# EGC442
# Class Notes
# 5/9/2023

**Baback Izadi**

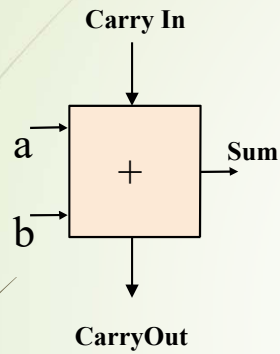Division of Engineering Programs

bai@engr.newpaltz.edu

# Final:

## Comprehensive

- Performance problems
- ALU design
- Data Path and control
- Pipelining design and hazard ✓
- Cache memory
- Virtual memory
- Parallel Computing ✓

1. review notes
   Zibodk
2. redo Quiz
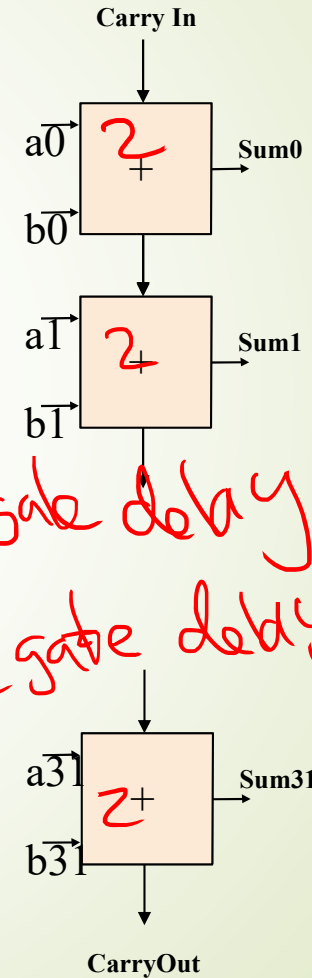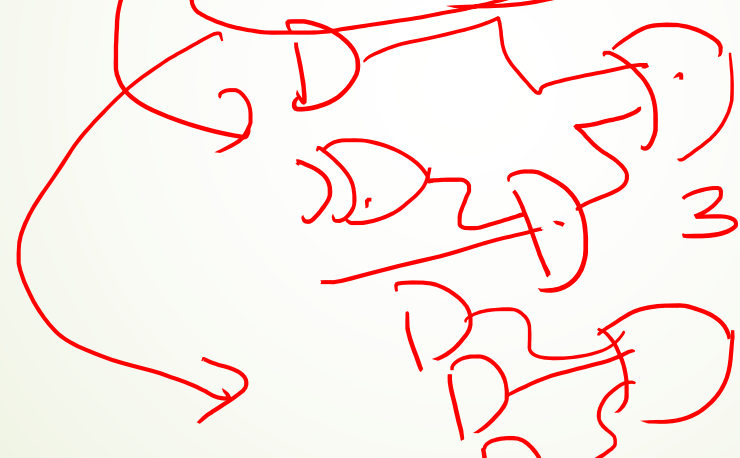3. redo tests
4. Aw's

# Making a faster adder Full Adder

Let's look at a 1-bit ALU for addition:

**Carry In**

a → [ + ] → **Sum**

b →

↓

**CarryOut**

$$Sum = a \oplus b \oplus c_{in}$$
$$C_{out} = a\,b + (a \oplus b)\,c_{in}$$

$$C_{out} = a\,b + a\,c_{in} + b c_{in}$$

What is the propagation delay of a 32-bit adder?

**Carry In**

a0 → [ **2** + ] → **Sum0**
b0 →

↓

a1 → [ **2** + ] → **Sum1**
b1 →

↓

a31 → [ **2** + ] → **Sum31**
b31 →

↓

**CarryOut**

$32 * 2 = 64$

3 gate delay

2 gate delay

# Problem with Ripple Carry

- Is a 32-bit ALU as fast as a 1-bit ALU?
- Is there more than one way to do addition?
  - two extremes: ripple carry and sum-of-products
- Can you see the ripple? How could you get rid of it?

$$c_1 = b_0c_0 + a_0c_0 + a_0b_0$$
$$c_2 = b_1c_1 + a_1c_1 + a_1b_1 \qquad c_2 =$$
$$c_3 = b_2c_2 + a_2c_2 + a_2b_2 \qquad c_3 =$$
$$c_4 = b_3c_3 + a_3c_3 + a_3b_3 \qquad c_4 =$$

Not feasible! Why?

# Carry-lookahead adder

- An approach in-between our two extremes
- $c_1 = b_0c_0 + a_0c_0 + a_0b_0 = (b_0 + a_0)c_0 + a_0b_0$
  - If we didn't know the value of carry-in, what could we do?
  - When would we always generate a carry?          $g_i = a_i b_i$
  - When would we propagate the carry?          $p_i = a_i + b_i$
- Did we get rid of the ripple?

