# CSE261: Computer Architecture

## 6. Arithmetic for Computers (2)

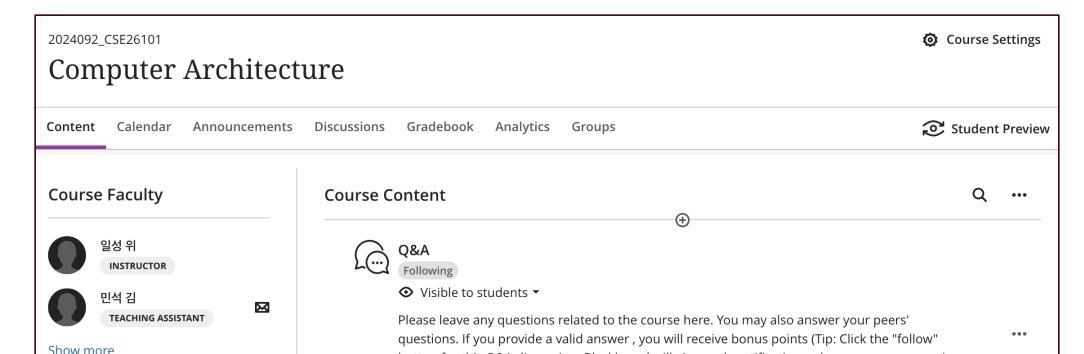Seongil Wi

# Notification: HW1

- Due date: **Today, 11:59PM**

# Notification: Midterm Exam

- Oct. 24 (Thursday)
- Class Time (1h 15m)

- T/F problems + Computation problems + Descriptive problems
- Closed book

# Notification: Q&A

- Avoid asking questions about issues that have "clear" answers

- Avoid asking duplicated questions

- Your instructors are not here to debug!)

- **Participation point**: If you answer other students' questions (before my answer), and if the answer is valid one, you will receive bonus point!

---

2024092_CSE26101

## Computer Architecture

⚙ **Course Settings**

Content   Calendar   Announcements   Discussions   Gradebook   Analytics   Groups

🔄 Student Preview

### Course Faculty

일성 위
**INSTRUCTOR**

민석 김   ✉
**TEACHING ASSISTANT**

Show more

### Course Content

🔍   ⋯

⊕

💬 **Q&A**
Following

👁 Visible to students ▾

Please leave any questions related to the course here. You may also answer your peers' questions. If you provide a valid answer , you will receive bonus points (Tip: Click the "follow"

⋯

# Integer Division

# Division

- Same as multiplication: just need **add/sub** and **shift operation**!

# Division

- Let's think about the division in grade school level

**Quotient**

**Divisor**

**Dividend**

**Remainder**

$$\begin{array}{r} 9 \\ 8 \overline{)74} \\ -72 \\ \hline 2 \end{array}$$

$$\begin{array}{r} 01001 \\ 1000 \overline{)01001010} \\ 0100 \\ 01001 \\ -\ 1000 \\ \hline 10 \\ 101 \\ 1010 \\ -\ 1000 \\ \hline 10 \end{array}$$

# Division Algorithm in HW Manner

**Quotient**

**Divisor**

**Dividend**

**Remainder**

```
loop len(divisor)+1 times:
  if divisor ≤ remainder:
    quotient bit = 1;
    remainder = remainder – divisor;
  else:
    quotient bit = 0;
  bring down next dividend bit to remainder;
```

$$
\begin{array}{r}
01001 \\
1000 \,\overline{)\,01001010} \\
0100 \\
01001 \\
-\ 1000 \\
\hline
10 \\
101 \\
1010 \\
-\ 1000 \\
\hline
10
\end{array}
$$

# Division Algorithm in HW Manner

Need 5 loops for completing the operation

Quotient

Divisor

Dividend

```
loop len(divisor)+1 times:
  if divisor ≤ remainder:
    quotient bit = 1;
    remainder = remainder – divisor;
  else:
    quotient bit = 0;
  bring down next dividend bit to remainder;
```

$$
\begin{array}{r}
01001 \\
1000\,\overline{)\,01001010} \\
0100 \\
01001 \\
-\ 1000 \\
\hline
10 \\
101 \\
1010 \\
-\ 1000 \\
\hline
10
\end{array}
$$

