# زبان و ساختار کامپیوتر

## فصل پنجم

## محاسبات کامپیوتری

# *Copyright Notice*

*Parts (text & figures) of this lecture are adopted from :*

📖 *D. Patterson & J. Hennessey, "Computer Organization & Design, The Hardware/Software Interface", 5th Ed., MK publishing, 2014*

ساختار و زبان کامپیوتر

# *Some Concepts*

- LSB and MSB
  - Least Significant Bit (LSB)
  - Most Significant Bit (MSB)
- Signed versus Unsigned
  - Unsigned (Assume all non-negative numbers)
    - Used usually for memory addresses
  - Signed
    - Using sign bit
    - Using two's complement notation
- Carry Out
- Overflow

# Number Representation

○ *Weighted number system*

  ● *Can be represented in any base (radix)*

    ○ *Value of $i^{th}$ digit "$d_i$" = $d_i * Base^i$*

    ○ *$0 =< d_i < Base$*

| 31 | 30 | 29 | … | i | … | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|----|----|----|----|
| $d_{31}$ | $d_{30}$ | $d_{29}$ | … | $d_i$ | … | $d_3$ | $d_2$ | $d_1$ | $d_0$ |

# *Base Examples*

## *Binary*

$$(101101)_2 = 1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$$

## *Octal*

$$(736.4)_8 = 7 \times 8^2 + 3 \times 8^1 + 6 \times 8^0 + 4 \times 8^{-1}$$
$$= 7 \times 64 + 3 \times 8 + 6 \times 1 + 4/8 = (478.5)_{10}$$

## *Hexadecimal*

$$(F3)_{16} = F \times 16 + 3 = 15 \times 16 + 3 = (243)_{10}$$

## *Decimal*

$$(7245)_{10} = 7 \times 10^3 + 2 \times 10^2 + 4 \times 10^1 + 5 \times 10^0$$

ساختار و زبان کامپیوتر

# *Which Base?*

○ *Number Representation*

- *Humans prefer base 10, why?*

- *Base 2 best works for computers, why?*

- *Base 10 inefficient for computers, why?*

# Decimal to Base i Conversion

○ Convert $65_{10}$ to base 5

○ Convert $19_{10}$ to base 2

# *Converting Fractions*

○ *Convert $0.4304_{10}$ to base 5 = $0.2034$*

$$.4304$$
$$\times \quad 5$$
$$\overline{2.1520}$$

$$.1520$$
$$\times \quad 5$$
$$\overline{0.7600}$$

$$0.7600$$

$$.7600$$
$$\times \quad 5$$
$$\overline{3.8000}$$

$$.8000$$
$$\times \quad 5$$
$$\overline{4.0000}$$

○ *Convert $0.34375_{10}$ to base 2 = $0.01011$*

# *Converting between power of 2 radices*

○ Convert $1101010100001111_2$

- to base 8

- to base 16

# *Binary to Decimal Conversion*

○ *Question:*

- *What is decimal value of this 32-bit number?*

  *1111 1111 1111 1111 1111 1111 1111 1000$_{two}$*

- *Depends on the notation*

  ○ *Signed*

  ○ *Unsigned*

# *Unsigned Numbers*

$$N = (d_{31}*2^{31}) + (d_{30}*2^{30}) + \ldots + (d_{1}*2^{1}) + (d_{0}*2^{0})$$

```
0000 0000 0000 0000 0000 0000 0000 0000_two = 0_ten
0000 0000 0000 0000 0000 0000 0000 0001_two = 1_ten
0000 0000 0000 0000 0000 0000 0000 0010_two = 2_ten
...                              ...
1111 1111 1111 1111 1111 1111 1111 1101_two = 4,294,967,293_ten
1111 1111 1111 1111 1111 1111 1111 1110_two = 4,294,967,294_ten
1111 1111 1111 1111 1111 1111 1111 1111_two = 4,294,967,295_ten
```

# *Signed Numbers (2's Complement)*

$$N = (d_{31} * -2^{31}) + (d_{30} * 2^{30}) + \ldots + (d_1 * 2^1) + (d_0 * 2^0)$$

```
0000 0000 0000 0000 0000 0000 0000 0000_two =  0_ten
0000 0000 0000 0000 0000 0000 0000 0001_two =  1_ten
0000 0000 0000 0000 0000 0000 0000 0010_two =  2_ten
...                                   ...
0111 1111 1111 1111 1111 1111 1111 1101_two =  2,147,483,645_ten
0111 1111 1111 1111 1111 1111 1111 1110_two =  2,147,483,646_ten
0111 1111 1111 1111 1111 1111 1111 1111_two =  2,147,483,647_ten
1000 0000 0000 0000 0000 0000 0000 0000_two = -2,147,483,648_ten
1000 0000 0000 0000 0000 0000 0000 0001_two = -2,147,483,647_ten
1000 0000 0000 0000 0000 0000 0000 0010_two = -2,147,483,646_ten
...                                   ...
1111 1111 1111 1111 1111 1111 1111 1101_two = -3_ten
1111 1111 1111 1111 1111 1111 1111 1110_two = -2_ten
1111 1111 1111 1111 1111 1111 1111 1111_two = -1_ten
```

# Other Signed Number Notations

○ Signed-Magnitude Notation

○ Ones' Complement Notation

○ Biased Notation

# Signed-Magnitude Notation

○ Signed Notation with Sign Flag

○ Most **positive** number

- 011 ... 1

○ Most **negative** number

- 111 ... 1

○ There are two zero's

- 000 ... 0

- 100 ... 0

○ Used in floating point representation (mantissa)

