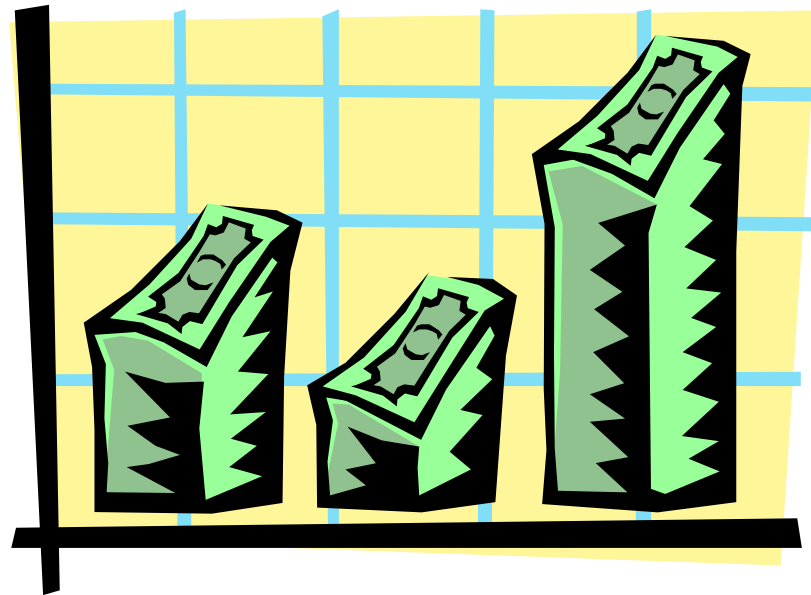# Computer Arithmetic

◆ ALU design

◆ Multiplication – Booth Encoding

◆ Floating-Point Numbers

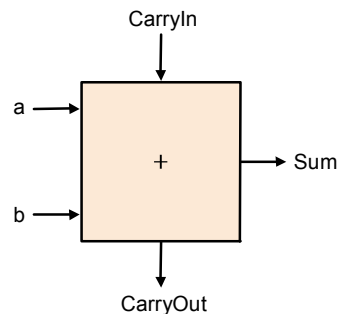# Full Adder

## Recall

## Adding 2's Complement Binary Numbers

$$P_3 \, P_2 \, P_1 \, P_0$$
$$Y_3 \, Y_2 \, Y_1 \, Y_0$$
$$+-----------$$

CarryIn

a

+ → Sum

b

CarryOut

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

$sum = a \, xor \, b \, xor \, c_{in}$

**Note 2-gate delays**

◆ N-bit ripple-carry adder: a series connection of n FA's

◆ Sign-extension: repeating sign bit to fill up higher bits,
  ○ e.g. 16 to 32bits

◆ Overflow:
  Carry into MSB(Most Significant Bit) is a sign-bit
  Carry out of MSB is also a sign-bit for (N+1) bit representation
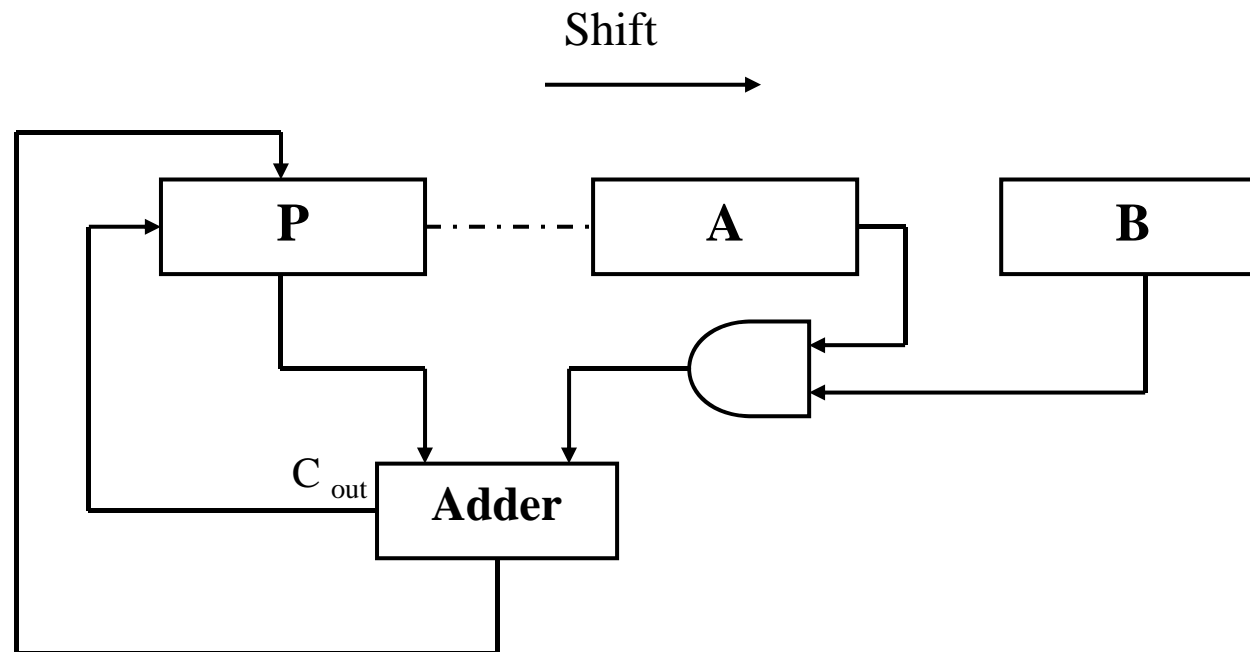      They should be the same, otherwise OVERFLOW.