$c_1 = g_0 + p_0c_0$

$c_2 = g_1 + p_1c_1$          $c_2 = g_1 + p_1 g_0 + p_1 p_0 c_0$
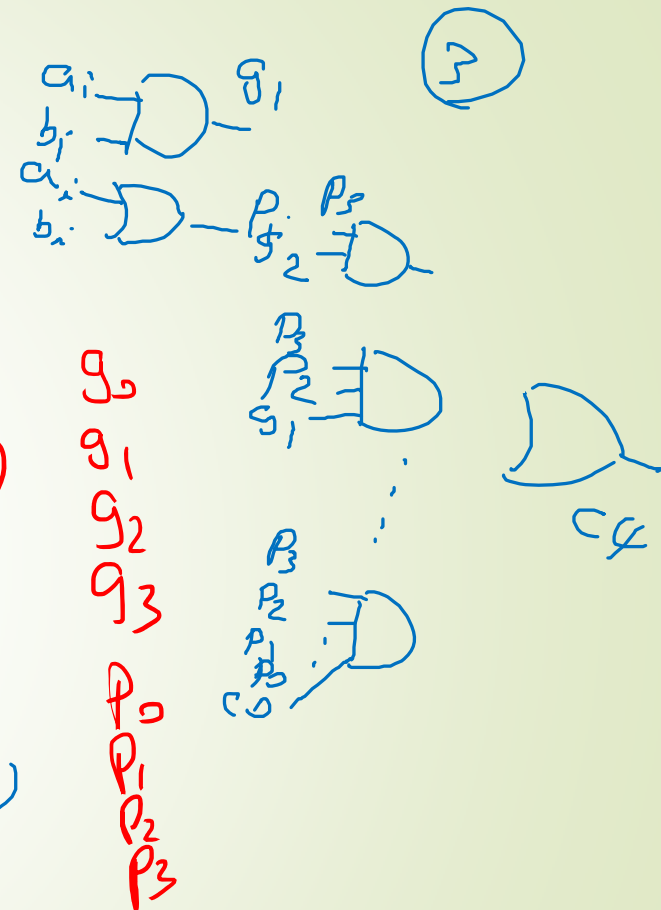
$c_3 = g_2 + p_2c_2$          $c_3 = g_2 + c_2 ( \quad )$

$c_4 = g_3 + p_3c_3$          $c_4 = g_3 + p_3 ( \quad )$

$c_4 = g_3 + p_3g_2 + p_3p_2 g_1 + p_3 p_2 p_1 g_0 + p_3 p_2 p_1p_0 c_0$

# Carry Look Ahead Design trick

$$c_{in}$$

$a_0$

$b_0$    g   p     S $_0$

$$c_1 = g_0 + p_0 \, c_0$$

$a_1$          S$_1$

$b_1$    g   p

$$c_2 = g_1 + p_1 \, g_0 + p_1 p_0 c_0$$

$a_2$          S

$b_2$    g   p     2

$$c_3 = g_2 + p_2 \, g_1 + p_2 \, p_1 \, g_0 + p_2 \, p_1 p_0 \, c_0$$

$a_3$          S$_3$

$b_3$    g   p
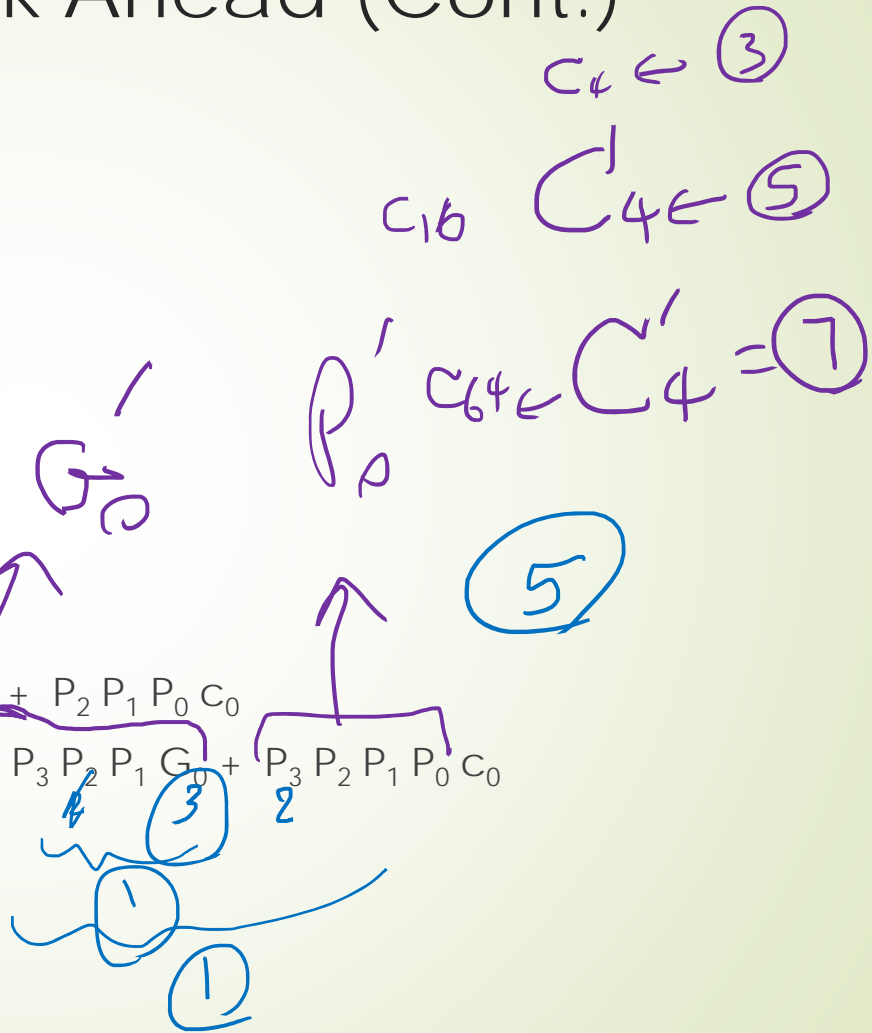
$$g = a \, b$$
$$p = a + b$$