# Ones' Complement Notation

○ Positive number *same* as two's complement

○ Negative number:

  ● *Invert* each bit in positive representation

○ There are two zero's in ones' complement

  ● *000...0*

  ● *111...1*

○ Most *positive* number

  ● *0111 ... 1*

○ Most *negative* number

  ● *1000 ... 0*

# *Biased Notation (Excess $2^{n-1}$)*

- *If n bits used for representation:*
  - *Add all numbers with $2^{n-1}$*
- *Zero represented by*
  - *100 ... 0*
- *Most negative number ($-2^{n-1}$)*
  - *000 ... 0*
- *Most positive number ($2^{n-1}-1$)*
  - *111 ... 1*

| N | Excess-4 | 2's Comp |
|---|----------|----------|
| -4 | 000 | 100 |
| -3 | 001 | 101 |
| -2 | 010 | 110 |
| -1 | 011 | 111 |
| 0 | 100 | 000 |
| 1 | 101 | 001 |
| 2 | 110 | 010 |
| 3 | 111 | 011 |

ساختار و زبان کامپیوتر

# Biased Notation (Excess $2^{n-1}-1$)

○ **If n bits used for representation:**

   ● Add all numbers with $2^{n-1}-1$

○ **Zero represented by**

   ● 011 … 1

○ **Most negative number ($-2^{n-1}+1$)**

   ● 000 … 0

○ **Most positive number ($2^{n-1}$)**

   ● 111 … 1

○ *Used in floating point representation (exponent)*

| N | Excess-3 | Excess-4 |
|---|----------|----------|
| -4 | - | 000 |
| -3 | 000 | 001 |
| -2 | 001 | 010 |
| -1 | 010 | 011 |
| 0 | 011 | 100 |
| 1 | 100 | 101 |
| 2 | 101 | 110 |
| 3 | 110 | 111 |
| 4 | 111 | - |

ساختار و زبان
کامپیوتر

# Signed Number Notations (Summary)

○ *Unbiased*

- *Positive Binary*          $N=+14$          *0 0001110*

- *Signed-Magnitude*       $-N=-14$         *1 0001110*

- *1's Complement* $(2^n-N-1)$  $-N=-14$   *1 1110001*

- *2's Complement* $(2^n-N)$   $-N=-14$    *1 1110010*

○ *Biased* $(2^{n-1}+N)$

- *Positive Binary*          $N=+14$          *1 0001110*

- *Negative Binary*        $M=-14$         *0 1110010*

ساختار و زبان کامپیوتر

# *Integer Addition / Subtraction*

```
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0  +        64
 0 0 0 0 0 0 0 0 0 0 1 0 1 0 1 0         + 42
 0 0 0 0 0 0 0 0 0 1 1 0 1 0 1 0           106
```

```
 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0  +        64
 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 0         - 42
 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 1 0          22
```

*2's compliment*

# Overflow Conditions for Add/Sub

ساختار و زبان کامپیوتر

| Operation | Operand A | Operand B | Result indicating overflow |
|-----------|-----------|-----------|----------------------------|
| $A + B$ | $\geq 0$ | $\geq 0$ | $< 0$ |
| $A + B$ | $< 0$ | $< 0$ | $\geq 0$ |
| $A - B$ | $\geq 0$ | $< 0$ | $< 0$ |
| $A - B$ | $< 0$ | $\geq 0$ | $\geq 0$ |

○ While adding signed numbers, an overflow occurs when

- Both operands have the same sign,

- but the result has the opposite sign

○ the carry into and out of the MSB differ

# One-bit Full Adder

Input and output specification for a 1-bit adder

| Inputs | | | Outputs | | |
|---|---|---|---|---|---|
| **a** | **b** | **CarryIn** | **CarryOut** | **Sum** | **Comments** |
| 0 | 0 | 0 | 0 | 0 | $0 + 0 + 0 = 00_{two}$ |
| 0 | 0 | 1 | 0 | 1 | $0 + 0 + 1 = 01_{two}$ |
| 0 | 1 | 0 | 0 | 1 | $0 + 1 + 0 = 01_{two}$ |
| 0 | 1 | 1 | 1 | 0 | $0 + 1 + 1 = 10_{two}$ |
| 1 | 0 | 0 | 0 | 1 | $1 + 0 + 0 = 01_{two}$ |
| 1 | 0 | 1 | 1 | 0 | $1 + 0 + 1 = 10_{two}$ |
| 1 | 1 | 0 | 1 | 0 | $1 + 1 + 0 = 10_{two}$ |
| 1 | 1 | 1 | 1 | 1 | $1 + 1 + 1 = 11_{two}$ |