# Multiplication

◆ More complicated than addition

    ○ accomplished via shifting and addition

◆ More time and more area

◆ Let's look at a version based on grade-school algorithm

```
    0010   (multiplicand)
  x 1011   (multiplier)
```

◆ Negative numbers:  convert and multiply

# Unsigned Integer Multiplication

- ◆ **Add shift.** **B x A => P || A.**

- ◆ LSB of A = 1 $\rightarrow$ P = P + B.

- ◆ Shift right  Carry & P & A  by 1 bit.

Shift

$\longrightarrow$

| | | |
|---|---|---|
| **P** | **A** | **B** |

$C_{out}$   **Adder**

# Unsigned Integer Multiplication - Example

◆ 0010x1011: BxA=P

| P | A | B |
|---|---|---|
| 0000 | 1011 | 0010 |
| + 0010 | | |
| 0010 | 1011 (and then shift right) | |
| 0001 | 0101 | |
| + 0010 | | |
| 0011 | 0101 (and then shift right) | |
| 0001 | 1010 | |
| + 0000 | | |
| 0001 | 1010 (and then shift right) | |
| 0000 | 1101 | |
| + 0010 | | |
| 0010 | 1101 (and then shift right) | |
| 0001 | 0110 | final answer in P||A. |

```
      0010
x    1011
----------
      0010
    0010
  0000
0010
-----------
0010110
```

# Unsigned Integer Multiplication

◆How about negative number multiplication?

$\rightarrow$      X . Y ; |X| . |Y|. <sup>g1</sup>

Check sign(X) $\neq$ sign(Y)    $\rightarrow$   - |X| . |Y|.

◆Checking sign is overhead!

g1        just do the previous slide algorithm with 2's complement, then it's OK

Booth is to reduce overhead not just sign bit check but also reducing additions.
ghlee, 2015-03-31

# Booth's Algorithm for 2's Complement. Numbers

◆ Consider $Y \cdot X$ , where X = 0111

Since, 0111 = 1000 – 0001,

$Y \cdot X = Y \cdot (1000) - Y \cdot (0001)$

◆ Subtract Y and Add 8 times Y (i.e. shift 3 times) . This req. 1 add, 1 sub. & 4 shifts.

◆ Recording the above 0111 = 1000 – 0001= 100 $\overline{1}$ is called **Booth Recording** (in radix 2)

# Booth Recording

◆ Let $X = X_{n-1} \cdots X_0$ for Y. X

At $i^{th}$ step of multiplication , ( $X_{-1} = 0$)

| $X_I$ | $X_{I-1}$ | Recording | Algorithm | Note |
|---|---|---|---|---|
| 1 | 0 | $\bar{1}$ | -Y & shift | Beginning of 1's |
| 1 | 1 | 0 | Shift | Middle of 1's |
| 0 | 1 | 1 | +Y & shift | End of 1's |
| 0 | 0 | 0 | shift | Middle of 0's |

# Example

Y. X → P || X   with Y = -5 ,  X = -6   in 4 bit

```
  P              X
0 0 0 0        1 0 1 0 0 0      ; X 0 = X -1 = 0  ( shift right )
0 0 0 0        0 1 0 1 0        ; X 1 = 1 ≠ X 0 = 0 ( subtract)
0 1 0 1        0 1 0 1          ; & shift right
0 0 1 0        1 0 1 0 1        ; X 2 = 0 , X 1 = 1 : Add
1 0 1 1

+ ----------------------
1 1 0 1        1 0 1 0          ; shift right
1 1 1 0        1 1 0 1 0        ; X 3 = 1 ,  X 2 =0 : (subtract)
0 0 1 1        1 1 0 1          ; & shift right
0 0 0 1        1 1 1 0

--------------------------     ; 30 = -5 x -6
```

# Example

$Y = 5$ , $X = -6$   in 4 bit

```
  P                X
0 0 0 0          1 0 1 0 0        ; X₀ = X₋₁ = 0  ( shift right )
0 0 0 0          0 1 0 1 0        ; X₁ = 1 ≠ X₀ = 0 ( subtract)
1 0 1 1          0 1 0 1          ; & shift right (with sign-bit)
1 1 0 1          1 0 1 0 1        ; X₂ = 0 , X₁ = 1 : Add
0 1 0 1
+ -----------------------
0 0 1 0          1 0 1 0          ; shift right
0 0 0 1          0 1 0 1 0        ; X₃ = 1 ,  X₂ =0 : (subtract)
1 0 1 1          0 1 0 1          ;
1 1 0 0          0 1 0 1          ; shift right
  --------------------------
1 1 1 0          0 0 1 0          ; -30 = 5 x -6
```

