, Jeanette Gerald: Program Restructuring for Virtual Memory.
IBM Systems Journal 10(3): 168-192 (1971)

February 4, 2010

CS152, Spring 2010



Typical Memory Reference Patterns





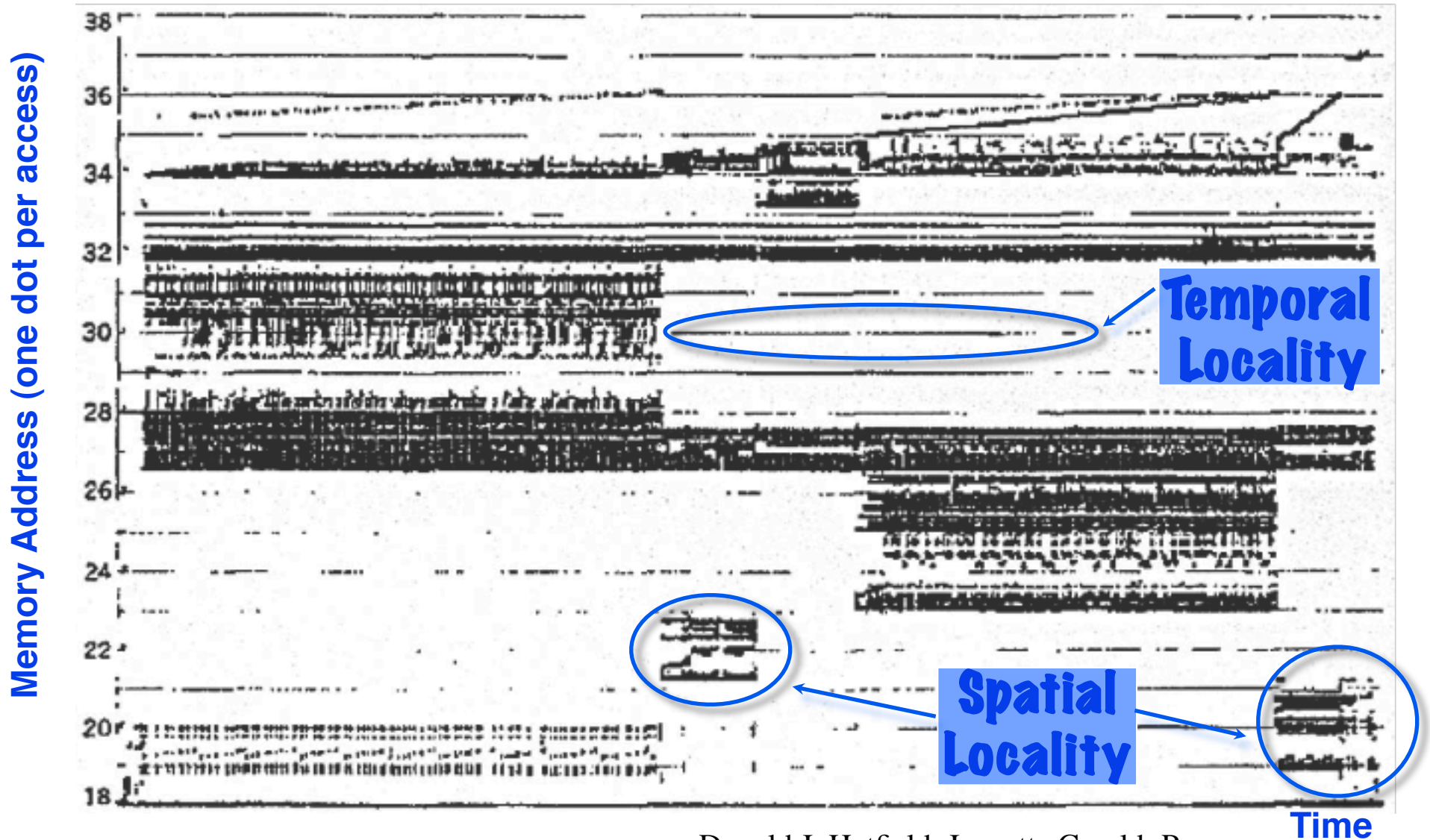
Common Predictable Patterns

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 Reference Patterns