Remainder

# Division Algorithm in HW Manner

Quotient

Divisor **1000** 01001010

0100

Dividend

(partial) Remainder

```
loop len(divisor)+1 times:
    if divisor ≤ remainder:
        quotient bit = 1;
        remainder = remainder – divisor;
    else:
        quotient bit = 0;
    bring down next dividend bit to remainder;
```

# Division Algorithm in HW Manner

Quotient

0

Divisor 1000 01001010

0100

Dividend

(partial) Remainder

```
loop len(divisor)+1 times:
  if divisor ≤ remainder:
    quotient bit = 1;
    remainder = remainder – divisor;
  else:
    quotient bit = 0;
  bring down next dividend bit to remainder;
```

# Division Algorithm in HW Manner

Quotient

$$0$$

Divisor → $1000\overline{)01001010}$

$$0100$$

Dividend

$$01001$$

(partial) Remainder

```
loop len(divisor)+1 times:
  if divisor ≤ remainder:
    quotient bit = 1;
    remainder = remainder – divisor;
  else:
    quotient bit = 0;
  bring down next dividend bit to remainder;
```

# Division Algorithm in HW Manner

Quotient

Divisor

Dividend

(partial) Remainder

```
0
1000 ) 01001010
       0100
       01001
```

```
loop len(divisor)+1 times:
    if divisor ≤ remainder:
        quotient bit = 1;
        remainder = remainder – divisor;
    else:
        quotient bit = 0;
    bring down next dividend bit to remainder;
```

# Division Algorithm in HW Manner

Quotient

$$01$$

Divisor $1000 | 01001010$

0100

01001

Dividend

(partial) Remainder

```
loop len(divisor)+1 times:
  if divisor ≤ remainder:
    quotient bit = 1;
    remainder = remainder – divisor;
  else:
    quotient bit = 0;
  bring down next dividend bit to remainder;
```

# Division Algorithm in HW Manner

Quotient

Divisor

Dividend

(partial) Remainder

```
01
1000 | 01001010
       0100
       01001
     -  1000
           1
```

```
loop len(divisor)+1 times:
  if divisor ≤ remainder:
    quotient bit = 1;
    remainder = remainder - divisor;
  else:
    quotient bit = 0;
  bring down next dividend bit to remainder;
```

# Division Algorithm in HW Manner

Quotient

Divisor

Dividend

```
loop len(divisor)+1 times:
    if divisor ≤ remainder:
        quotient bit = 1;
        remainder = remainder – divisor;
    else:
        quotient bit = 0;
    bring down next dividend bit to remainder;
```

```
           01
1000 | 01001010
       0100
       01001
     -  1000
          10
```

(partial)
Remainder

# Division Algorithm in HW Manner

**Quotient**

**Divisor**

**Dividend**

**Remainder**

```
loop len(divisor)+1 times:
  if divisor ≤ remainder:
    quotient bit = 1;
    remainder = remainder – divisor;
  else:
    quotient bit = 0;
  bring down next dividend bit to remainder;
```

$$
\begin{array}{r}
01001 \\
1000 \overline{)01001010} \\
0100 \\
01001 \\
-\ 1000 \\
\hline
10 \\
101 \\
1010 \\
-\ 1000 \\
\hline
10
\end{array}
$$

# Division Hardware

In HW, two values are **subtracted** to compare their values
i.e., remainder = remainder – divisor

Quotient

Divisor

Dividend

If remainder ≥ 0, the result of the subtraction is used **directly as the remainder**

If remainder < 0,
**restore the remainder** by adding the divisor

Remainder

```
loop len(divisor)+1 times:
  if divisor ≤ remainder:
    quotient bit = 1;
    remainder = remainder – divisor;
  else:
    quotient bit = 0;
  bri
```

```
      01001
1000 |01001010
      0100
      01001
       10
      101
     1010
   - 1000
       10
```

# Division Hardware

Start

1. Subtract the Divisor register from the Remainder register and place the result in the Remainder register

Test Remainder

Remainder ≥ 0          Remainder < 0

2a. Shift the Quotient register to the left, setting the new rightmost bit to 1

2b. Restore the original value by adding the Divisor register to the Remainder register and placing the sum in the Remainder register. Also shift the Quotient register to the left, setting the new least significant bit to 0

3. Shift the Divisor register right 1 bit

33rd repetition?

