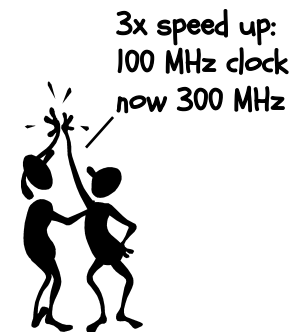


WHERE DOES THIS LEAVE US



Overall we can now nearly triple the clock rate. Instructions have a throughput of one-per-clock with the following caveats:

1. Taken branches take 2 cycles.
2. Loads and store take 1 cycle.



You can pipeline an ARM CPU even more. There exist ARM implementations with 7, 8, and 9 pipeline stages. But the overhead of bypass paths and stall cases increase.

REALITY VS SPECMANSHIP



Assuming approximately 10% of instructions executed are branches, and of those 80% of the time they are taken, and 12.5% of instruction executed are loads or stores, what sort of real speed up do we expect?

$$\text{Perf}_{\text{before}} = (100) * 1 = 100 \text{ Clocks} * 10 * 10^{-9} \text{ sec/clock} = 1000 * 10^{-9} \text{ secs}$$

$$\text{Perf}_{\text{after}} = (10)(0.8) * 2 + 12.5 * 2 + 79.5 * 1 = 120.5 \text{ Clocks}$$

$$120.5 * 3.333 * 10^{-9} \text{ sec/clock} = 401.666 * 10^{-9} \text{ secs}$$

$$\text{Speedup} = \frac{\text{Perf}_{\text{before}}}{\text{Perf}_{\text{after}}} = 1000 / 401.666 = 2.490 \times$$

NEXT TIME



It appears memory access time is our real bottleneck. What tricks can be applied to improving CPU performance in this case?

- Interleaving
- Block-transfers
- Caching

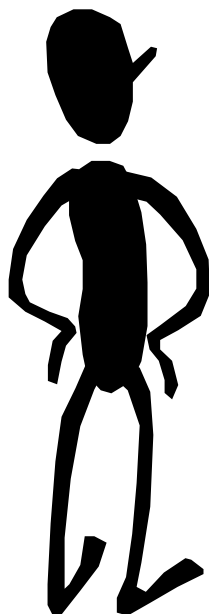


MEMORY HIERARCHY + CACHING



It makes me look faster,
don't you think?

Still in your Halloween
costume?

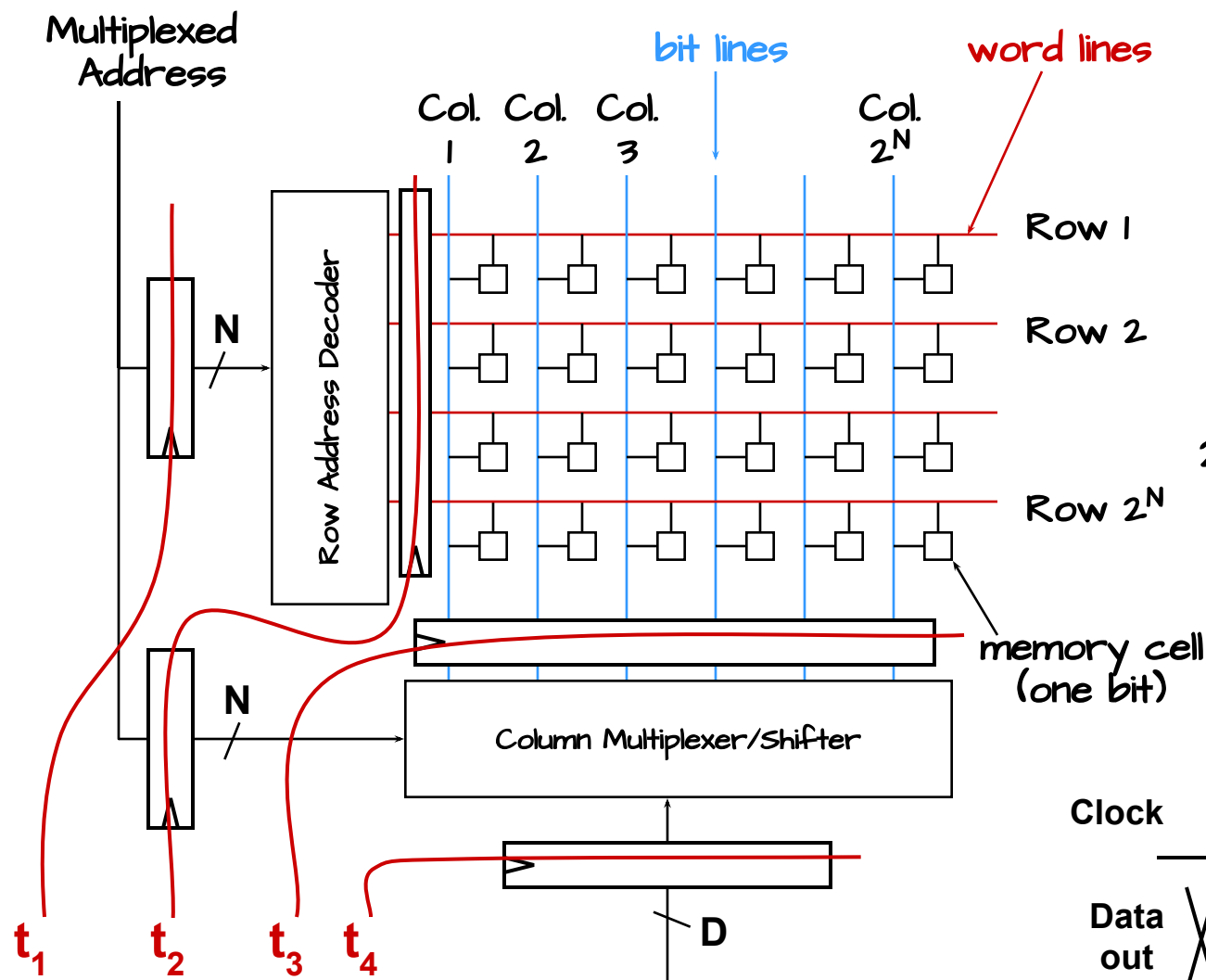


- Memory Flavors
- Principle of Locality
- Memory Hierarchies
- Caches
- Associativity
- Write-through
- Write-back

