

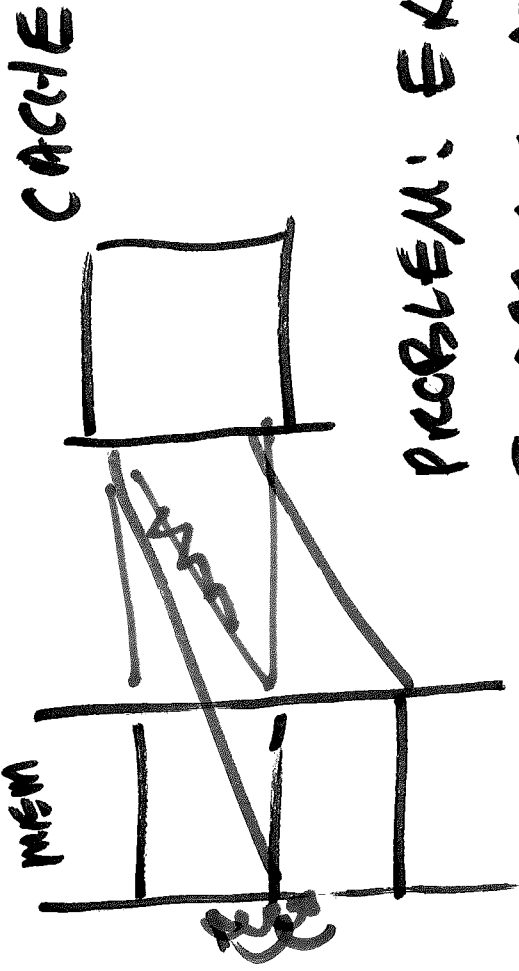
CS 451 / 551 / ECE 541

ADVANCED
COMPUTER ARCHITECTURE

SESSION no. 12

University of Idaho

DIRECT-MAPPED



PROBLEM: EXCESSIVE
SWAPPING AT BLOCK
BOUNDARIES.

DISADVANTAGES (CHALLENGES)