$$C_4 = \ldots$$

# 16 Bit Carry Look Ahead

- $g_i = a_i\, b_i$

- $p_i = a_i + b_i$

- $c_1 = g_0 + p_0\, c_0$

- $c_2 = g_1 + p_1\, g_0 + p_1 p_0 c_0$

- $c_3 = g_2 + p_2\, g_1 + p_2\, p_1\, g_0 + p_2\, p_1 p_0\, c_0$

- $c_4 = g_3 + p_3 g_2 + p_3 p_2\, g_1 + p_3\, p_2\, p_1\, g_0 + p_3\, p_2\, p_1 p_0\, c_0$

$G_0$

- $G_0 = g_3 + p_3\, g_2 + p_3 p_2\, g_1 + p_3\, p_2\, p_1\, g_0$     $P_0 = p_3\, p_2\, p_1\, p_0$

- $G_1 = g_7 + p_7\, g_6 + p_7\, p_6\, g_5 + p_7\, p_6\, p_5\, g_4$     $P_1 = p_7\, p_6\, p_5\, p_4$

- $G_2 = g_{11} + p_{11}\, g_{10} + p_{11}\, p_{10}\, g_9 + p_{11}\, p_{10}\, p_9\, g_8$     $P_2 = p_{11}\, p_{10}\, p_9\, p_8$

- $G_3 = g_{15} + p_{15}\, g_{14} + p_{15}\, p_{14}\, g_{13} + p_{15}\, p_{14}\, p_{13}\, g_{12}$     $P_3 = p_{15}\, p_{14}\, p_{13}\, p_{12}$

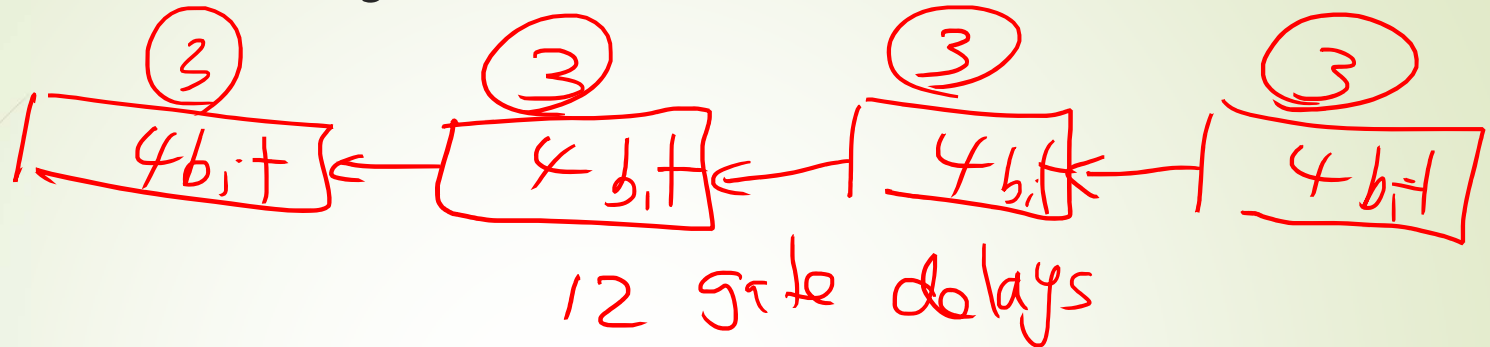# 16 Bit Carry Look Ahead (Cont.)

$C_1 = G_0 + P_0C_0$

$C_2 = G_1 + P_1C_1$     $C_2 =$

$C_3 = G_2 + P_2C_2$     $C_3 =$

$C_4 = G_3 + P_3C_3$     $C_4 =$

$C_1 = C_4 = G_0 + P_0 \, c_0$

$C_2 = C_8 = G_1 + P_1 \, G_0 + P_1 P_0 \, c_0$

$C_3 = C_{12} = G_2 + P_2 \, G_1 + P_2 P_1 G_0 + P_2 P_1 P_0 \, c_0$

$C_4 = C_{16} = G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0 + P_3 P_2 P_1 P_0 \, c_0$