# Booth Recording

◆ Booth recording takes care of negative multiplier

◆ In the example

◆ X=-6=1010 $\longrightarrow$ $\bar{1}1\ \bar{1}0$

$\qquad\qquad\qquad\qquad\qquad$ Subtract 2 times

$\qquad\qquad\qquad\qquad\qquad$ Add 4 times

$\qquad\qquad\qquad\qquad\qquad$ Subtract 8 times

Note

Hope to have less add/sub:

alternating 0 and 1 makes worst case g2

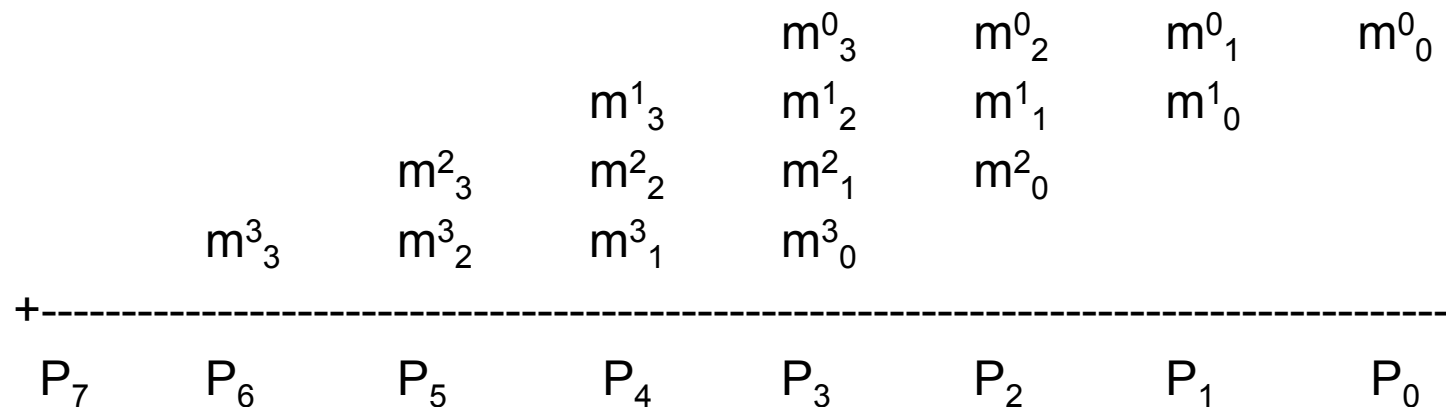g2          so $-6=1010$ is the worst case!

# Faster Array Multiplier

Instead of doing the adding partial products linearly,
i.e. one after another partial product,
can be done in parallel by grouping partial products

Each of group has three partial products
        Wallace Tree Multiplier
Each group of two partial products
        Binary Tree Multiplier

|  |  |  |  | $m^0_3$ | $m^0_2$ | $m^0_1$ | $m^0_0$ |
|---|---|---|---|---|---|---|---|
|  |  |  | $m^1_3$ | $m^1_2$ | $m^1_1$ | $m^1_0$ |  |
|  |  | $m^2_3$ | $m^2_2$ | $m^2_1$ | $m^2_0$ |  |  |
|  | $m^3_3$ | $m^3_2$ | $m^3_1$ | $m^3_0$ |  |  |  |

$+$----------------------------------------------------------------------------------

| $P_7$ | $P_6$ | $P_5$ | $P_4$ | $P_3$ | $P_2$ | $P_1$ | $P_0$ |
|---|---|---|---|---|---|---|---|

# Division

◆ Can be done in reverse of multiplication

of shift and add

shift-left and subtract

Faster Division uses

◆ the fact x/y = x*(1/y) iterative guess of 1/y

◆ multiple quotient bits with precalculated table: SRT division

read Text 3.4(p.178~182)

# Floating Point Numbers

◆ We need a way to represent

- o numbers with fractions, e.g., 3.1416

- o very small numbers, e.g., .000000001

- o very large numbers, e.g., 31,557,600,000,000

◆ Normalized Representation:

- o No leading zeros and one significant digit before decimal
  point:        (sign)m.xyz…. X $10^e$

e.g.    $+3.15576 \times 10^9$
        sign: +
        significand: 3.15576
        exponent: 9
        number base: 10, i.e decimal

# Floating Point Numbers

For computer floating point:

Number base: 2 ($2^n$)

Sign: 0 (+) or 1(-)

Exponent: actual exponent + bias

Significand:   binary number (with Sign bit)

$$(-1)^{sign} \cdot \text{significand} \cdot 2^{exponent}$$

- bias : to make smallest exponent 0
    comparing exponents and checking 0 more
    convenient with biased representation
- more bits for significand gives more accuracy
- more bits for exponent increases range
- Mantissa (or Fraction): normalized with no significant
digit left to the decimal point $0.315576 \times 10^{10}$

# Floating Point Numbers

- $(-1)^{sign} \times (significand) \times 2^{exponent}$

- Issues
    - Format : which is which and how many bits, number base ($2^n$)
    - Handling of exceptions
        - Overflow, Underflow, very small numbers
    - Rounding
        - How to round up numbers?