1. MAPPING ADDRESSES - COMPLICATED

UNIQUE TAG FOR EACH MEMORY

LOCATION

SEARCH: ADDRESS = KEY

PARALLEL

e.g. 512 CACHE ENTRIES \Rightarrow

512 COMPARATORS + LOGIC

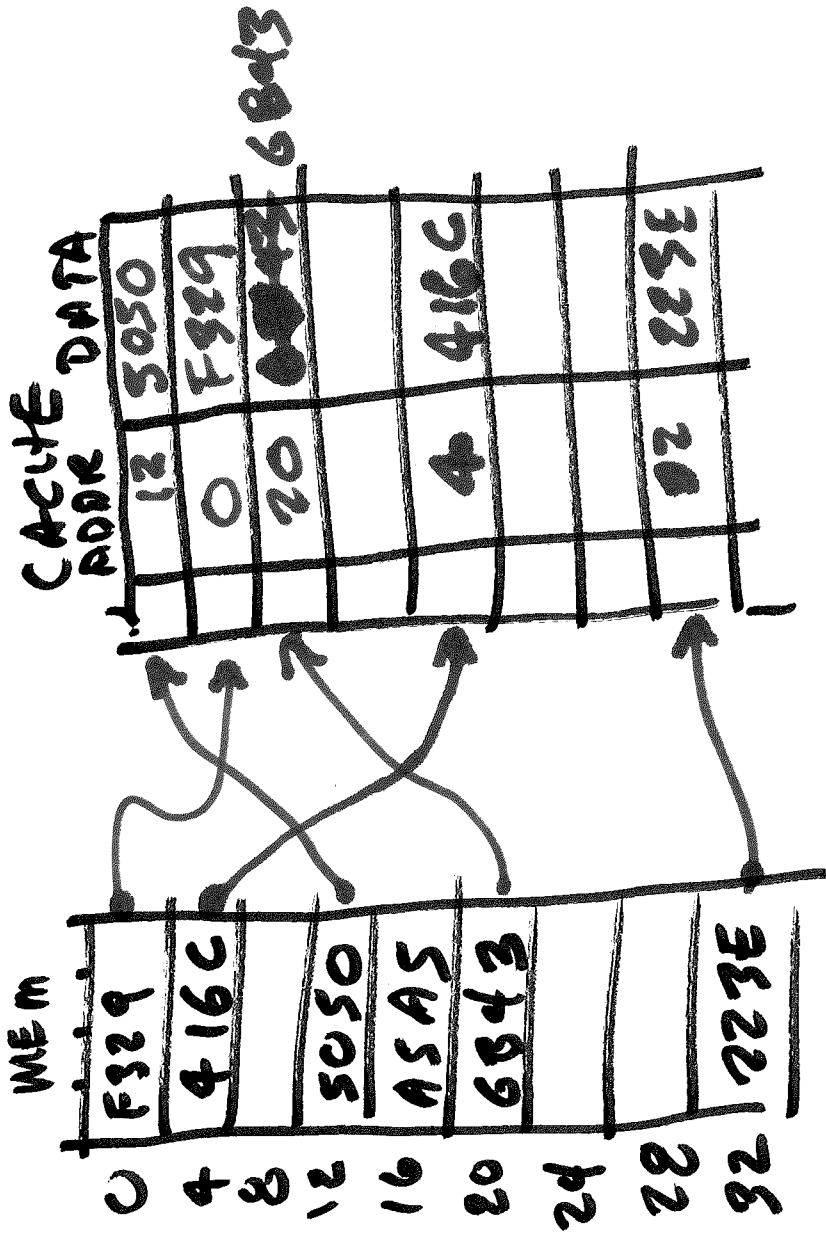
CONTENT-ADDRESSABLE MEMORY (CAM)

2. LOSS ADVANTAGES OF BLOCK TRANSFER

EACH CACHE TRANSACTION (MEME)

IS ONE WORD.