Midterm #2 on Wednesday

- Open notes & internet
- Must be in SNO14, or previously arranged and monitored location

TRICKS FOR INCREASING THROUGHPUT



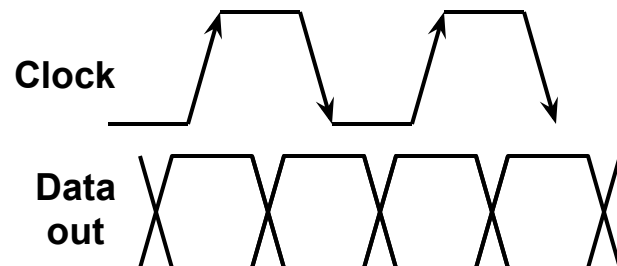
The first thing that should pop into your mind when asked to speed up a digital design...

PIPELINING

Synchronous DRAM (SDRAM)

20ns reads and writes
(\$5 per Gbyte)

Double Data Rate Synchronous DRAM (DDR)



ANOTHER TRICK



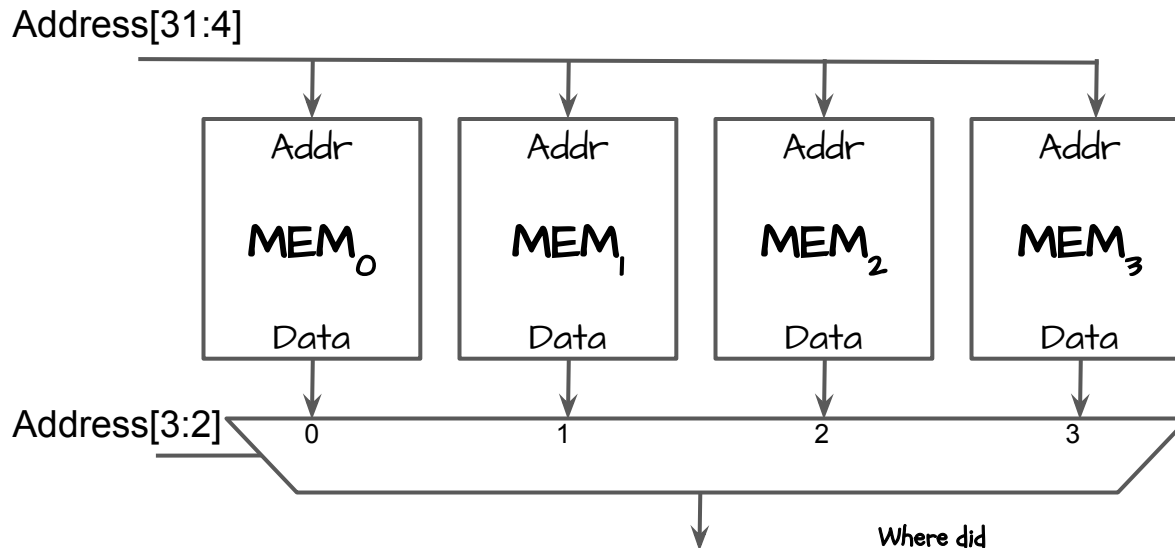
The second thing that should try when asked to speed up a digital design...

Interleaving

Accessing 4 memories at the same time has 4x the throughput. Also, every 4th word is in a different memory.

A limitation of both pipelining and interleaving is their assumption that addresses are sequential!