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

February 4, 2010

CS152, Spring 2010



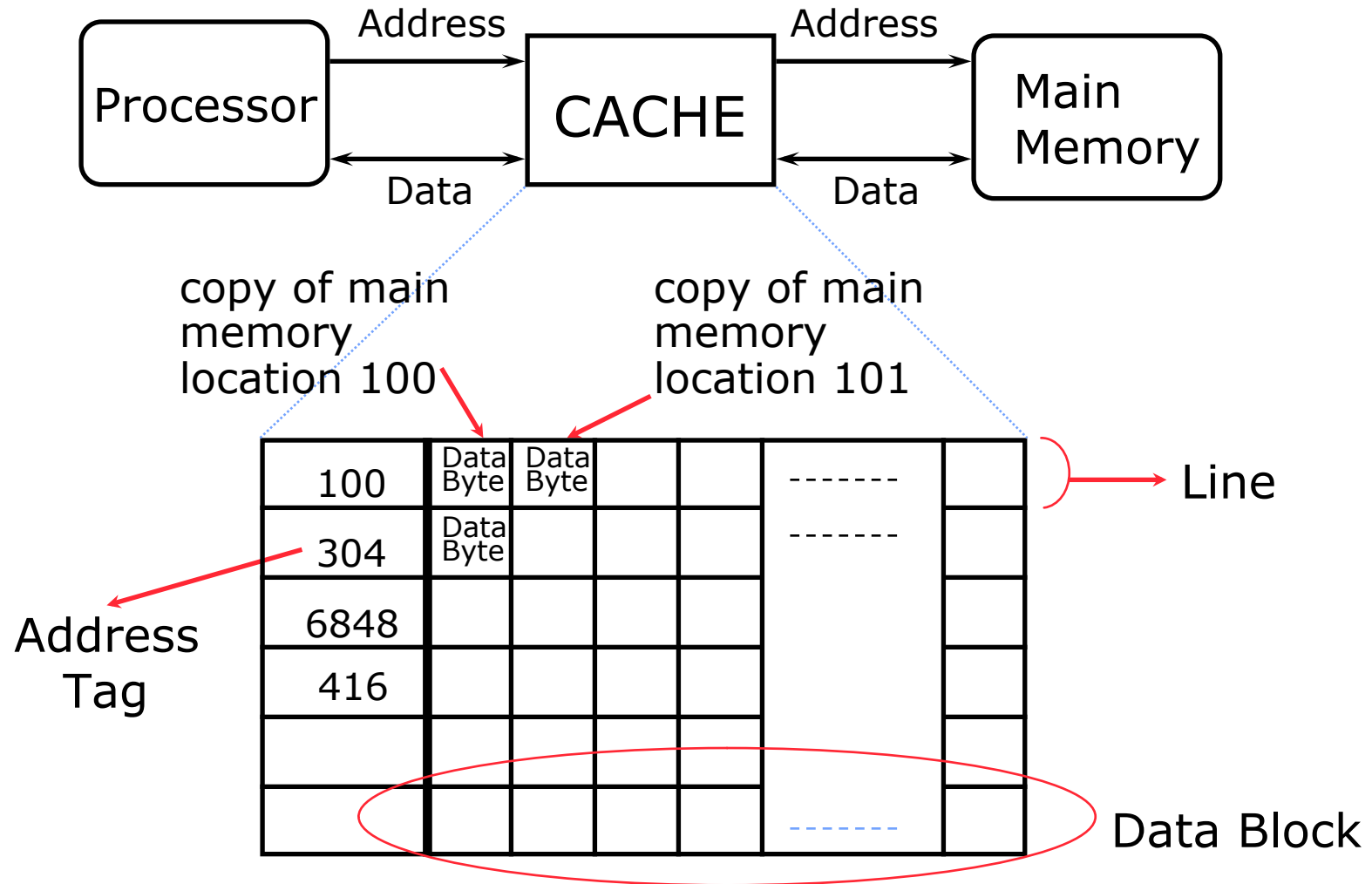
Caches

Caches exploit both types of predictability:

- Exploit temporal locality by remembering the contents of recently accessed locations.
- Exploit spatial locality by fetching blocks of data around recently accessed locations.