③ ③ ③ ③

4 bit ← 4 bit ← 4 bit ← 4 bit

12 gate delays

32 Bit

\* 8X3 = 24 gate delay
    if we use 4 bit CLA

\* ripple    32 \* 2 = 64

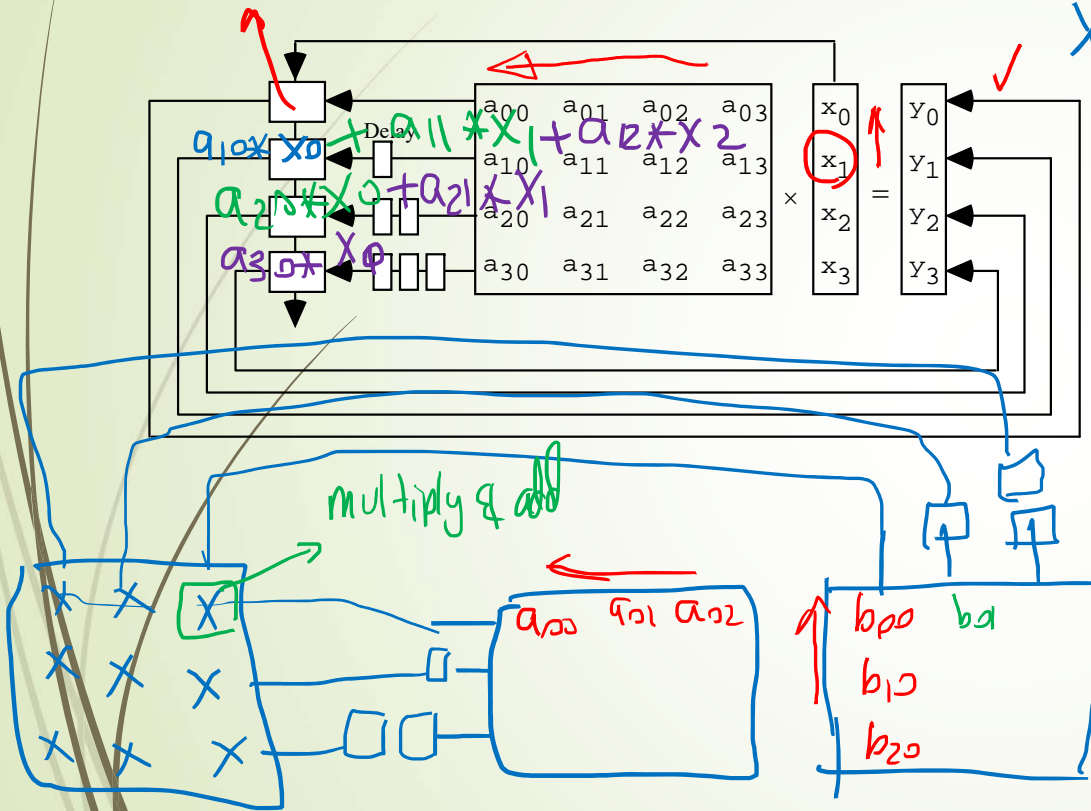\* 2 level CLA → 16 bit CLA  ⑤

10 gate delay   16 bit CLA ← 16 bit CLA

SIMD    multiply add

6) The following diagram depicts a vector processor for matrix – vector multiplier. Devise a vector processor for a matrix – matrix multiplication.

$a_{00} * X_0 + a_{01} * X_1 + a_{02} * X_2 + a_{03} * X_3$

$Y_0 = a_{00} * X_0 + a_{01} * X_1 + a_{02} * X_2 + a_{03} * X_3$

$a_{10} * X_0 + a_{11} * X_1 + a_{12} * X_2$

$a_{20} * X_0 + a_{21} * X_1$

$a_{30} * X_0$

Delay

$$\begin{bmatrix} a_{00} & a_{01} & a_{02} & a_{03} \\ a_{10} & a_{11} & a_{12} & a_{13} \\ a_{20} & a_{21} & a_{22} & a_{23} \\ a_{30} & a_{31} & a_{32} & a_{33} \end{bmatrix} \times \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