$$Sum = a \oplus b \oplus CarryIn$$

$$CarryOut = a.b + CarryIn.(a \oplus b)$$

# Ripple-Carry Adder

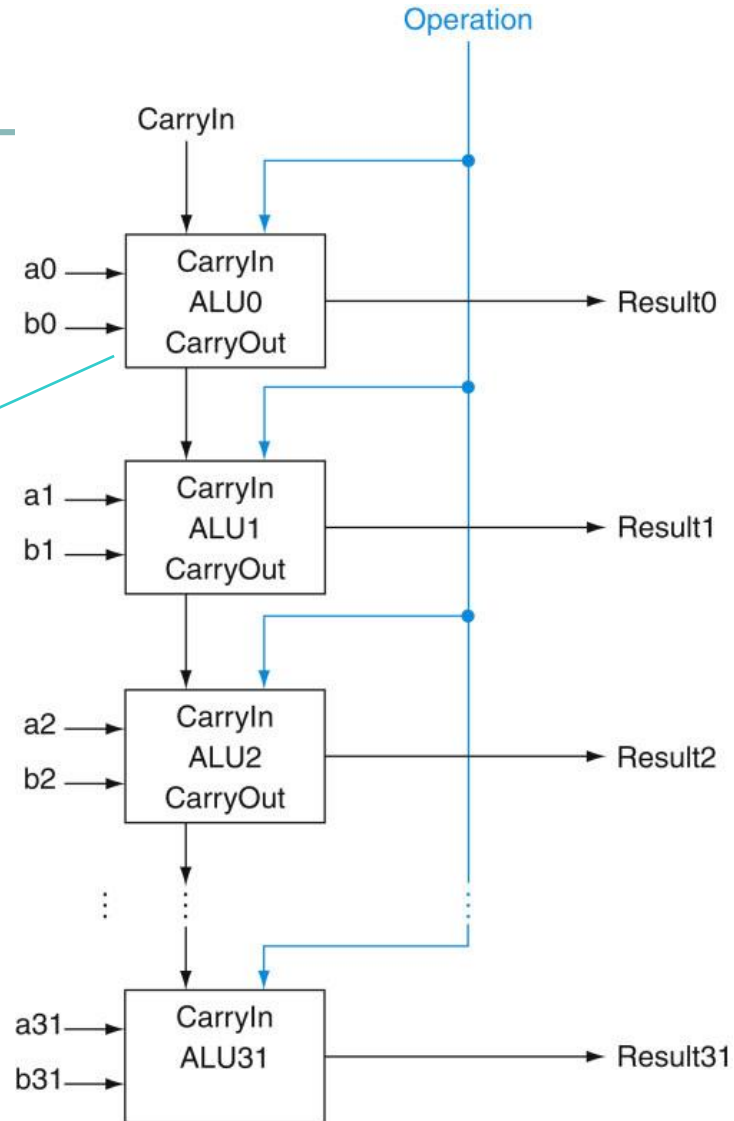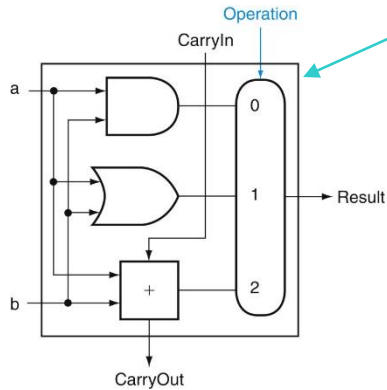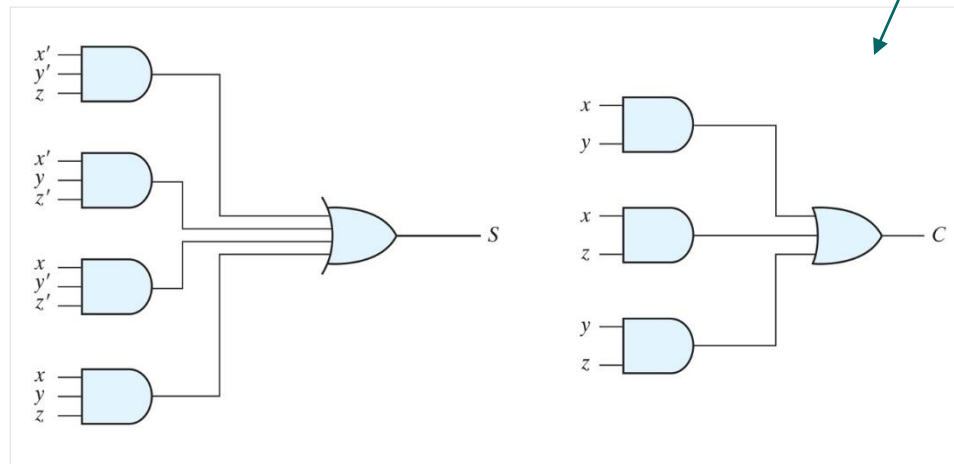# One-bit ALU