Inside a Cache





Cache Algorithm (Read)

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

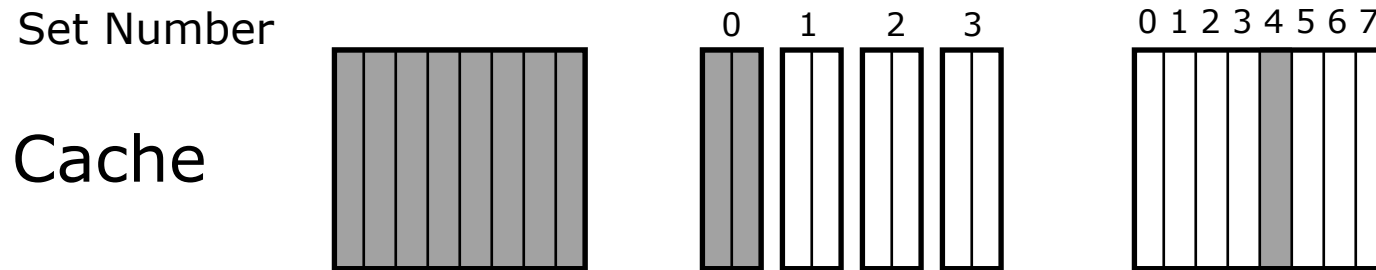
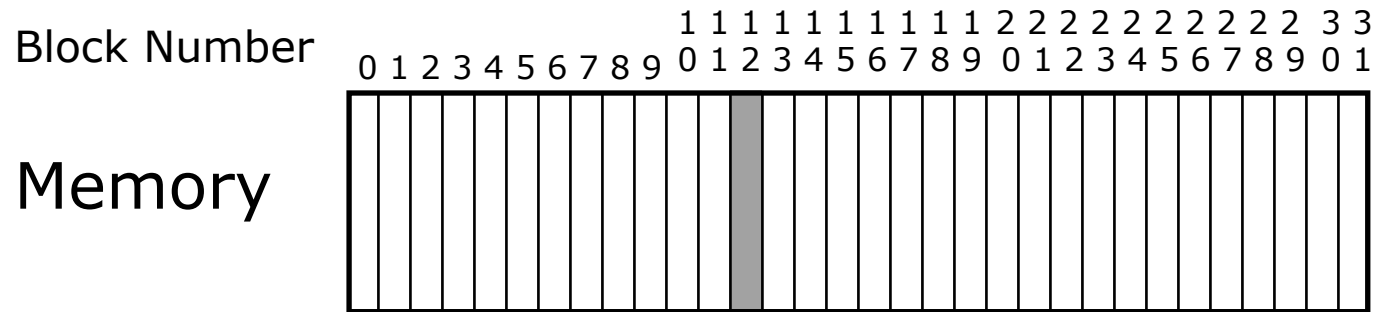
Wait ...

Return data to processor
and update cache

Q: Which line do we replace?