multiply & add

X

$a_{00} \quad a_{01} \quad a_{02}$

$b_{00} \quad b_{01}$
$b_{10}$
$b_{20}$

$$\begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{10} & a_{12} \\ a_{20} & a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} b_{00} & b_{01} & b_{02} \\ b_{10} & b_{11} & b_{12} \\ b_{20} & b_{21} & b_{22} \end{bmatrix}$$
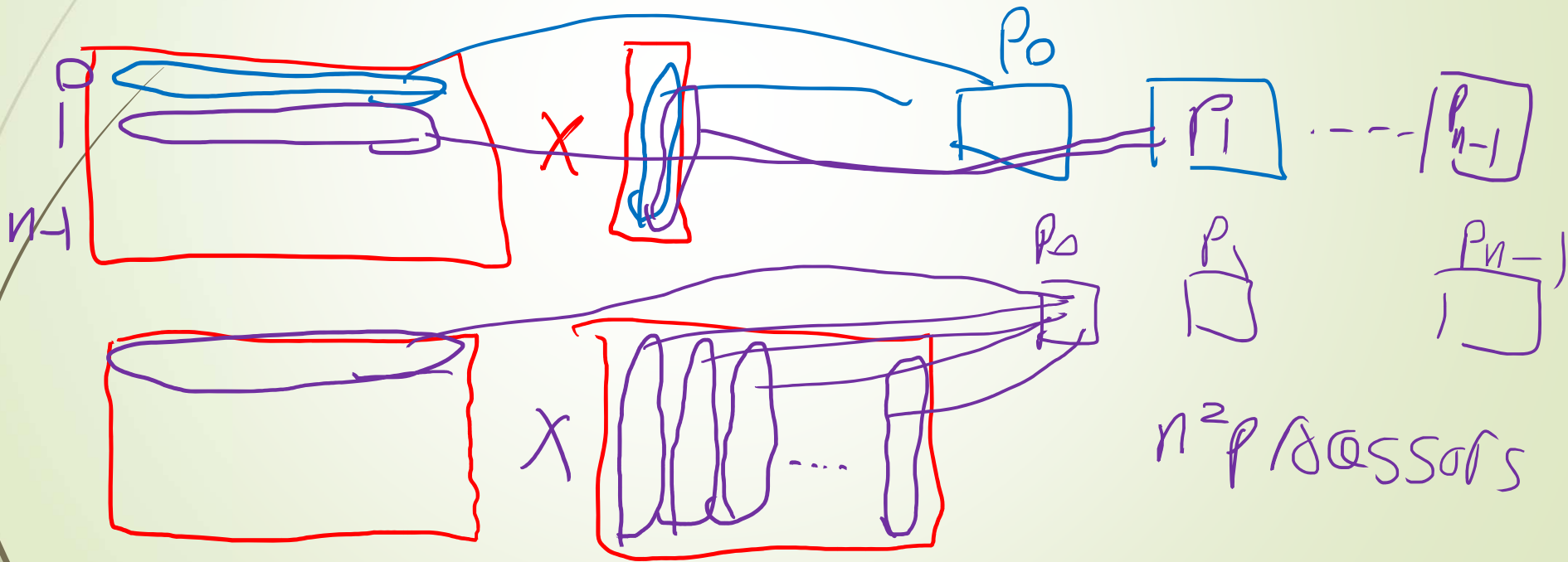
# Direct Mapped Cache

- Taking advantage of spatial locality:



Handwritten annotations:

$2^{12} = 4K$ rows

$4K \times 4 \times 4$ way Bytes

$64K$ cache

word size → 4 Bytes

direct + 4 way

# 4 way direct cache

word size 4 bytes

8 rows

128 64 32 16 8 4 2 1

total cach

2 $A_6 A_5 A_4$ $A_3 A_2$ $A_1 A_0$

$32 - 7 = 25$

tag

M

24 0 0 1

001 x x

28 001

H

index

select word 12 000 11

M

11          10          01          00

000

001

010

011

100

101

110

111

| tag 140 | 8 136 | 4 132 | 0 128 |
|---------|-------|-------|-------|
| 28      | 24    | 20    | 16    |
|         |       |       |       |
|         |       |       |       |
|         |       |       |       |
|         |       |       |       |
|         |       |       |       |
|         |       |       |       |

4 H

8 H

132 M    000 01

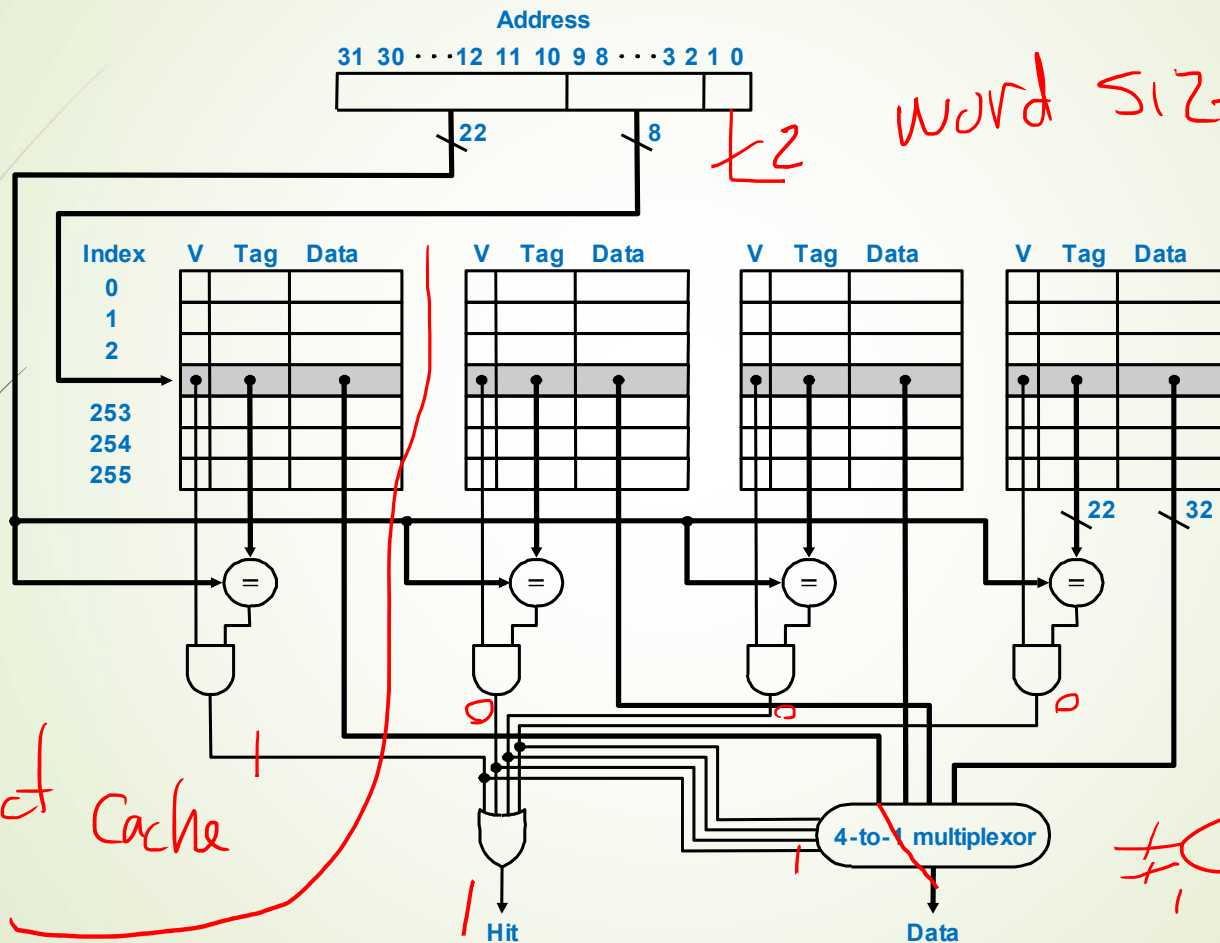# Decreasing Miss Ratio with Associativity