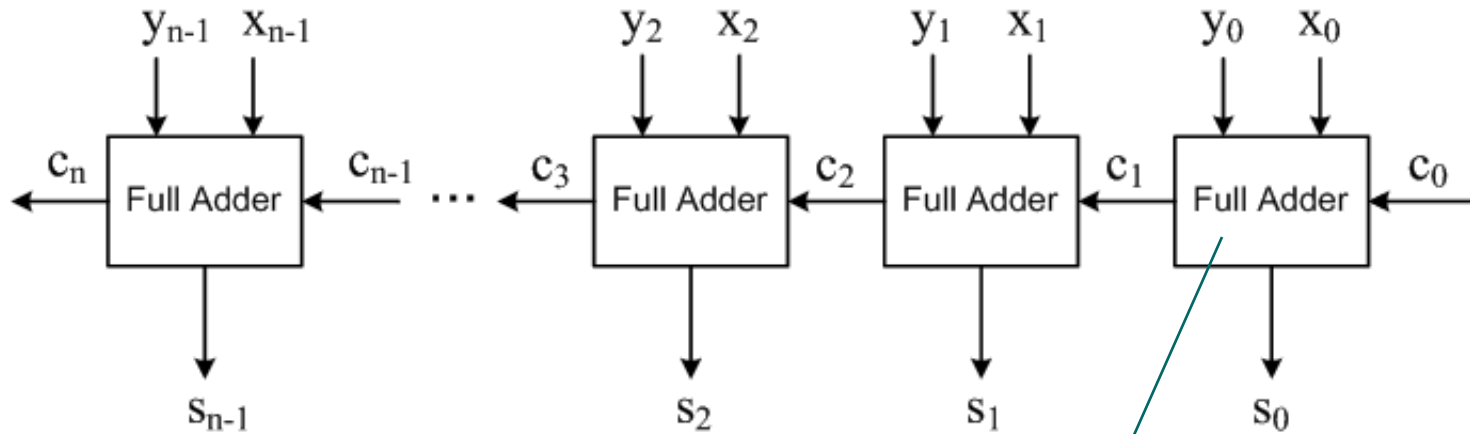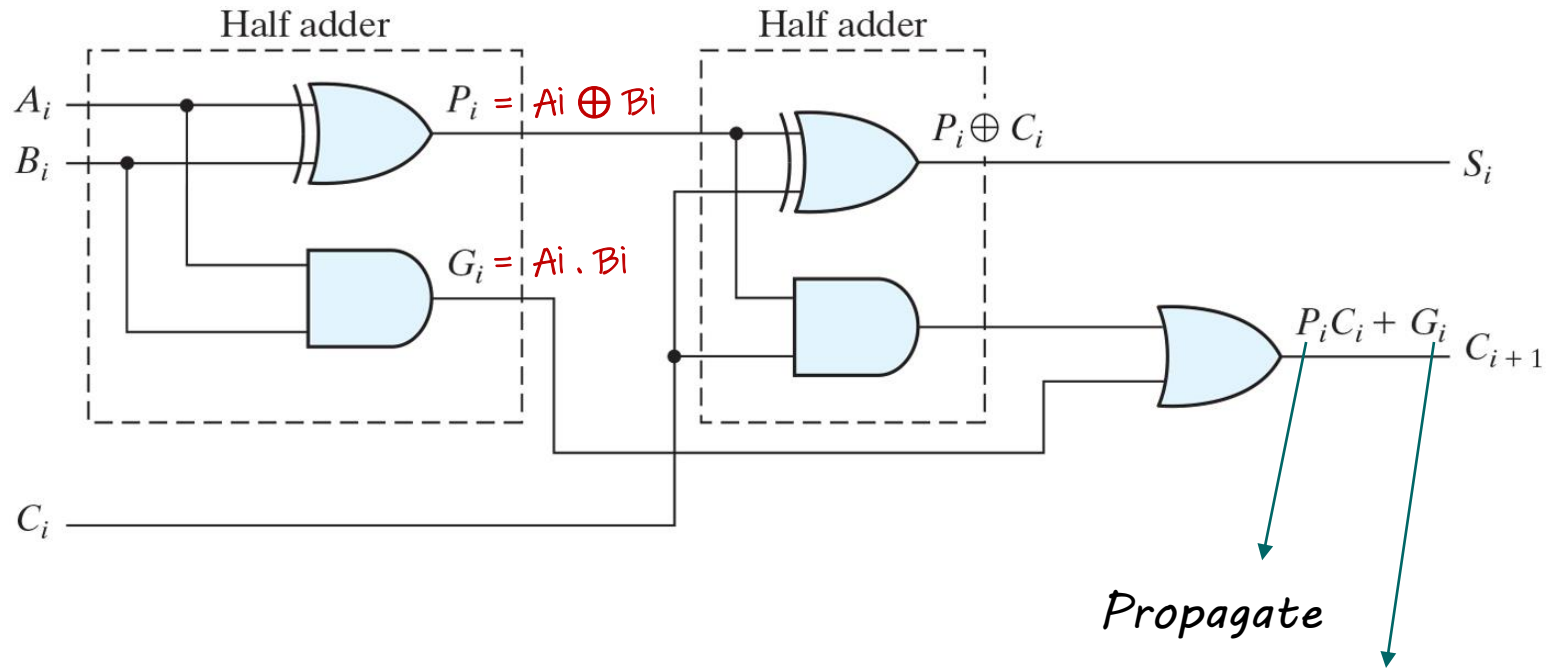*Performs AND, OR, and addition on a and b or a and b̲*