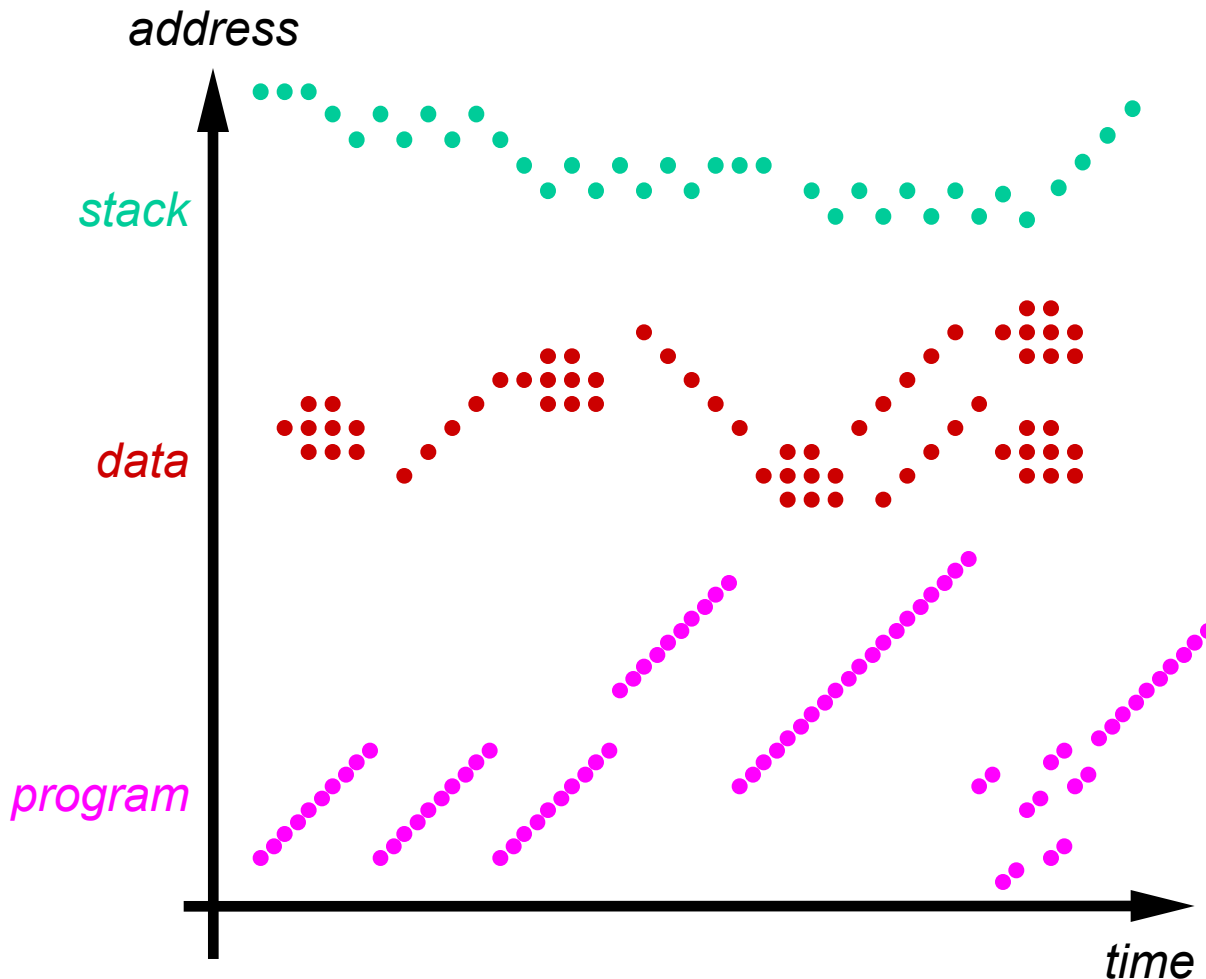
Which is approximately true!



If only the lower order addresses change, we need only wait the T_{pd} of the mux.



TYPICAL MEMORY REFERENCE PATTERNS



MEMORY TRACE -

A temporal sequence of memory references (addresses) from a real program.

TWO KEY OBSERVATIONS:

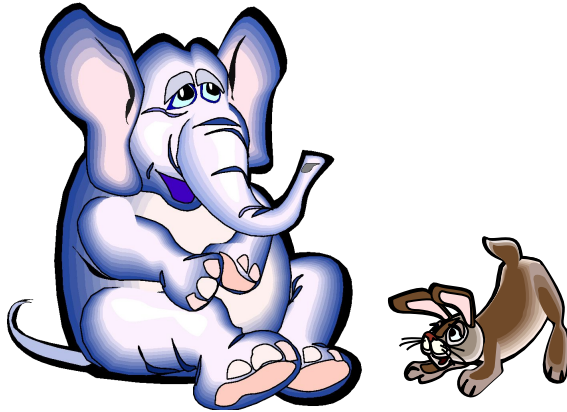
TEMPORAL LOCALITY -

If an item is referenced, it will tend to be referenced again soon

SPATIAL LOCALITY -

If an item is referenced, nearby items will tend to be referenced soon.

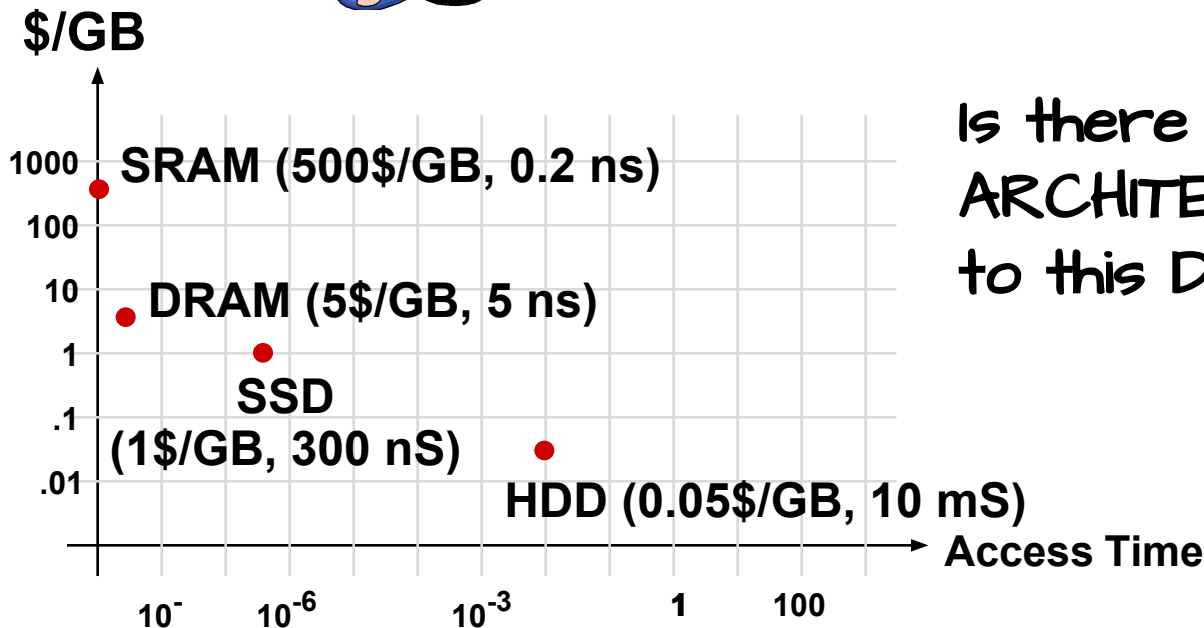
ALL MEMORIES AREN'T CREATED EQUAL



Quantity vs Speed...