Placement Policy



Fully
Associative
anywhere

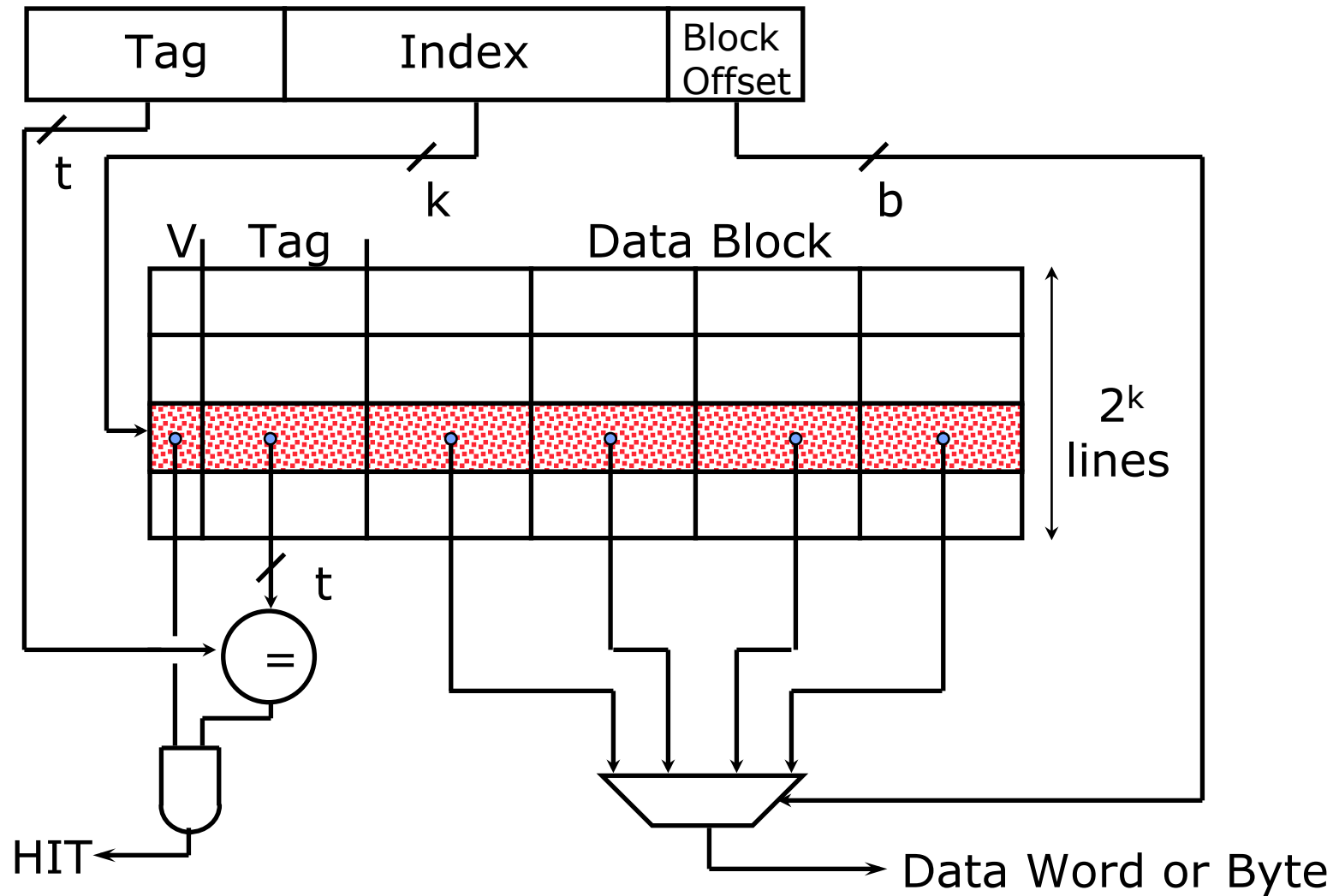
(2-way) Set
Associative
anywhere in
set 0
($12 \bmod 4$)