# 32-bit ALU

# *Reminder: Ripple-Carry Adder*

# Carry Generate / Propagate

Half adder        Half adder

$A_i$

$B_i$

$P_i = A_i \oplus B_i$

$G_i = A_i . B_i$

$P_i \oplus C_i$        $S_i$

$P_i C_i + G_i$        $C_{i+1}$

$C_i$

Propagate

Generate

$$Sum = (a \oplus b) \oplus CarryIn$$
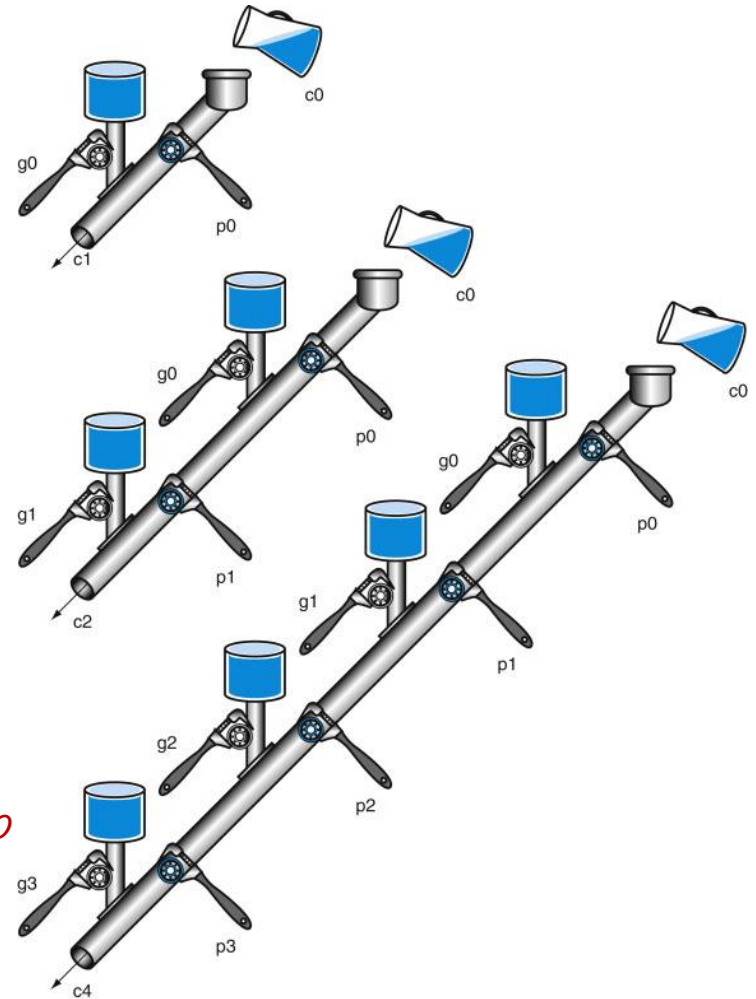
$$CarryOut = a.b + CarryIn.(a \oplus b)$$
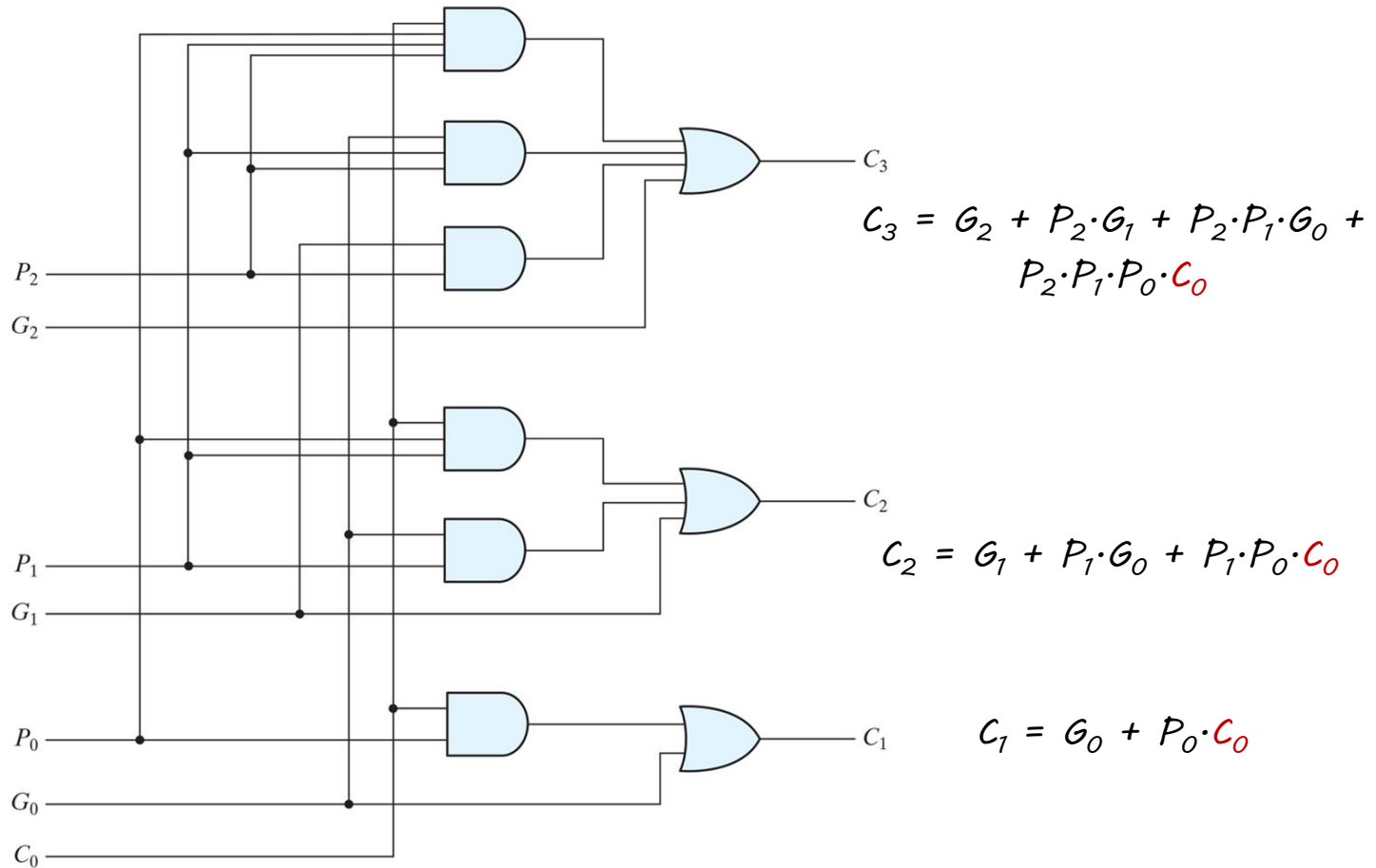
# CLA vs. Pipes