Memory systems can be either:

- BIG and SLOW...
- or
- SMALL and FAST.



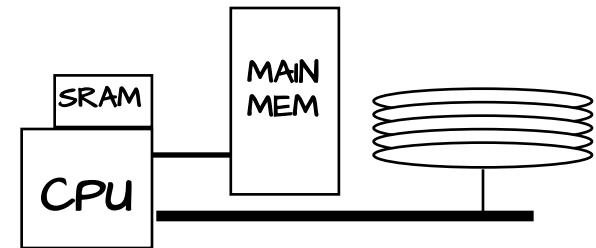
Is there an
ARCHITECTURAL solution
to this DILEMMA?

EXPLOITING THE MEMORY HIERARCHY



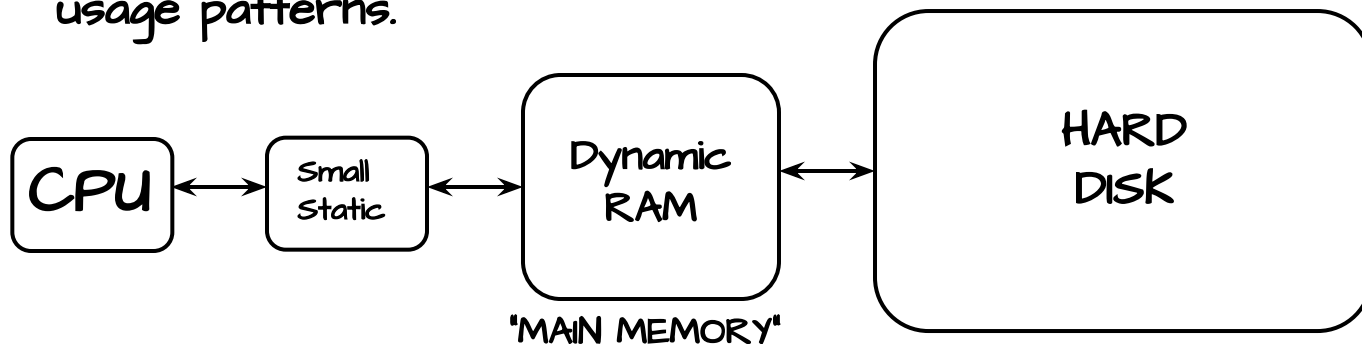
Approach 1 (Cray, others): Expose Hierarchy

- Registers, Main Memory, Disk each available as storage alternatives;
- Tell programmers: "Use them wisely"



Approach 2: Hide Hierarchy

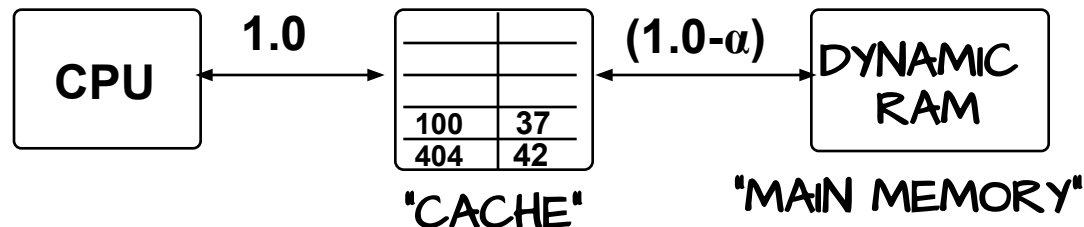
- Programming model: SINGLE kind of memory, single address space.
- Machine AUTOMATICALLY assigns locations to fast or slow memory, depending on usage patterns.





THE CACHE CONCEPT:

PROGRAM-TRANSPARENT MEMORY HIERARCHY



Cache contains TEMPORARY COPIES of selected main-memory locations... eg. $\text{Mem}[100] = 37$

GOALS:

1) Improve the **average access** time

α HIT RATIO: Fraction of refs found in CACHE.


$(1 - \alpha)$ MISS RATIO: Remaining references.

$$t_{\text{ave}} = \alpha t_c + (1 - \alpha)(t_c + t_m) = t_c + (1 - \alpha)t_m$$

2) Transparency (compatibility, programming ease)

Challenge:
Make the hit ratio, α , as high as possible.

Why, on a miss, do I incur the access penalty for both main memory and cache?



HOW HIGH OF A HIT RATIO?



Suppose we can easily build an on-chip static memory with a 800 ps access time, but the fastest dynamic memories that we can buy for main memory have an average access time of 10 ns. How high of a hit rate do we need to sustain an average access time of 1 ns?

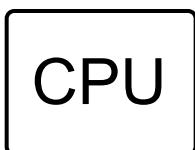
$$\text{Solve for } \alpha: t_{ave} = t_c + (1-\alpha)t_m$$

$$\alpha = 1 - (t_{ave} - t_c)/t_m = 1 - (1-0.8)/10 = 98\%$$

Wow, caches really need to be good! And they are!



BASIC CACHE ALGORITHM



ON REFERENCE TO $\text{Mem}[X]$: Look for X among cache tags...