Associativity: Reducing cache misses by more flexible placement of blocks

One Way
Assoc

2 Way

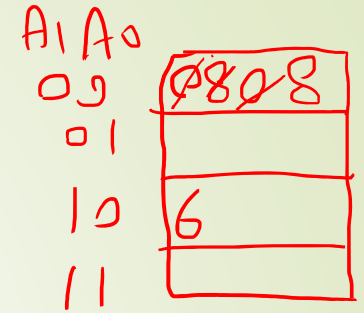| Direct | Set Associative | Fully Associative |

Data

Cache Location

# 4-Way Associative Cache Organization

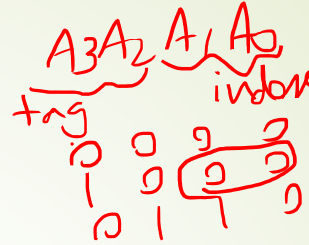# Associative Caches

- Fully associative
  - Allow a given block to go in any cache entry
  - Requires all entries to be searched at once
  - Comparator per entry (expensive)
- *n*-way set associative
  - Each set contains *n* entries
  - Block number determines which set
    - (Block number) modulo (#Sets in cache)
  - Search all entries in a given set at once
  - *n* comparators (less expensive)

# Associativity Example

*(handwritten annotations: A3 A2 A1 A0 tag index; A1 A0 with binary 00, 01, 10, 11; grid with 0808, 6)*

- Assume word size = 1 byte
- Compare 4-block caches
  - Direct mapped, 2-way set associative, fully associative
  - Block access sequence: 0, 8, 0, 6, 8
- Direct mapped

| Block address | Cache index | Hit/miss | Cache content after access | | | |
|---|---|---|---|---|---|---|
| | | | 0 | 1 | 2 | 3 |
| 0 | 0 | miss | **Mem[0]** | | | |
| 8 | 0 | miss | **Mem[8]** | | | |
| 0 | 0 | miss | **Mem[0]** | | | |
| 6 | 2 | miss | Mem[0] | | **Mem[6]** | |
| 8 | 0 | miss | **Mem[8]** | | Mem[6] | |

# Associativity Example

- 2-way set associative

| Block address | Cache index | Hit/miss | Cache content after access | | | |
|---|---|---|---|---|---|---|
| | | | Set 0 | | Set 1 | |
| 0 | 0 | miss | Mem[0] | | | |
| 8 | 0 | miss | Mem[0] | Mem[8] | | |
| 0 | 0 | hit | Mem[0] | Mem[8] | | |
| 6 | 0 | miss | Mem[0] | Mem[6] | | |
| 8 | 0 | miss | Mem[8] | Mem[6] | | |

## Fully associative

| Block address | | Hit/miss | Cache content after access | | | |
|---|---|---|---|---|---|---|
| 0 | | miss | Mem[0] | | | |
| 8 | | miss | Mem[0] | Mem[8] | | |
| 0 | | hit | Mem[0] | Mem[8] | | |
| 6 | | miss | Mem[0] | Mem[8] | Mem[6] | |
| 8 | | hit | Mem[0] | Mem[8] | Mem[6] | |

# Spectrum of Associativity

- Assume one byte word size.
- For a cache with 8 entries. Assume using the "least recently used" replacement strategy 7, 4, 17, 12, 4, 25, 7, 13