$c_1 = g_0 + p_0 \cdot c_0$

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

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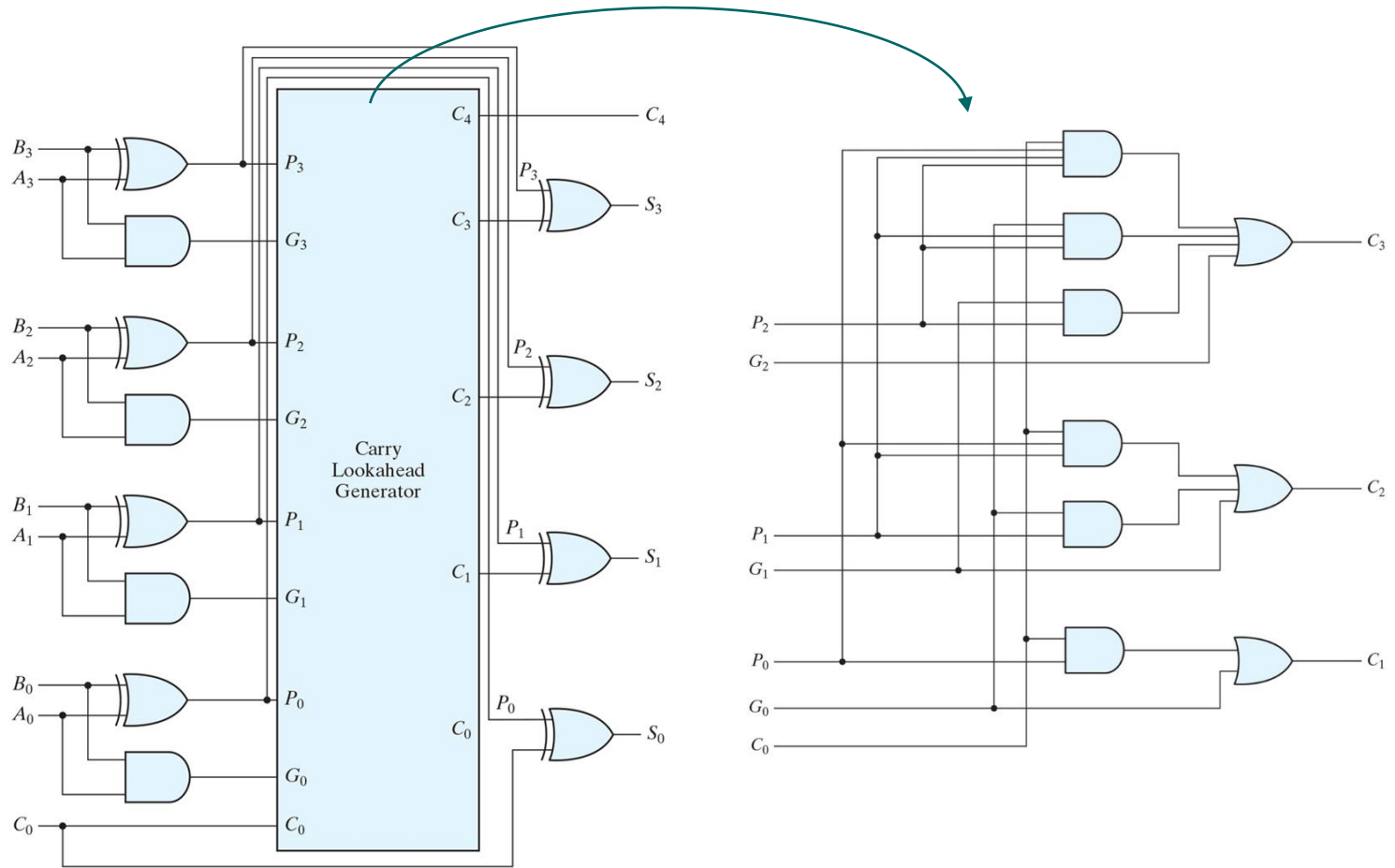$\quad\quad + p_2 \cdot p_1 \cdot p_0 \cdot c_0$

$c_4 = g_3 + p_3 \cdot g_2 + p_3 \cdot p_2 \cdot g_1$
$\quad\quad + p_3 \cdot p_2 \cdot p_1 \cdot g_0 + p_3 \cdot p_2 \cdot p_1 \cdot p_0 \cdot c_0$
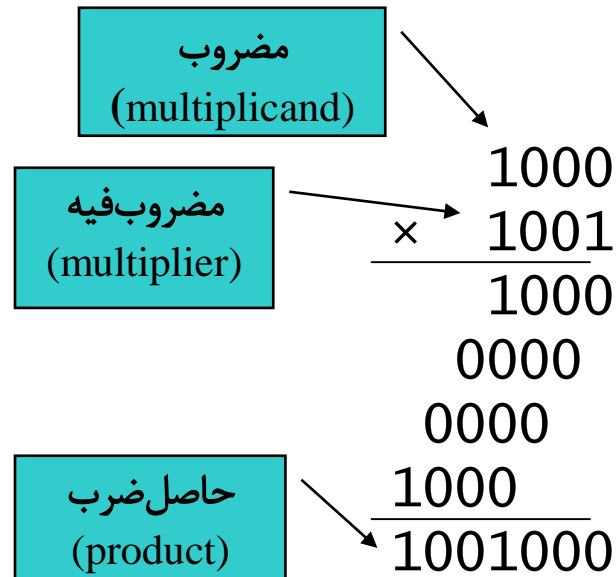
# CLA Generate/Propagate Circuit