HIT: $X == \text{TAG}(i)$, for some cache line i

READ: return $\text{DATA}(i)$

WRITE: change $\text{DATA}(i)$;
Start Write to $\text{Mem}(X)$

"X" here is a memory address.



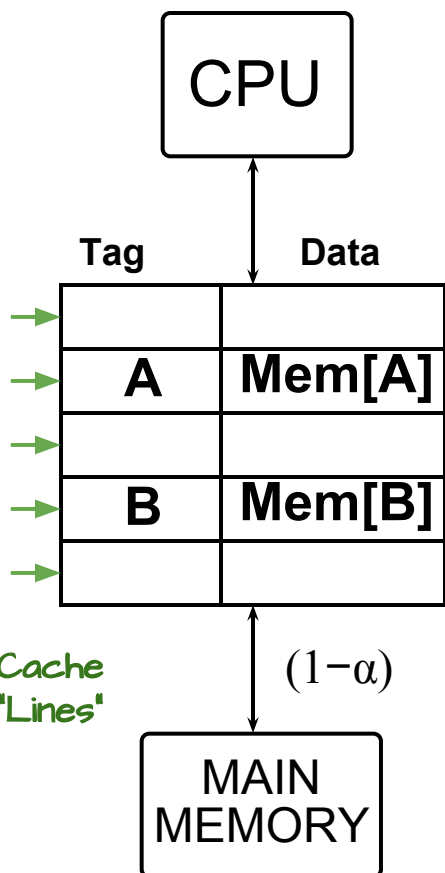
MISS: X not found in any TAG of the cache

REPLACEMENT SELECTION:

Select some LINE k to hold $\text{Mem}[X]$ (**Allocation**)

READ: Read $\text{Mem}[X]$
Set $\text{TAG}(k)=X$, $\text{DATA}(k)=\text{Mem}[X]$

WRITE: Start Write to $\text{Mem}(X)$
Set $\text{TAG}(k)=X$, $\text{DATA}(k)=\text{new Mem}[X]$



Cache-lines might contain multiple sequential words from memory, thus amortizing the number of tag bits per data bits.

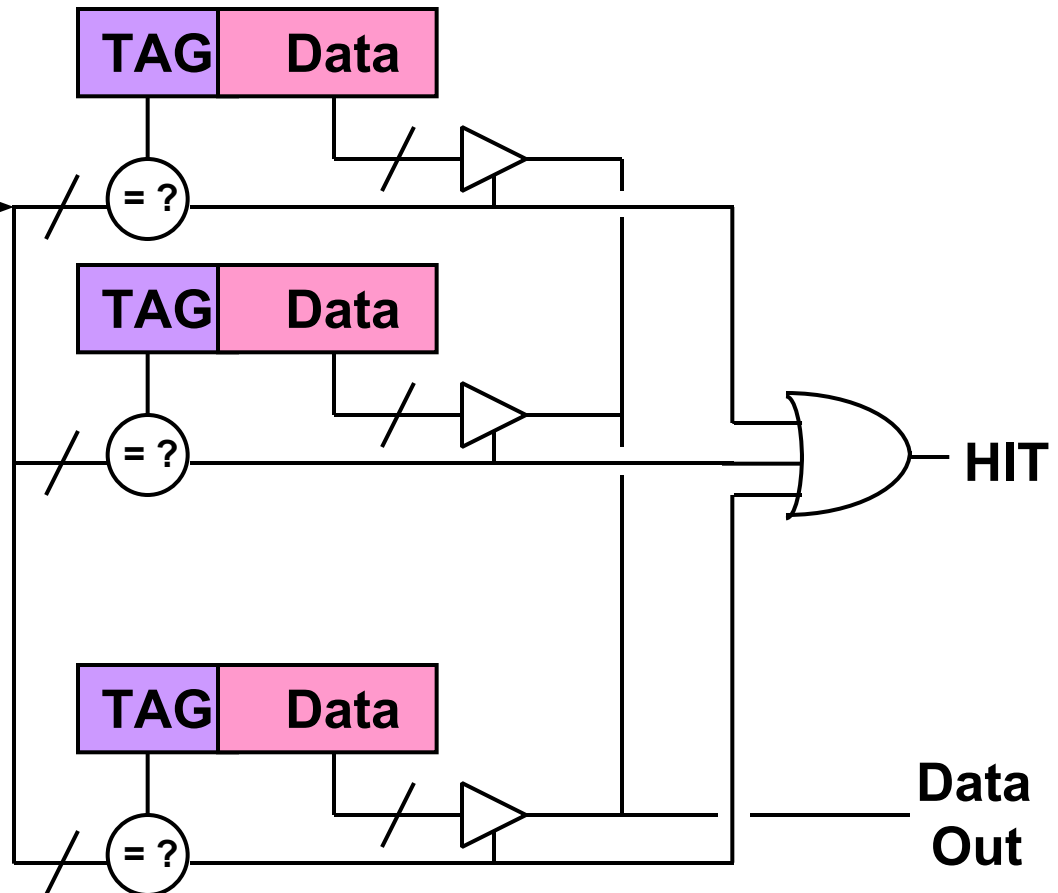
SEARCHING FOR TAGS



Associativity: Degree of parallelism used to lookup tags

Fully-Associative Cache:

Incoming
Address



The extreme in
associativity:

All TAGS are searched
in parallel

Data items from **any**
address can be located in
any cache line

THE OTHER EXTREME

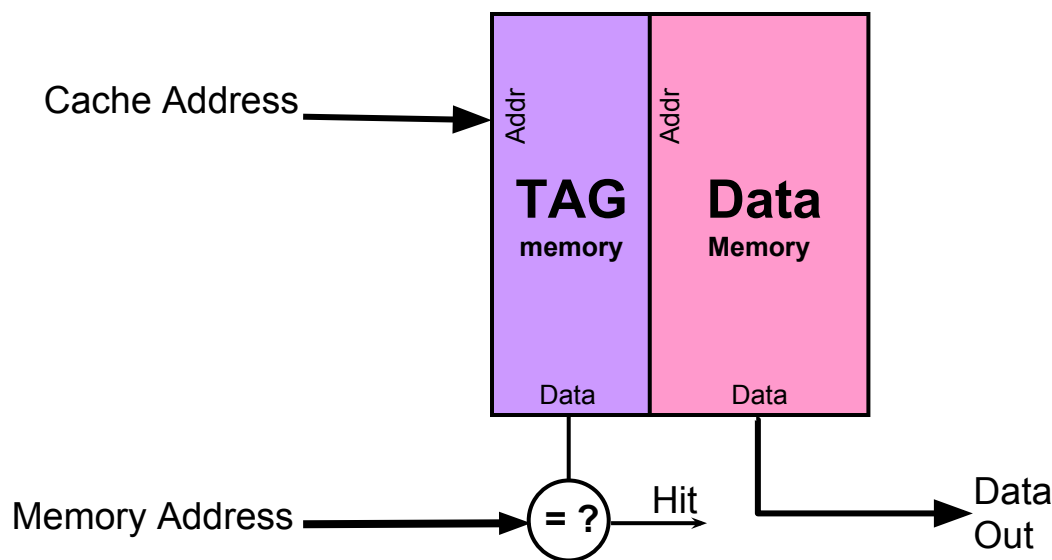


Direct-mapped: If it is in cache it is in exactly one place

Non-associative or "one-way" associative. No parallelism.

Uses only one **comparator** and ordinary RAM for tags:

Low-cost leader:



Direct-mapped caches require a means for translating "Memory Addresses" to "Cache Addresses". A simple **hash** function.

DIRECT-MAPPED EXAMPLE



With 8-byte lines, 3 low-order bits determine the byte within the line.

With 4 cache lines, the next 2 bits can be used to decide which line to use

$$1024_{10} = 1000000\mathbf{00}000_2 \rightarrow \text{line} = 00_2 = 0_{10}$$

$$1000_{10} = 011111\mathbf{01}000_2 \rightarrow \text{line} = 01_2 = 1_{10}$$

$$1040_{10} = 100000\mathbf{10}000_2 \rightarrow \text{line} = 10_2 = 2_{10}$$

Cache			
Line 0	1024	44	99
Line 1	1000	17	23
Line 2	1040	1	4
Line 3	1016	29	38
Tag		Data	

Memory

1000	17
1004	23
1008	11
1012	5
1016	29
1020	38
1024	44
1028	99
1032	97
1036	25
1040	1
1044	4



DIRECT-MAPPED MISS

What happens when we now ask for address 1008?

$$1008_{10} = 0111111\mathbf{10}000_2 \rightarrow \text{line} = 10_2 = 2_{10}$$

but earlier we put 1040 there...

$$1040_{10} = 100000\mathbf{10}000_2 \rightarrow \text{line} = 10_2 = 2_{10}$$

Cache			
Line 0	1024	44	99
Line 1	1000	17	23
Line 2	1008	11	5
Line 3	1016	29	38
Tag		Data	

Memory	
1000	17
1004	23
1008	11
1012	5
1016	29
1020	38
1024	44
1028	99
1032	97
1036	25
1040	1
1044	4

FULLY-ASSOC. VS. DIRECT-MAPPED



Fully-associative N-line cache:

- N tag comparators, registers used for tag/data storage (\$\$\$)
- Location A can be stored in ANY of the N cache lines; no "collisions"
- Needs a replacement strategy to pick which line to use when loading new word(s) into cache



COLLISIONs occur when there are multiple items that we'd like to keep cached, we have room, but our management policies only keeps a subset of them.