## IEEE 754 floating point standard

# IEEE 754 floating point standard

o single precision: 8 bit exponent, 24 bit significand (23+1)

o double precision: 11 bit exponent, 53 bit significand (52+1)

o Do not represent 1 before the binary point: It's always 1!

o Exponent Bias: to make smallest negative exponent 0

    v 127 for single precision; 0 ~ 255 with 255 (all 1's) reserved

    v 1023 for double precision

        v all 1's for exponent: reserved for special cases

## Note: some specials in IEEE754

o zero: zero significand and zero exponent

o 0/0: NaN (nozero significand, all 1's for exponent)

o n/0: infinity (zero significand, all 1's for exponent)


o Denormal: close to zero (nonzero significand with 0 exponent)

    v Without denormal, more possibility of Underflow


See Figure 3.13 Text p.187

# IEEE 754 floating-point standard

◆ Example:

- o decimal: $-.75 = -3/4 = -3/2^2$
- o binary: $-.11 = -1.1 \times 2^{-1}$
- o floating point: exponent = -1 + bias (= 127) = 126 = 01111110

- o IEEE single precision:
  1**01111110**10000000000000000000000

  sign exponent(w.bias) fraction(+1=significand)

# Floating Point Complexities

◆ **Operations are somewhat more complicated**

◆ Needs to compare exponents and align them by shifting significand

◆ Needs to work on significand and exponent separately

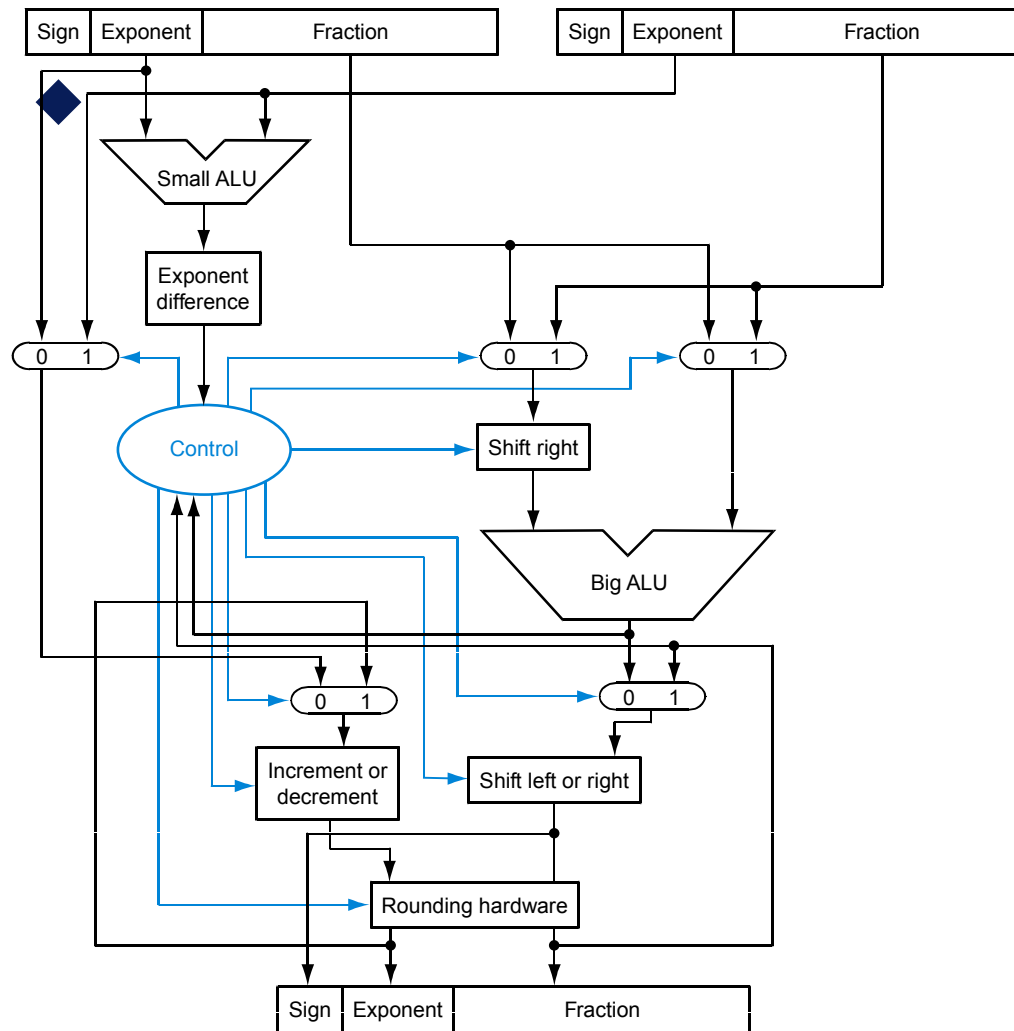◆ Needs to normalize the result for correct representation