**One-way set associative (direct mapped)**

| Block | Tag | Data |
|-------|-----|------|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |

**Two-way set associative**

| Set | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

**Four-way set associative**

| Set | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|-----|------|-----|------|
| 0 | | | | | | | | |
| 1 | | | | | | | | |

**Eight-way set associative (fully associative)**

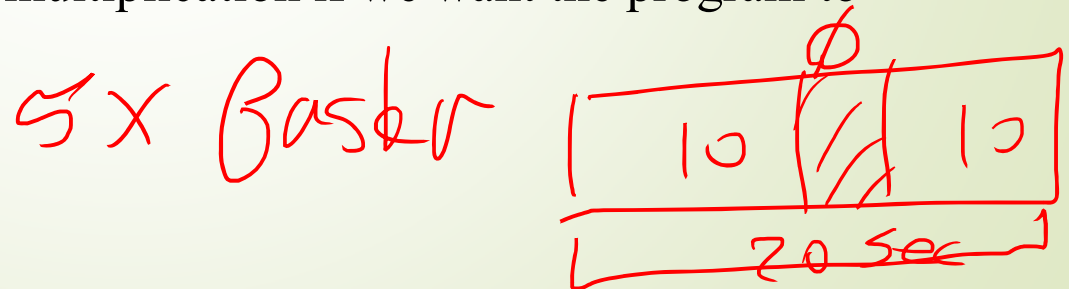| Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|
| | | | | | | | | | | | | | | | |

# Pitfall: Amdahl's Law

$$T_{improved} = \frac{T_{affected}}{\text{improvement factor}} + T_{unaffected}$$

$$\frac{80}{5} = 16\,X$$

10    80 sec    10    ⑤

100 sec    2.5 sec
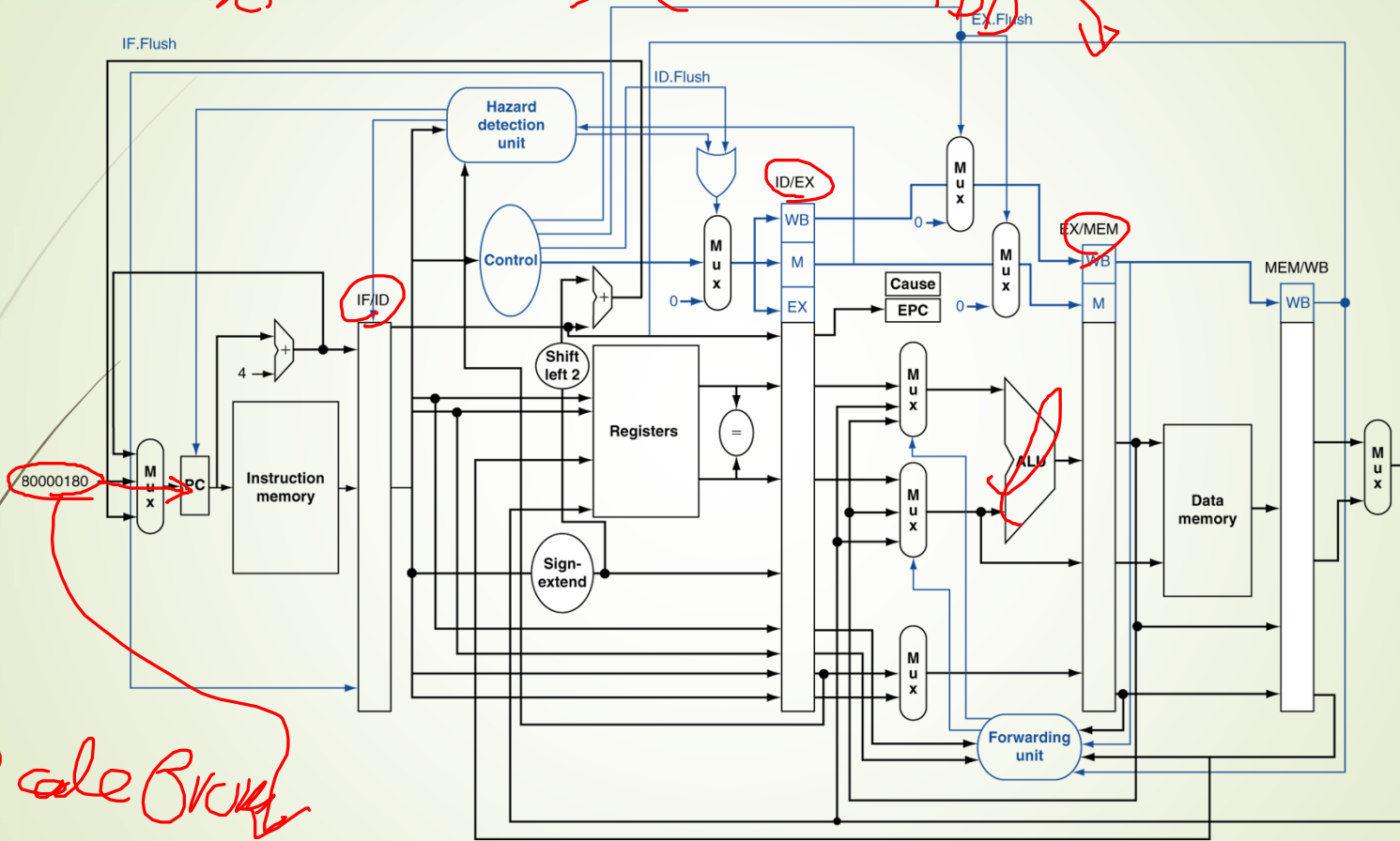
Example: Suppose a program runs in 100 seconds on a machine, with multiply responsible for 80 seconds of this time.   How much do we have to improve the speed of multiplication if we want the program to run 4 times faster?