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

block 12
can be placed



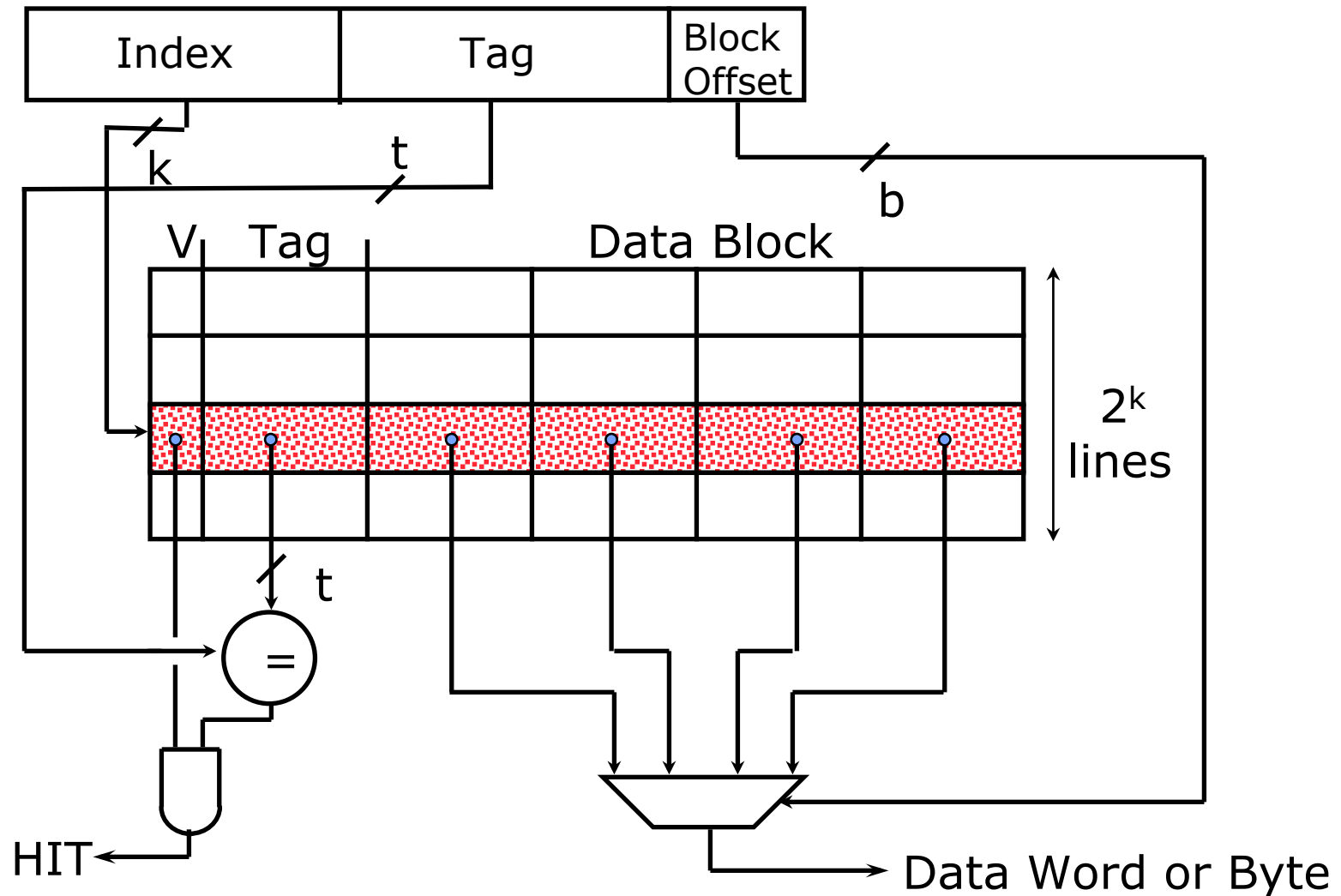
Direct-Mapped Cache





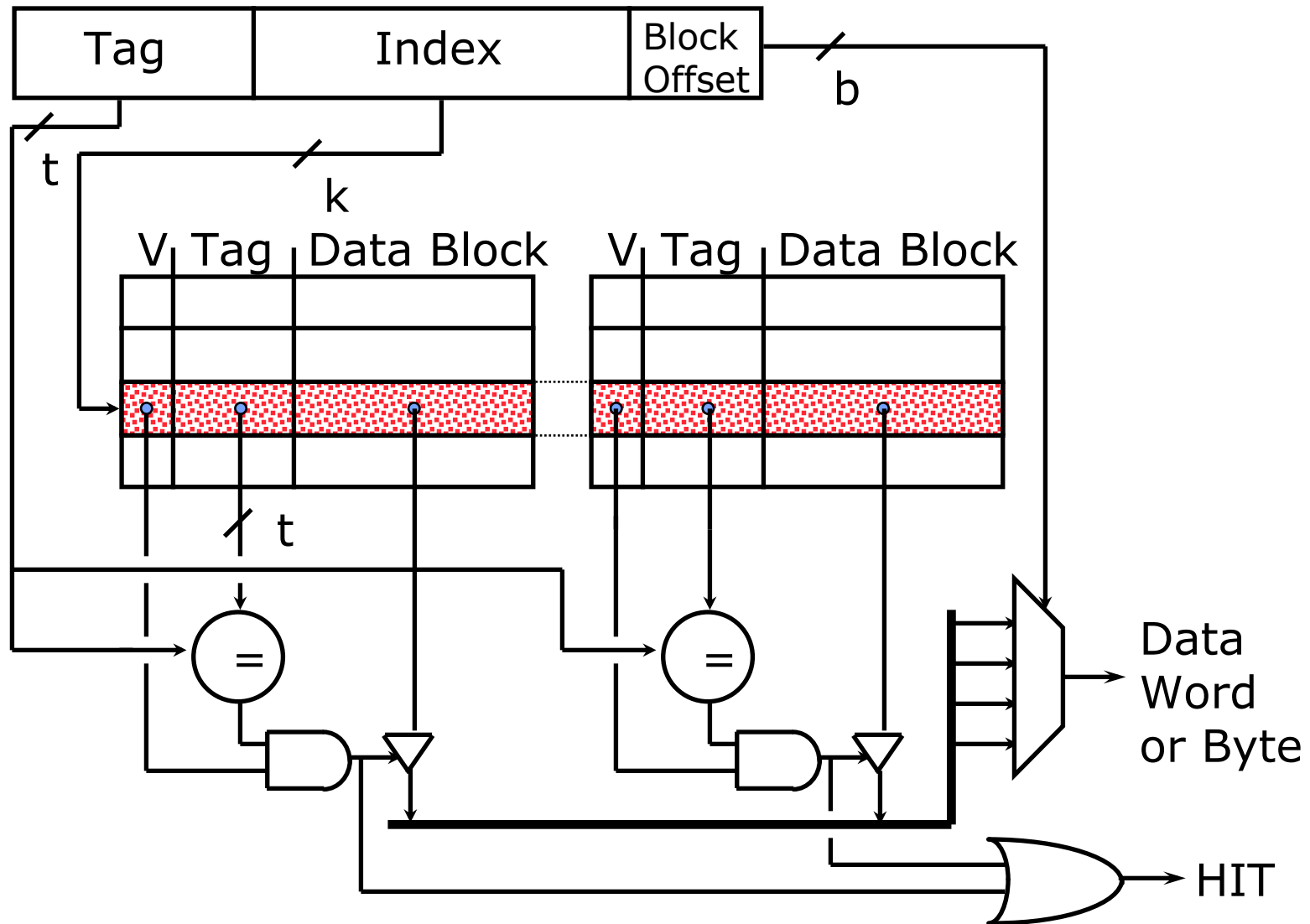
Direct Map Address Selection

higher-order vs. lower-order address bits



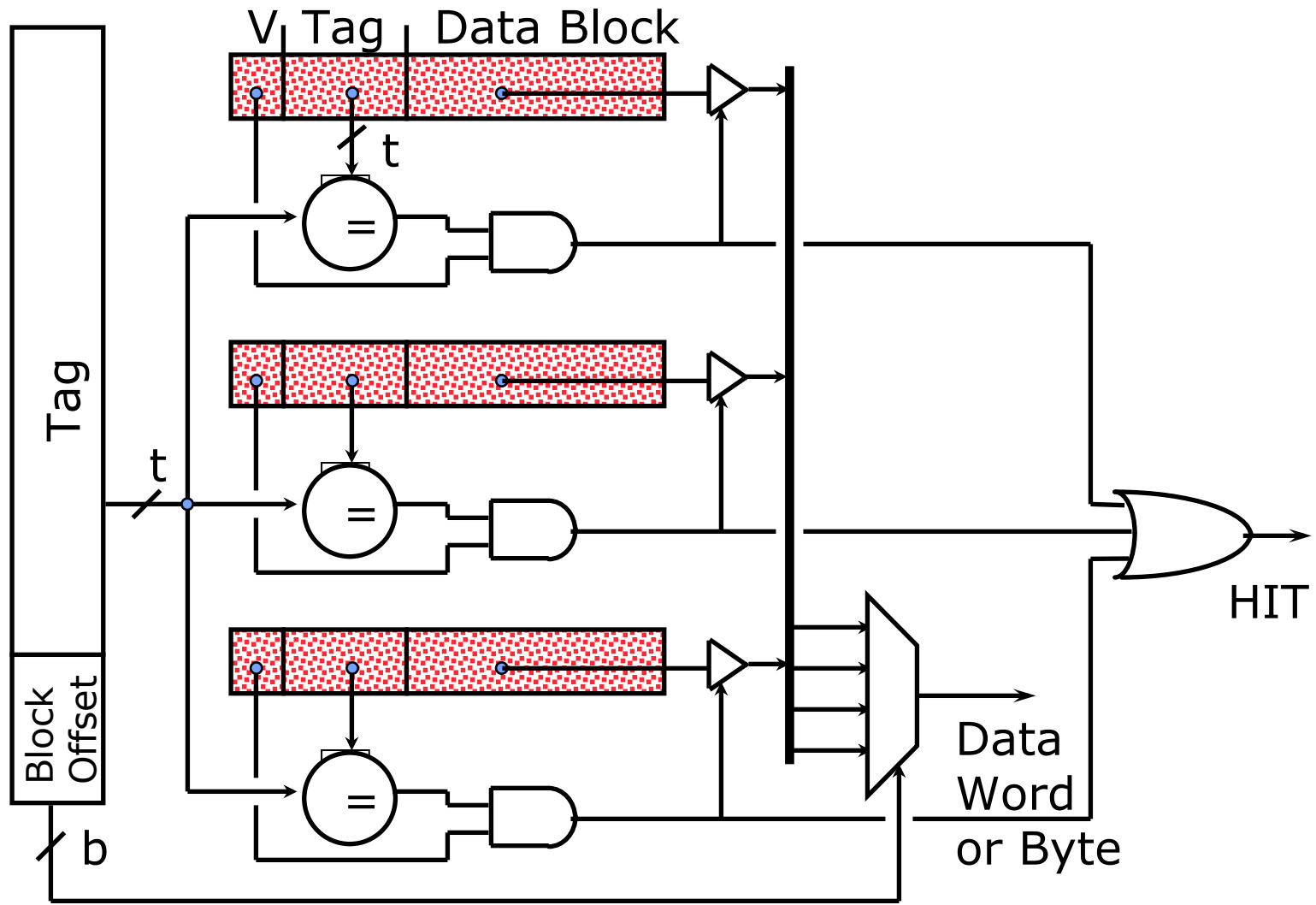


2-Way Set-Associative Cache





Fully Associative Cache





Replacement Policy

In an associative cache, which block from a set should be evicted when the set becomes full?

- 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 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

This is a second-order effect. Why?

Replacement only happens on misses



Acknowledgements

- These slides contain material developed and copyright by:
 - Arvind (MIT)
 - Krste Asanovic (MIT/UCB)
 - Joel Emer (Intel/MIT)
 - James Hoe (CMU)
 - John Kubiatowicz (UCB)
 - David Patterson (UCB)
- MIT material derived from course 6.823
- UCB material derived from course CS252