~~Direct-Mapped~~
FULLY-ASSOCIATIVE CACHE



ADVANTAGE: AVOIDS BOUNDARY PROBLEM

16 bit mem

3. IN REPLACING A CACHE ITEM -
 WHICH ITEM TO REPLACE?
 "REPLACEMENT STRATEGY"

A. LRU - LEAST RECENTLY USED.
 . EXPLITS TEMPORAL LOCALITY
 . LOSES SPATIAL LOCALITY

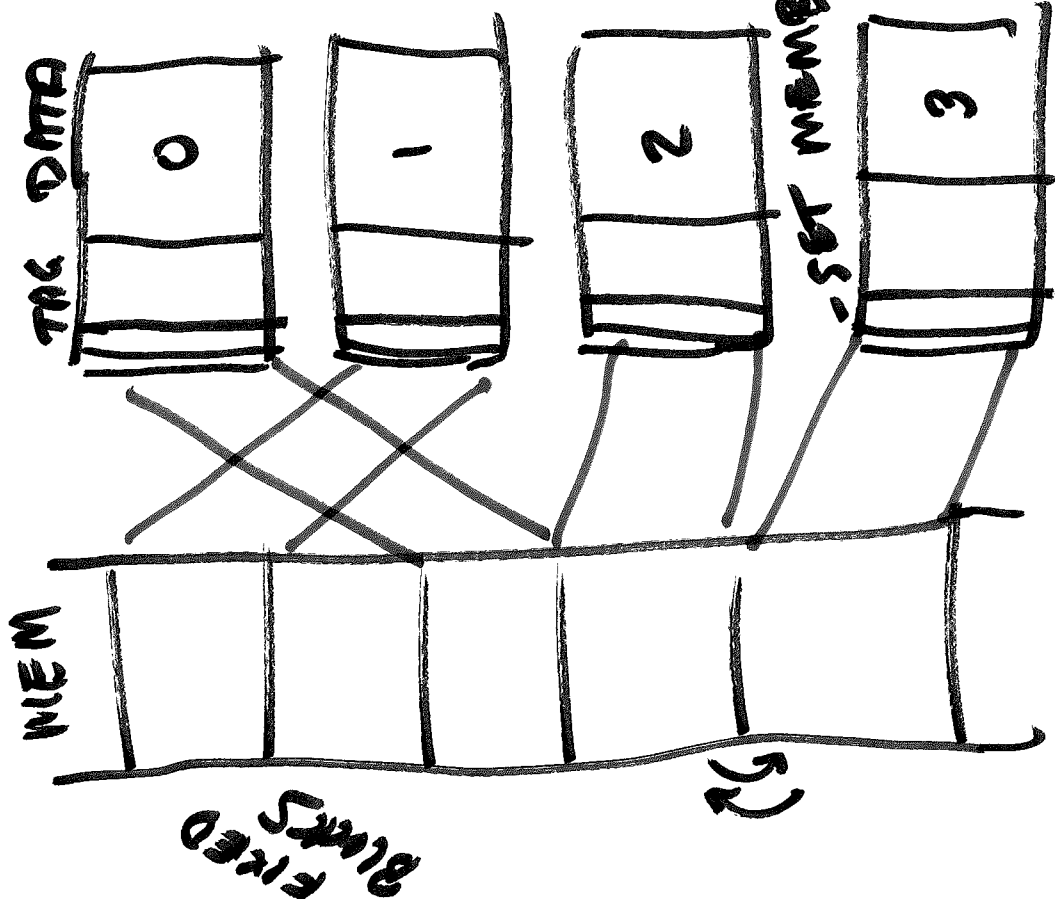
Q: How to keep track of order
 in which cache entries were
 used?
 LINEAR FEEDBACK
 PRIORITY SEARCH. SHIFT REGISTER

. HARDWARE COST ↑ (LFSR)