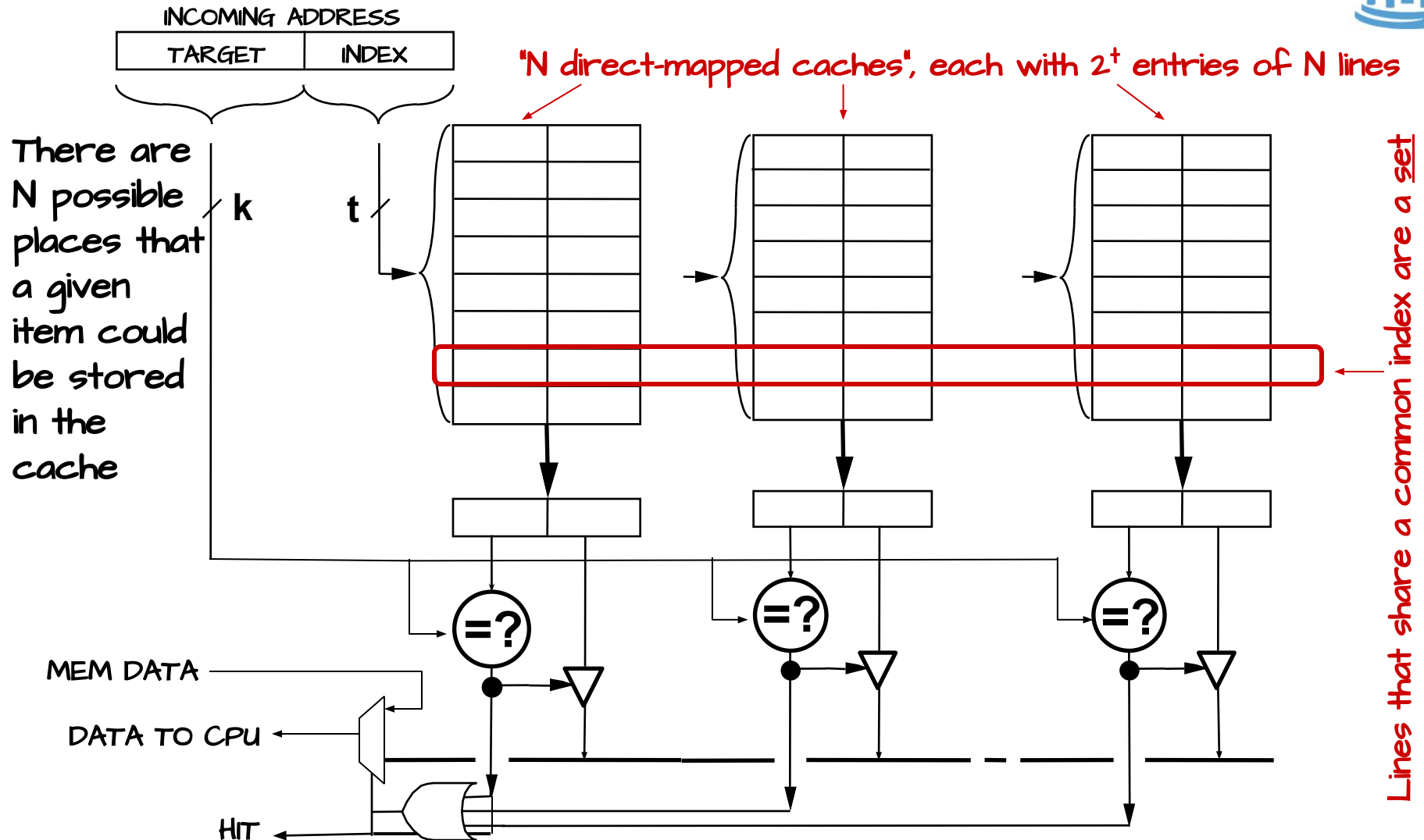
Direct-mapped N-line cache:

- One tag comparator, SRAM used for tag/data storage (\$)
- Location A is stored in a SPECIFIC line of the cache determined by its address; address "collisions" possible
- Replacement strategy not needed: each word can only be cached in one specific cache line

Is there something in-between?