$C_3 = G_2 + P_2 \cdot G_1 + P_2 \cdot P_1 \cdot G_0 + P_2 \cdot P_1 \cdot P_0 \cdot C_0$

$C_2 = G_1 + P_1 \cdot G_0 + P_1 \cdot P_0 \cdot C_0$

$C_1 = G_0 + P_0 \cdot C_0$

# CLA 4-bit Adder

# *Multiplication Approach (1ˢᵗ ver)*



مضروب
**(**multiplicand**)**

مضروب‌فیه
(multiplier)

حاصل‌ضرب
(product)

```
      1000
   ×  1001
      1000
     0000
    0000
   1000
  1001000
```

Length of product is the
sum of operand lengths



Multiplicand — Shift left — 64 bits

64-bit ALU

Product — Write — 64 bits

Control test

Multiplier Shift right — 32 bits

# *Multiplication Algorithm (1ˢᵗ ver)*

# *Multiplication (2ⁿᵈ ver)*

منصار و زبان کامپیوتر

**مضروب**
(multiplicand)

**مضروب‌فیه**
(multiplier)

**حاصل‌ضرب**
(product)

$$
\begin{array}{r}
1000 \\
\times\ 1001 \\
\hline
1000 \\
0000 \\
0000 \\
1000 \\
\hline
1001000
\end{array}
$$

Multiplicand

32 bits

32-bit ALU

Shift right

Product

Write

Control test

64 bits

Multiplier's Bits

# Division

مقسوم (dividend)

مقسوم‌علیه (divisor)

```
1001010 |1000
-1000    1001
    10
   101
  1010
 -1000
    10
```

خارج قسمت (quotient)

باقی‌مانده (remainder)

$$\text{Dividend} = \text{Quotient} \times \text{Divisor} + \text{Remainder}, \qquad |\text{Remainder}| < |\text{Divisor}|$$

# Real Numbers

○ **Numbers with Fractions:**

    ○ 3·14159...

    ○ 2·17

    ○ 0·0000001

    ○ $1·25 * 10^{-12}$

    ○ $1·43 * 10^{+12}$

○ **Representation in computers:**

    ● Fixed point

    ● Floating point

# Fixed-Point Representation

○ A real Example:

- $d_{23}d_{22}...d_1d_0 \cdot f_0f_1f_2f_3f_4f_5f_6f_7$

- 24-bit: integer bits

- 8-bit: fraction bits

○ Application

- Used in CPUs with no floating-point unit

  ○ Embedded microprocessors and microcontrollers

- Digital Signal Processing (DSP) applications

ساختار و زبان کامپیوتر

# *Fixed-Point Representation* *(cont.)*

○ Consider 5-Bit Representation

- $d_2d_1d_0\cdot f_1f_0$

- $(d_2 \times -2^2)+(d_1 \times 2^1)+(d_0 \times 2^0)+(f_1 \times 2^{-1})+(f_0 \times 2^{-2})$

○ Largest positive number?

○ Smallest positive number?

○ Largest magnitude negative number?

○ Smallest magnitude negative number?

# Fixed-Point Representation (cont.)

○ Arithmetic:

out of range (overflow)

- 011.11 + 011.11 = 111.10

out of range (underflow)

- 010.10 × 000.10 = 000001.0100

- 000.01 × 000.01 = 000000.0001

- 011.01 × 011.01 = **?**

# *Fixed-Point Representation* (cont.)

○ *Arithmetic:*

out of range
(overflow)

- *011.11 + 011.11 = 111.10*

out of range
(underflow)

- *010.10 × 000.10 = 000001.0100*

- *000.01 × 000.01 = 000000.0001*

- *011.01 × 011.01 = 001010.1001*

Both
overflow &
underflow
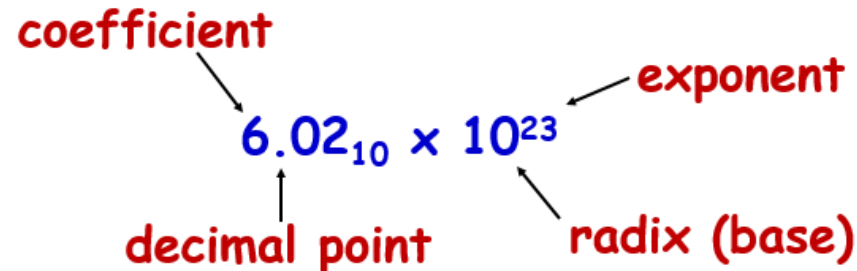
# *Fixed-Point Representation* *(cont.)*