B. RANDOM - LOSE TEMPORAL LOCALITY
 BUT ~~SHARPE~~ SIMPLE

SET - ASSOCIATIVE

~~CACHE~~ CACHE



4-WAY SET ASSOCIATIVE

RESOLVE BOUNDARY PROBLEM

HARDWARE MUCH SIMPLER

4 COMPARATORS

REPLACEMENT STRATEGY

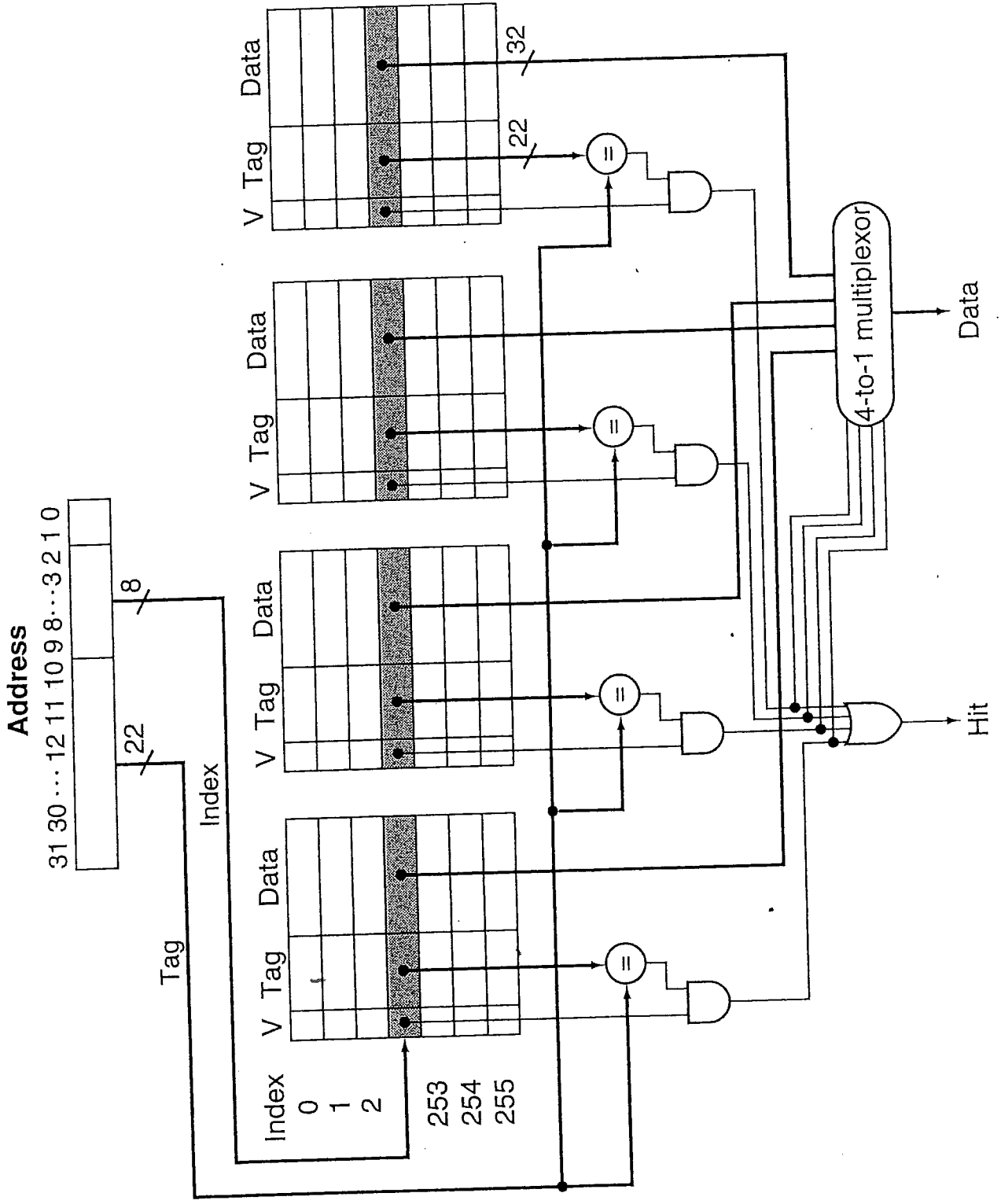
LRU

BLOCK TRANSFER

SPATIAL LOCALITY

6

B-24



CACHE PERFORMANCE

PURPOSE OF CACHE

GOAL: SPEED AVG MEMORY ACCESS

CONSTRAINT: AFFORDABLE COST

LOCALITY - SPATIAL & TEMPORAL


MAIN MEM: DRAM

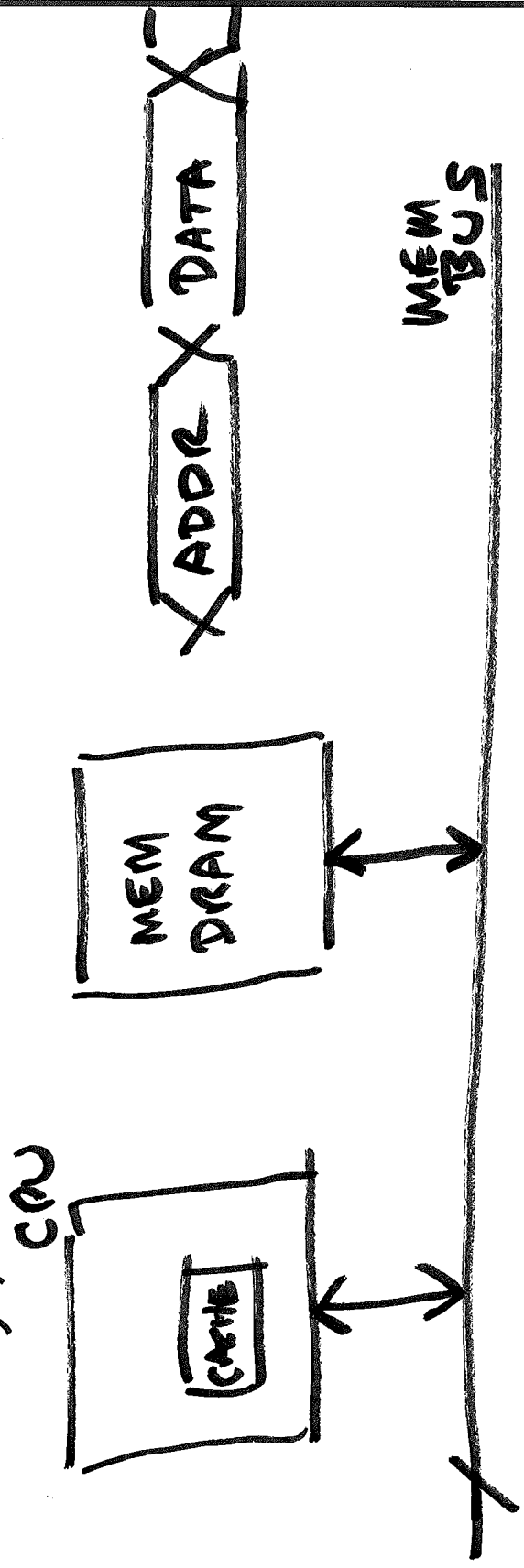
Cheap, large capacity, slower

CACHE: SRAM

Expensive, small cap, faster

SPATIAL
TEMPORAL





MISS RATE = 1 - HIT RATE

. FRACTION OF TIMES DATA NOT FOUND IN CACHE

MISS PENALTY = PERFORMANCE COST IN CACHE MISS

UNITS: MEMORY (SYSTEM) BUS CYCLES SLOWER THAN CPU CYCLES

ASSUME: 4 WORD CACHE BLOCK (LINE)
 WORD ACCESS, 4 BYTES (WORD)
 MISS PENALTY

$$MP = 1 + 4 \times 15 + 4 \times 1 = 65 \text{ BUS CYCLES}$$

ADDR TRANSFER | 15 | words | ADDR TRANSFER
 LATECY / WORD | DATA TRANSFER

BYTES TRANSFERRED

$$\text{BANDWIDTH MISS: } \frac{4 \times 4}{65} = .25 \text{ BYTES / CYCLE}$$

How to improve this?