5 X Faster

10 ∅ 10

20 sec

# Pipeline with Exceptions

# Pitfall: Amdahl's Law

$$T_{improved} = \frac{T_{affected}}{\text{improvement factor}} + T_{unaffected}$$
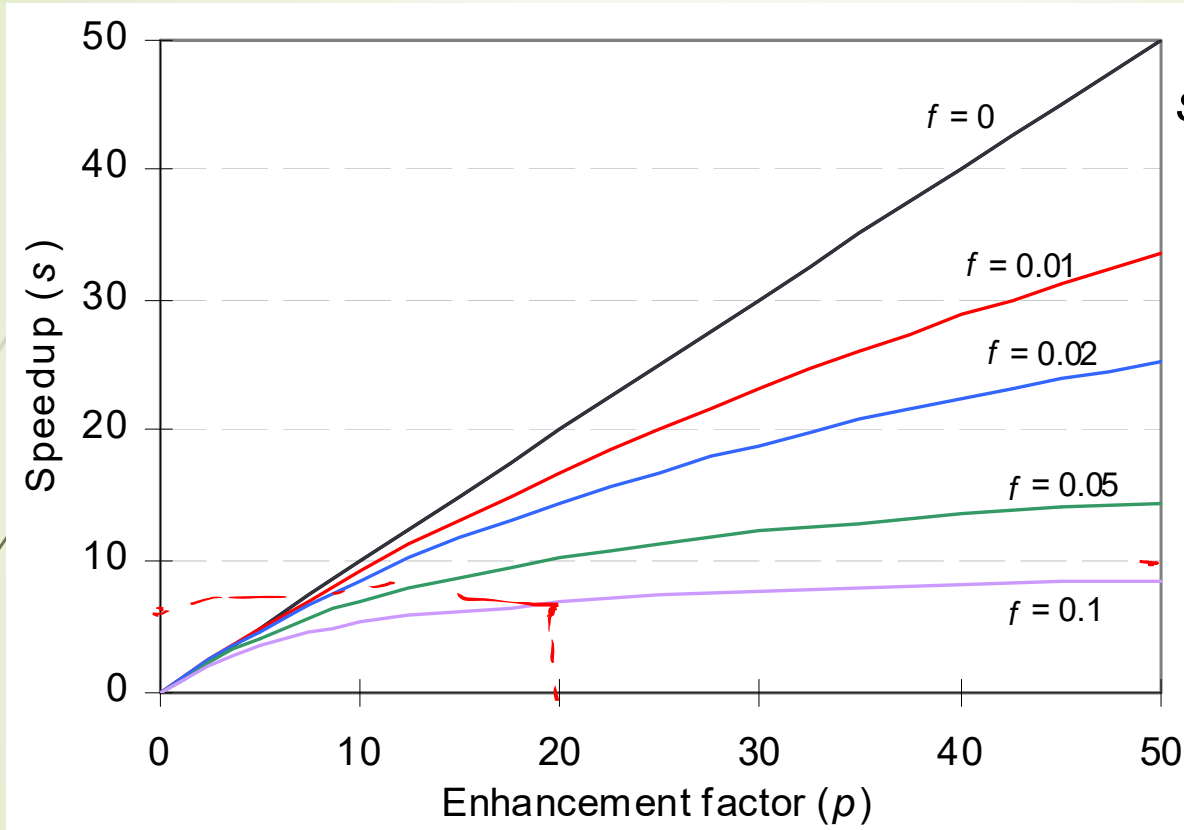
$$\frac{100}{4} = \frac{80}{n} + 20$$

- n = 16

- How much improvement in multiply performance to get 5× overall?

$$20 = \frac{80}{n} + 20$$

- Can't be done!

- Improving an aspect of a computer and expecting a proportional improvement in overall performance

- Corollary: make the common case fast

# Amdahl's Law – Parallel Processing Version



$$s = \frac{\text{Exec Time 1 Processor}}{\text{Exec Time P Processors}}$$

$f$ = fraction sequential

$p$ = speedup of the rest

$$s = \frac{1}{f + (1-f)/p}$$

$$\leq \min(p, 1/f)$$

10%

$f = .1$

$$s = \frac{1}{.1 + \frac{.9}{p}}$$

$$s = \frac{1}{.1 + \frac{.9}{20}}$$

Limit on speed-up according to Amdahl's law.