e.g. (read Text p. 198)
Add 0.5 and –0.4375 in binary Floating Point with 4-bit significand

$0.5 = 1.000 \times 2^{-1}$
$-0.4375 = -7/16 = -1.110 \times 2^{-2}$

1. Shift lesser exponent operand: $-0.111 \times 2^{-1}$
2. Add significands: $1.000 - 0.111 = 0.001$
3. Normalize $1.000 \times 2^{-3} \times 2^{-1} = 1.000 \times 2^{-4}$

# Floating point addition

# Floating Point Complexities

◆In addition to overflow we can have "underflow"

◆Accuracy can be a big problem

   ○ IEEE 754 keeps two extra bits, guard and round – representing two additional bits for better precision plus one sticky bit (the bit moved out when shifted for exponent alignment) representing any non-zero bit beyond the precision (read Text p.206&208)

   ○ four rounding modes: up/down, truncate, nearest even.

   ○ positive divided by zero yields "infinity"

   ○ zero divide by zero yields "not a number"

   ○ other complexities

```
0.5
0.51
0.501
0.5000003
```

◆Implementing the standard can be tricky

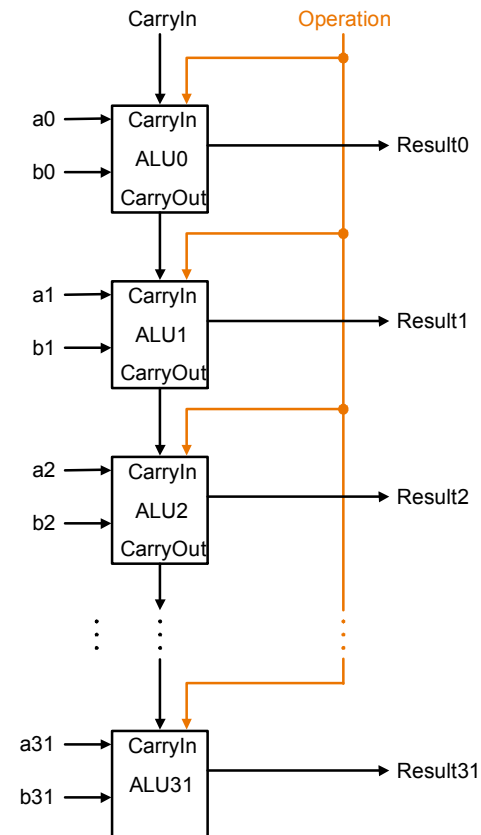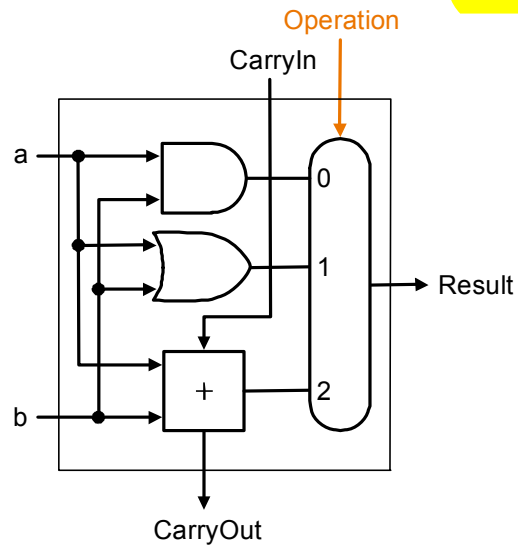   ○ see text for description of 80x86 and Pentium bug!

# MIPS ALU

◆ We can build an ALU to support the MIPS instruction set

   o key idea:  use multiplexor to select the output we want

   o we can efficiently perform subtraction using two's complement

   o we can replicate a 1-bit ALU to produce a 32-bit ALU

◆ Important points about hardware

   o all of the gates are always working

   o the speed of a gate is affected by the number of inputs to the gate

   o the speed of a circuit is affected by the number of gates in series
      (on the "critical path" or the "deepest level of logic")