1. MAKE MEMORY SYSTEM WIDER

eg. Bus is 2 words wide

MISS PENALTY

~~MP = 1 + 4x~~

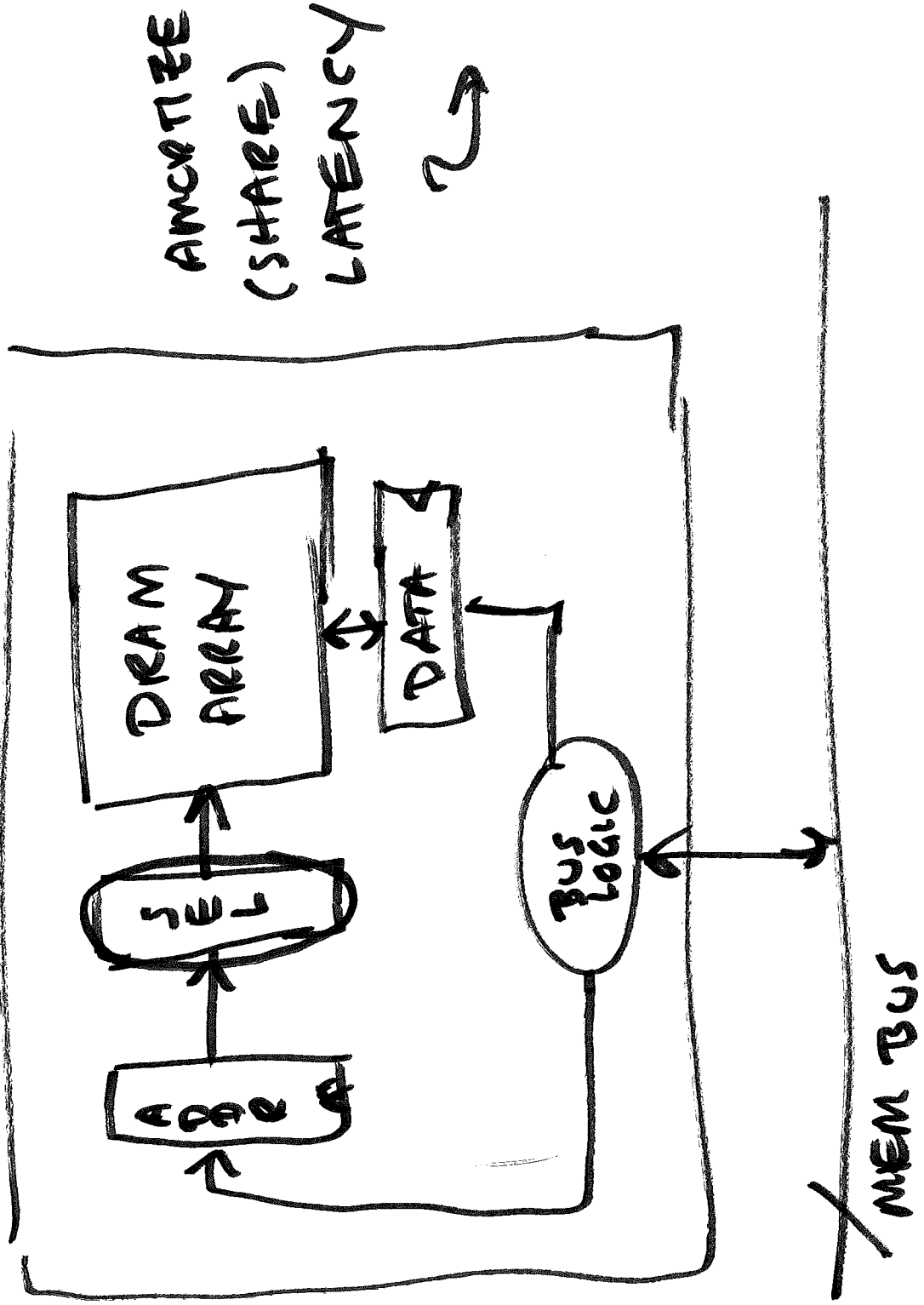
$$MP = 1 + 2 \times 15 + 2 \times 1 = 33 \text{ BUS CYCLES}$$