○ Pros

  ● Simple hardware

  ● Fast computation

  ● Different precisions at different applications

    ○ 24bits/8bits , 18bits/14bits, 8bits/24bits

○ Cons

  ● Low precision

  ● Small range

# *Scientific Notation (Decimal)*



○ *Normalized Form:*

- *Exactly one non-zero digit to left of decimal point*

○ *Alternatives to representing 0.0000000012*

- *Normalized:*     $1.2 \times 10^{-9}$

- *Not normalized:*   $0.12 \times 10^{-8}$, $12.0 \times 10^{-10}$

# *Normalized Scientific Notation (Binary)*

**mantissa (fraction)**

**exponent**

$$1.0_{two} \times 2^{-1}$$

**binary point**

**radix (base)**

# *Floating-Point Notation*

○ *Floating Point Notation Consists of:*

- *Fraction (F): 23 bits*

- *Exponent (E): 8 bits*

- *Fraction Sign bit (S)*

- *Also called, single precision floating-point*

○ $N = (-1)^S \times (1+F) \times 2^E$

| 31 | 30 | ... | 24 | 23 | 22 | 21 | ... | 1 | 0 |
|----|----|----|----|----|----|----|----|----|----|
| S | Exponent | | | | Fraction | | | | |

# *Floating-Point Notation (cont.)*

○ *Pros (compared to fixed-point)*

- *Very Wide Range*
- *More precision bits*

○ *Cons (compared to fixed-point)*

- *Arithmetic operation more complicated*
- *HW more complicated*

| 31 | 30 | ... | 24 | 23 | 22 | 21 | ... | 1 | 0 |
|----|----|----|----|----|----|----|----|----|----|
| S | Exponent | | | | Fraction | | | | |

# *Floating-Point Notation* (cont.)

○ $N = (-1)^S \times (1+F) \times 2^E$

○ *Precision* versus *Range*

- *More precision ➔ smaller range?*

- *Wider range ➔ less precision?*

○ True for fixed-point

○ Not necessarily correct for floating point

| 31 | 30 | ... | 24 | 23 | 22 | 21 | ... | 1 | 0 |
|----|----|-----|----|----|----|----|-----|---|---|
| S | Exponent | | | | Fraction | | | | |

# *Floating-Point Notation* *(cont.)*

○ *Overflow:*

- *Exponent too large to fit in "Exponent" field*

○ *Underflow:*

- *Non-zero fraction so small to represent*

- *Negative exponent too large to fit*

| 31 | 30 | ... | 24 | 23 | 22 | 21 | ... | 1 | 0 |
|----|----|-----|----|----|----|----|-----|---|---|
| S | Exponent | | | | Fraction | | | | |

# IEEE 754 – Single Precision

○ *Signed-magnitude* notation for fraction (mantissa)

○ *Biased* (Excess $2^{n-1}-1$) notation for exponent

○ $E_{min}$=00000001

○ $E_{max}$=11111110

| 31 | 30 | … | 24 | 23 | 22 | 21 | … | 1 | 0 |
|----|----|---|----|----|----|----|---|---|---|
| S | Exponent | | | | Fraction | | | | |

$$N = (-1)^S * (1 + F) * 2^E$$

○ E=00000000 reserved for zero

○ E=11111111 reserved for *infinity* & *NaN*

○ Smallest positive no: 1.17549435 E-38

○ Largest positive no: 3.4028235 E38

# IEEE 754 – Double Precision

○ Two words long (64 bits)

○ Reduced chances of overflow/underflow

○ Format

- Sign bit (S)

- Fraction (F): 52 bits

- Exponent (E): 11 bits

○ A bias of 1023 in Exponential part

# *More on IEEE 754 Standard*

○ *Single precision (32bits)/Double precision (64bits)*

○ *Normalized/ Denormalized forms*

○ *Standard definitions for <span style="color:red">zero</span>, <span style="color:red">infinity</span>, <span style="color:red">NaN</span>*

○ *Check:* *https://www.h-schmidt.net/FloatConverter/IEEE754.html*

| Single precision | | Double precision | | Object represented |
|---|---|---|---|---|
| Exponent | Fraction | Exponent | Fraction | |
| 0 | 0 | 0 | 0 | 0 |
| 0 | Nonzero | 0 | Nonzero | ± denormalized number |
| 1–254 | Anything | 1–2046 | Anything | ± floating-point number |
| 255 | 0 | 2047 | 0 | ± infinity |
| 255 | Nonzero | 2047 | Nonzero | NaN (Not a Number) |

# Denormalized Forms

○ *An attempt to squeeze every last bit of precision from a floating-point operation*

○ *Smallest pos. single precision normalized no:*

- *$1.00000000000000000000000 \times 2^{-126}$*

○ *Smallest single precision denormalized no:*

- *$0.00000000000000000000001 \times 2^{-126} = 1.0 \times 2^{-149}$*

| 31 | 30 | ... | 24 | 23 | 22 | 21 | ... | 1 | 0 |
|----|----|-----|----|----|----|----|-----|---|---|
| S | Exponent | | | | Fraction | | | | |

# Concluding Remarks

- *Bits have no inherent meaning*

  - *Interpretation depends on the operations applied*

- *Computer representations of numbers*

  - *Finite range and precision*

  - *Need to account for this in programs*

- *Bounded range and precision*

  - *Operations can overflow and underflow*

# Outlines

○ *Weighted Number System*

○ *Signed Number Representation*

- *2's Compliment/ 1's Compliment*

- *Signed-Magnitude Notation*

- *Biased Notation*

○ *Arithmetic Operations*

- *Addition/ Subtraction/ Multiplication/ Division*

○ *Real Numbers*

- *Fixed Point / Floating Point Representation*

- *IEEE 754 Standard*