◆ Our primary focus:  comprehension,  however,

   o Clever changes to organization can improve performance
      (similar to using better algorithms in software)

   o we'll look at examples for addition and multiplication

# Building a 32 bit ALU

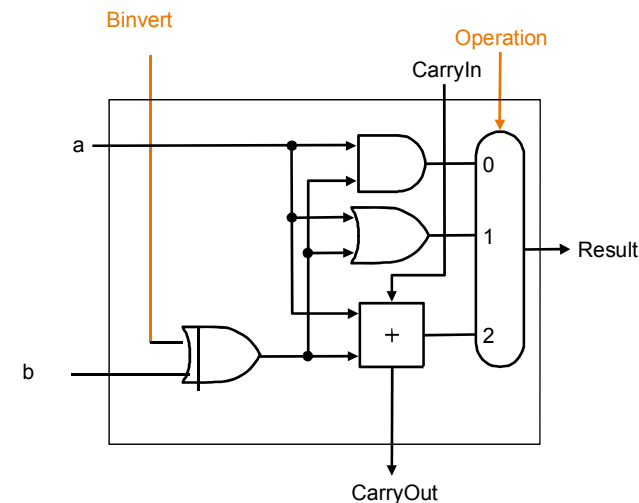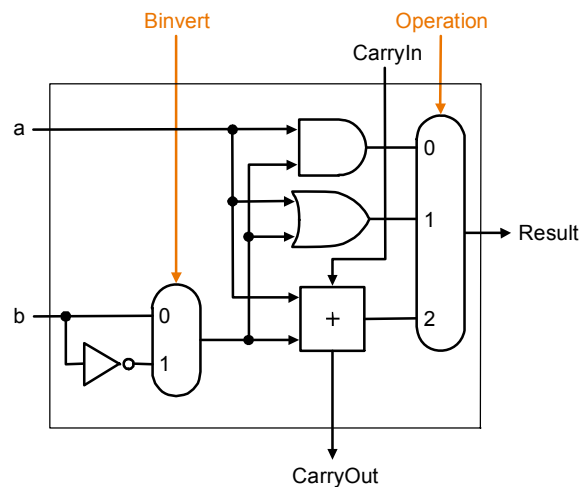Support AND, OR, Add/Sub

Key: Mux

# Tailoring the ALU to the MIPS

◆ Adding more and select through Mux per opcode

o Need to support test for equality (beq $t5, $t6, 0x100)

v use subtraction:  (a-b) = 0 implies a = b

| Name | Format | Example | | | | | | Comments |
|------|--------|---------|---|---|---|---|---|----------|
| | | 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits | |
| add | R | 0 | 2 | 3 | 1 | 0 | 32 | add $1,$2,$3 |
| sub | R | 0 | 2 | 3 | 1 | 0 | 34 | sub $1,$2,$3 |
| addi | I | 8 | 2 | 1 | 100 | | | addi $1,$2,100 |
| addu | R | 0 | 2 | 3 | 1 | 0 | 35 | subu $1,$2,$3 |
| and | R | 0 | 2 | 3 | 1 | 0 | 36 | and $1,$2,$3 |
| or | R | 0 | 2 | 3 | 1 | 0 | 37 | or $1,$2,$3 |
| lw | I | 35 | 2 | 1 | 100 | | | lw  $1,100($2) |
| sw | I | 43 | 2 | 1 | 100 | | | sw $1,100($2) |
| beq | I | 4 | 1 | 2 | 25 | | | beq $1,$2,100 |
| j | J | 2 | 2500 | | | | | j 10000 |

# What about subtraction  (a – b)  ?

◆ Two's complement approach:  just negate b and add.

◆ A clever solution: negate+1 = 2's complement
   negate signal = carry-in for least significant bit position



$$B \text{ xor } 1 = \text{not } B$$

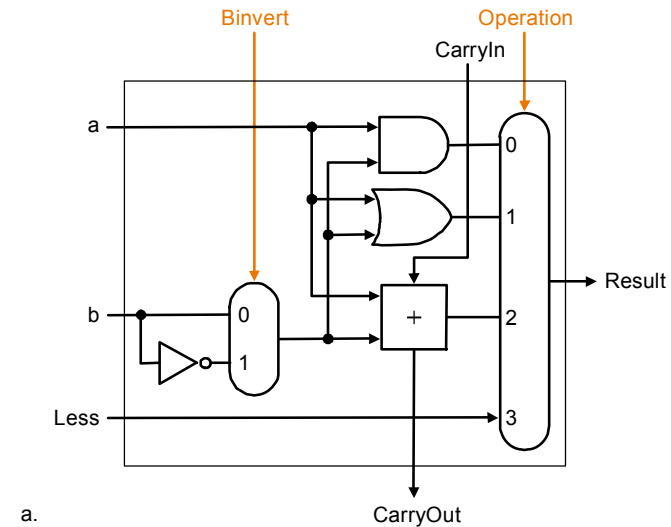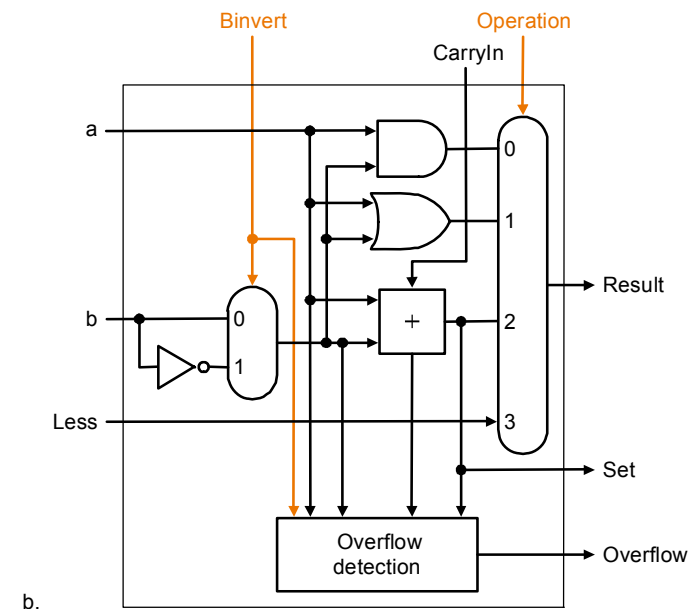# Tailoring the ALU to the MIPS