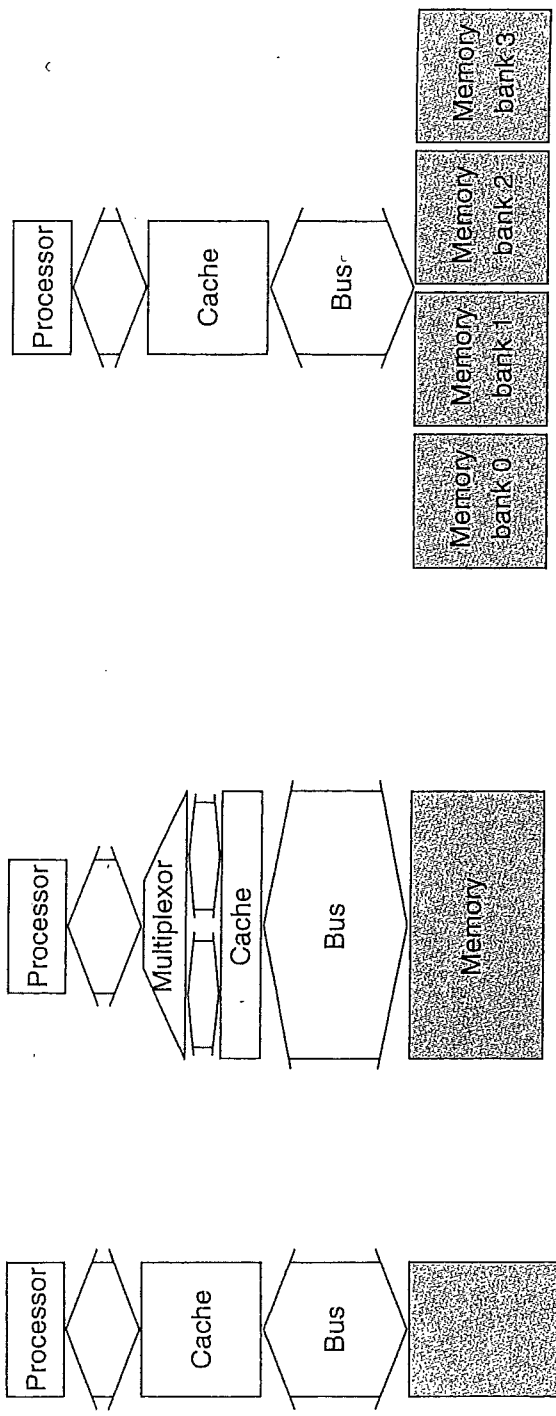


BANDWIDTH: $\frac{4 \times 4 \text{ BYTES}}{33 \text{ CYCLES}} = .48 \text{ BYTES / CYCLE}$

COST: AREA (INTERCONNECT)

2. MEMORY BANKS - INTERLEAVING





c. Interleaved memory organization

b. Wider memory organization

a. One-word-wide memory organization

FIGURE 5.11 The primary method of achieving higher memory bandwidth is to increase the physical or logical width of the memory system. In this figure, memory bandwidth is improved two ways. The simplest design, (a), uses a memory where all components are one word wide; (b) shows a wider memory, bus, and cache; while (c) shows a narrow bus and cache with an interleaved memory. In (b), the logic between the cache and processor consists of a multiplexor used on reads and control logic to update the appropriate words of the cache on writes.

MEMORY BANKS

- EACH MEM MODULE HELDS A DIFFERENT WORD IN CACHE LINE
- SAME ADDRESS (+ OFFSET)
- ALL MODULES ACCESSING INTERNALLY SIMULTANEOUSLY

$$MP = 1 + 1 \times 15 + 4 \times 1 = 20 \text{ BUS CYCLES}$$

University of Idaho

MEASURING CACHE PERFORMANCE

CPU TIME:

$(\text{CPU CYCLES} + \text{MEM STALL CYCLES}) \times$

CLOCK CYCLE TIME

• NO CACHE MISS \Rightarrow NO STALL CYCLES

IF \rightarrow DEC \rightarrow EX/ADDR \rightarrow MEM \rightarrow WB

ASSUME READ & WRITE HAVE SAME COST

MEMORY STALL CYCLES =

MEMORY ACCESES \times MISS RATE \approx MISS PENALTY
PROGRAM

MEM BUS CYCLE = $N \times$ CPU CYCLES

EX ASSUME INSTRUCTION CACHE & DATA CACHE

INSTR CACHE MISS RATE: 2%
DATA CACHE " : 4%

CPI = 2 CLOCKS / INSTR.

MISS PENALTY: 100 CPU CYCLES

I = INSTR. COUNT PER PROGRAM

INSTR MISS CYCLES: $I \times 2\% \times 100 = 2 \times I$

DATA MISS CYCLES: $I \times 4\% \times 100 = 4 \times I$
= 3.44 I

LOADS & STORES ARE 36% INSTRUCTIONS

> 3 CYCLES OF STACK / INSTR.