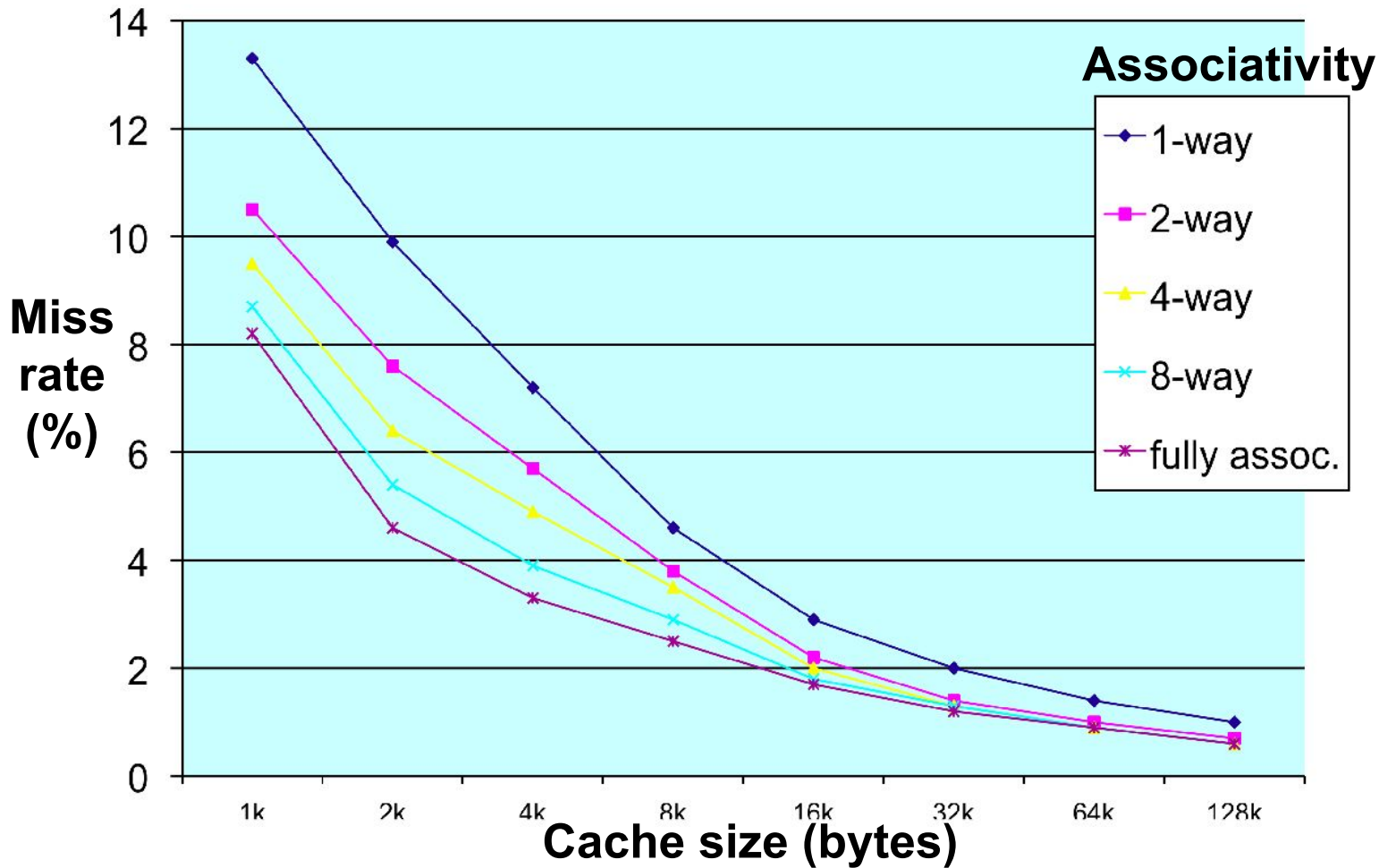


N-WAY SET-ASSOCIATIVE CACHE





ASSOCIATIVITY VS. MISS RATE



8-way is (almost) as effective as fully-associative

HANDLING WRITES



Observation: Most (80+%) of memory accesses are READs, but writes are essential. How should we handle writes?

Policies:

- WRITE-THROUGH: CPU writes are cached, but also written to main memory (stalling the CPU until write is completed). Memory always holds "the truth".
- WRITE-BACK: CPU writes are cached, but not immediately written to main memory. Memory contents can become "stale".

Additional Enhancements:

- WRITE-BUFFERS: For either write-through or write-back, writes to main memory are buffered. CPU keeps executing while writes are completed (in order) in the background.

What combination has the highest performance?



WRITE-THROUGH

ON REFERENCE TO $\text{Mem}[X]$: Look for X among tags...

HIT: $X == \text{TAG}(i)$, for some cache line i

READ: return $\text{DATA}[i]$

WRITE: change $\text{DATA}[i]$; **Start Write to $\text{Mem}[X]$**

MISS: X not found in TAG of any cache line

REPLACEMENT SELECTION:

Select some line k to hold $\text{Mem}[X]$

READ: Read $\text{Mem}[X]$

Set $\text{TAG}[k] = X$, $\text{DATA}[k] = \text{Mem}[X]$

WRITE: **Start Write to $\text{Mem}[X]$**

Set $\text{TAG}[k] = X$, $\text{DATA}[k] = \text{new Mem}[X]$



WRITE-BACK

ON REFERENCE TO $\text{Mem}[X]$: Look for X among tags...

HIT: $X = \text{TAG}(i)$, for some cache line i

READ: return $\text{DATA}(i)$

WRITE: change $\text{DATA}(i)$; start Write to $\text{Mem}[X]$

MISS: X not found in TAG of any cache line

REPLACEMENT SELECTION:

Select some line k to hold $\text{Mem}[X]$

Write Back: Write $\text{Data}(k)$ to $\text{Mem}[\text{Tag}[k]]$

READ: Read $\text{Mem}[X]$

Set $\text{TAG}[k] = X$, $\text{DATA}[k] = \text{Mem}[X]$

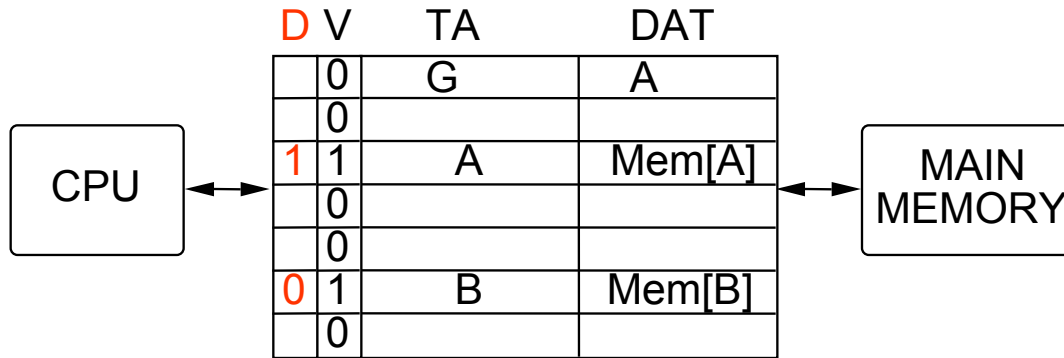
WRITE: Start Write to $\text{Mem}[X]$

Set $\text{TAG}[k] = X$, $\text{DATA}[k] = \text{new Mem}[X]$

WRITE-BACK W/ "DIRTY" BITS



Dirty and Valid bits are per line not per set



What if the cache has a block-size larger than one?

A) If only one word in the line is modified, we end up writing back ALL words

ON REFERENCE TO Mem[X]: Look for X among tags...

HIT: $X = TAG(i)$, for some cache line i

READ: return DATA(i)

WRITE: change DATA(i); Start Write to Mem[X] $D[i]=1$

MISS: X not found in TAG of any cache line

REPLACEMENT SELECTION:

Select some line k to hold Mem[X]

if $D[k] == 1$ the Write Data(k) to Mem[Tag[k]]

READ: Read Mem[X]; Set TAG[k] = X, DATA[k] = Mem[X], $D[k]=0$

WRITE: Start Write to Mem[X] $D[k]=1$

Set TAG[k] = X, DATA[k] = new Mem[X], Read Mem[X]



B) On a MISS, we need to READ the line BEFORE we WRITE it.

CACHE DESIGN SUMMARY



Various design decisions that affect cache performance

- Block size, exploits spatial locality, saves tag H/W, but, if blocks are too large you can load unneeded items at the expense of needed ones
- Write policies
- Write-through - Keeps memory and cache consistent, but high memory traffic
- Write-back - allows memory to become STALE, but reduces memory traffic

No simple answers, in the real-world cache designs are based on simulations using memory traces.