◆ **Need to support test for equality (beq $t5, $t6, L)**

  ○ use subtraction:  (a-b) = 0 implies a = b

◆ **to support the set-on-less-than instruction**

  **(slt $r1, $r2, $r3)**

  ○ remember:  slt is an arithmetic instruction

  ○ produces a 1 if rs < rt and 0 otherwise

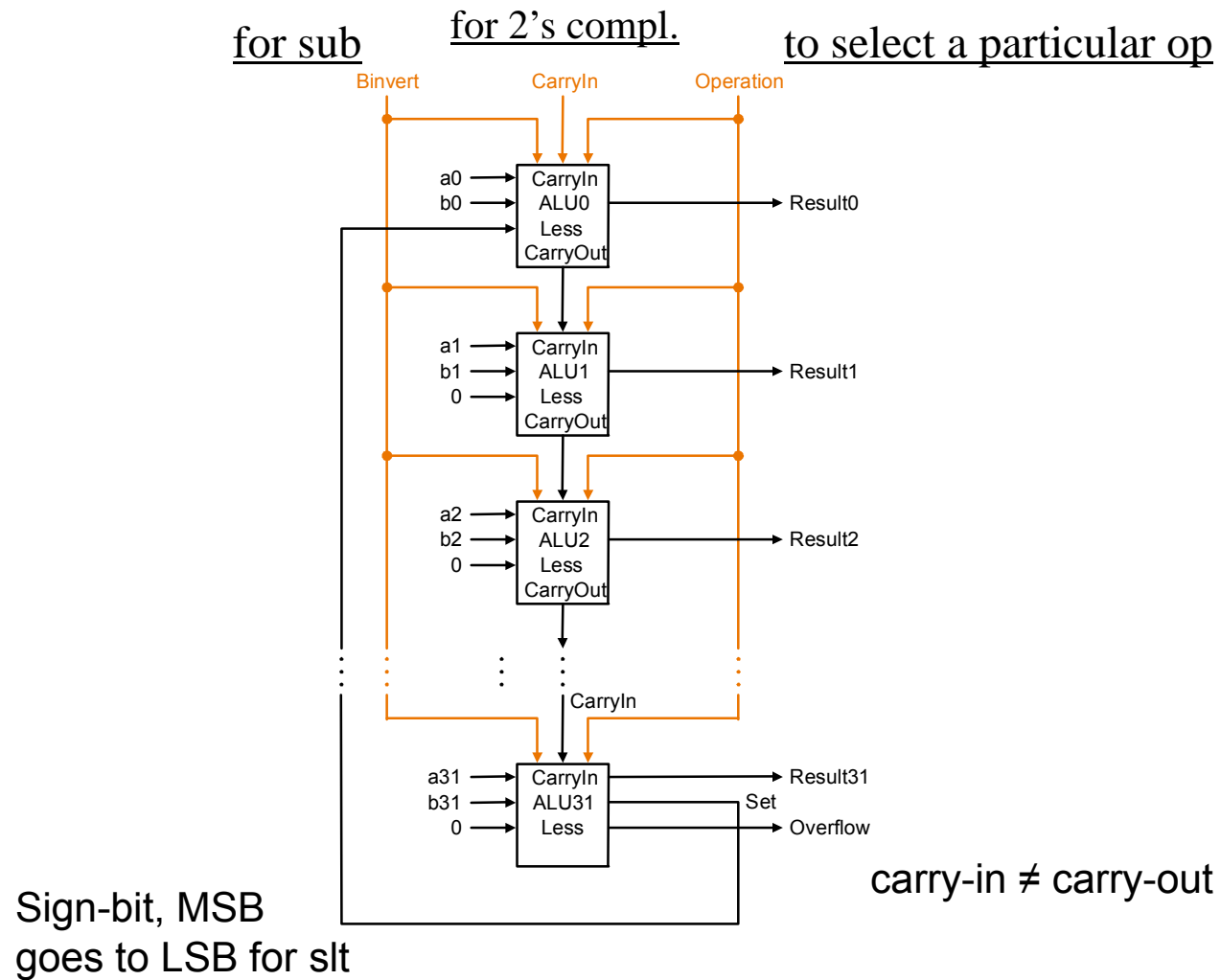  ○ use subtraction:  (a-b) < 0 implies a < b

# Supporting slt

◆ Can we figure out the idea?

a.