$$\text{ACTUAL CPI} = 2 + 3.44 = 5.44$$

1
IDEM PENALTY
(CPI)

PERFECT CACHE? NO MISSES

$$\frac{5.44}{2} = 2.72 \times \text{BETTER PERFORMANCE}$$

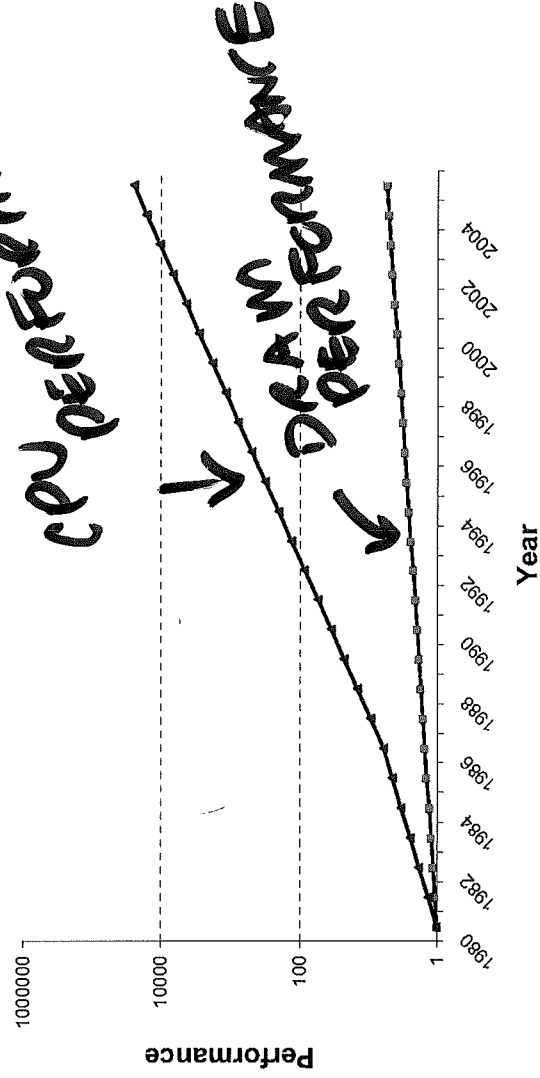
AND THAN
REAL CACHE W/ STACKS.

DRAM Access Latency

- Access times are a speed of light issue
- Memory technology is also changing
 - SRAM are getting harder to scale
 - DRAM is no longer cheapest cost/bit
- Power efficiency is an issue here as well

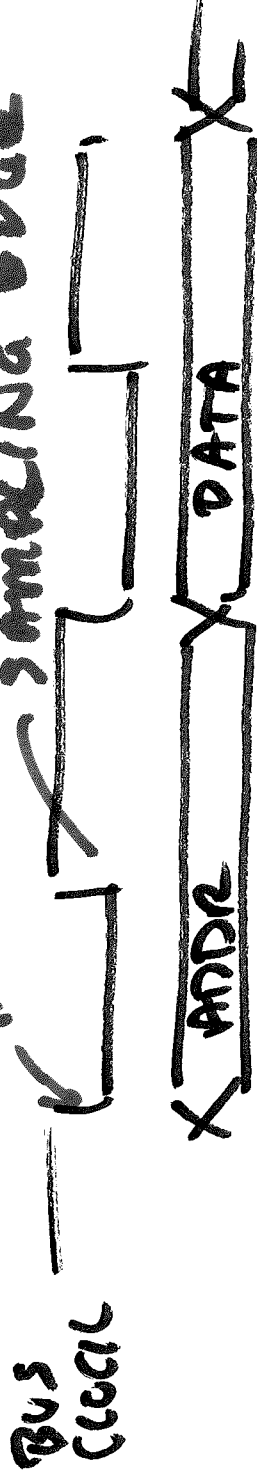
Images removed due to copyright restrictions.

μProc 60%/yr. (2X/1.5yr)
 DRAM 9%/yr. (2X/10 yrs)



DDR: DOUBLE RATE DRAM

ASSERTION EDGE
SAMPLING EDGE



DDR X ADDR X DATA X ADDR X DATA X

DDR2, DDR3 - JEDEC STANDARDS

SDRAM -

- DRAM WITH A WRAPPER TO MAKE IT BEHAVE EXTERNALLY LIKE SRAM
- DRAM DENSITY, SRAM SIMPLICITY