◆ If a<b then (a-b)<0

So, if after subtraction, the result is Negative, i.e. MSB=1

Then

LSB of Result =MSB

and other bits of Result = 0

b.

# Supporting slt & overflow



for sub     for 2's compl.     to select a particular op

Binvert    CarryIn    Operation

a0 → CarryIn ALU0 Less CarryOut → Result0

a1 → b1 → 0 → CarryIn ALU1 Less CarryOut → Result1

a2 → b2 → 0 → CarryIn ALU2 Less CarryOut → Result2

CarryIn

a31 → b31 → 0 → CarryIn ALU31 Less → Result31, Set, Overflow

Sign-bit, MSB goes to LSB for slt

carry-in ≠ carry-out

# Test for equality

• *Note:  zero is 1 when the result is zero!*

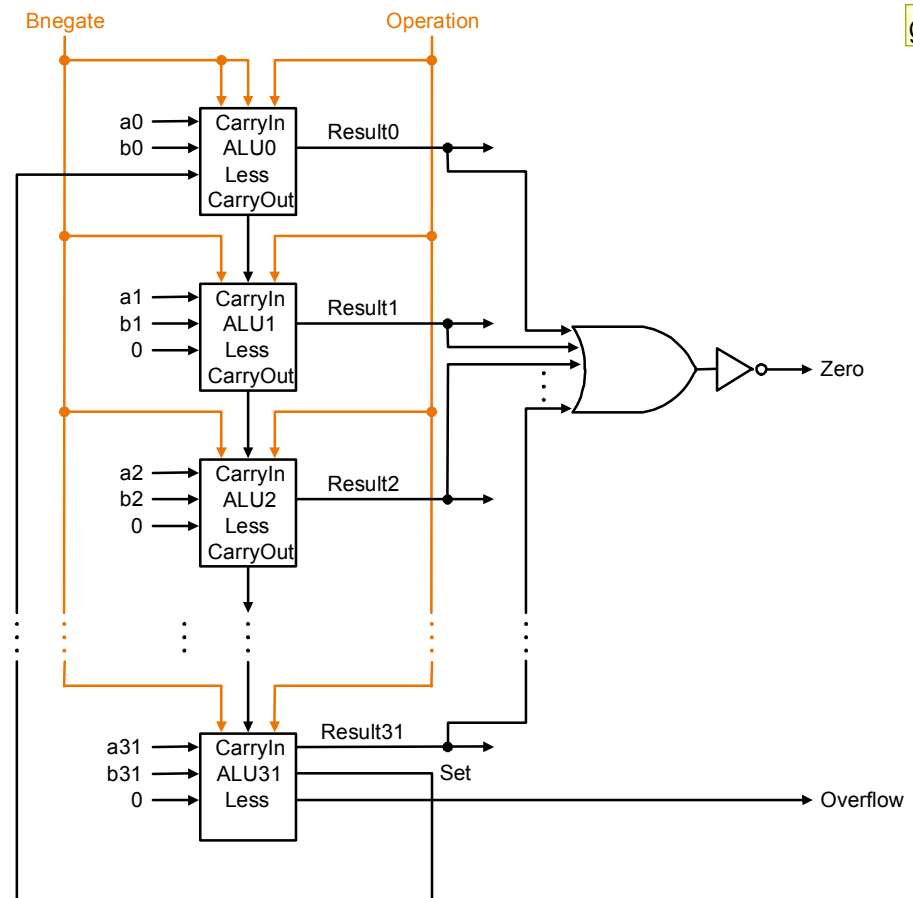◆ Notice control lines:

```
0000 = and
0001 = or
0010 = add
0110 = subtract
0111 = slt
```

(see Sec. 4.4. p.247)

g3          for beq, subtraction yields 0. all bit positions will be o. so negating them makes 1.
for sub and slt we need subtraction, so Bnegate will be 1 and it will feed to the lsb as a carry-in to have two's comlement.
ghlee, 2015-03-31

# Chapter Three Summary

◆ Computer arithmetic is constrained by limited precision

- May Not be asscociative (a+b)+c ≠ a+(b+c)

◆ Bit patterns have no inherent meaning but standards do exist
  - two's complement
  - IEEE 754 floating point

◆ Computer instructions determine "meaning" of the bit patterns

For Fast Adder, read the supplemental slides posted!

**We are ready to move on (and implement the processor)!**