No: < 33 repetitions

Yes: 33 repetitions

Done

Divisor
Shift right
64 bits

64-bit ALU

Remainder
Write
64 bits

Quotient
Shift left
32 bits

Control test

# Division Hardware: Example

```
                  Quotient
                  01001
Dividend
1000 | 01001010
            0100
            01001
          -  1000
               10
              101
             1010
           - 1000
               10
Remainder
```

# Example: Initial State

Initially divisor in left half

Quotient

Dividend

Divisor

```
             01001
   1000|01001010
             0100
             01001
         -    1000
                10
               101
              1010
            - 1000
                10
```

Remainder

**1000**0000



Divisor
Shift right

8 bits

8-bit ALU

Remainder
**01001010** Write

8 bits

Control test

0000

Quotient
Shift left

4 bits

# Example: Initial State

Initially divisor in left half

Quotient

Dividend

```
         01001
1000|01001010
      0100
      01001
    -  1000
         10
        101
       1010
      - 1000
         10
```

Divisor

Remainder

1000 0000

Divisor
Shift right

8 bits

8-bit ALU

Remainder
01001010  Write

8 bits

Control test

Start with 0

0000

Quotient
Shift left

4 bits

Start with the dividend

# Example: 1st Iteration - 1

remainder = remainder – divisor;

Quotient

Dividend

```
        01001
  1000 01001010
       0100
       01001
     -  1000
         10
         101
         1010
       - 1000
           10
```

Divisor

Remainder

1000**0000**

Divisor
Shift right

8 bits

01001010    10000000

sub

8-bit ALU

0000

Quotient
Shift left

4 bits

Remainder
**01001010** Write

Control
test

8 bits

# Example: 1st Iteration - 1

```
remainder = remainder – divisor;
```

Quotient

Dividend

```
      01001
1000 │01001010
      0100
      01001
    -  1000
         10
        101
       1010
      - 1000
         10
```

Divisor

Remainder

**10000000**



Divisor
Shift right

8 bits

**01001010    10000000**

**sub**

8-bit ALU

**11001010**

Remainder
**11001010** Write

**1**

Control test

**0000**

Quotient
Shift left

4 bits

8 bits

**Quotient**

**Dividend**

```
      01001
    01001010
    0100
      01001
   - 1000
        10
       101
      1010
    - 1000
        10
```

**Divisor**

**Remainder**

1000 0000

Remainder ≥ 0        **Test Remainder**        Remainder < 0

`if divisor ≤ remainder:`

Divisor
Shift right

8 bits

8-bit

Remainder
**1**1001010 Write

8 bits

Remainder < 0
(divisor > remainder)

0000

Quotient
Shift left

4 bits

Control
test

**1 (negative)**

# Example: 1st Iteration – 2b

**Quotient**

**Dividend**

```
        01001
1000 | 01001010
        0100
        01001
      -  1000
            10
           101
          1010
        - 1000
            10
```

**Divisor**

**Remainder**

1000**0000**

Divisor
Shift right

8 bits

**11001010    10000000**

**add**

8-bit ALU

**01001010**

Remainder
**01001010** Write

**1**

8 bits

Control test

**0000**

Quotient
Shift left

4 bits

Restore the original remainder

# Example: 1st Iteration – 2b

Quotient

Dividend

01001

1000 | 01001010

Divisor

0100

01001

- 1000

10

101

1010

- 1000

10

Remainder

10000000

quotient bit = 0;

Divisor
Shift right

8 bits

11001010    10000000

0000

add

8-bit ALU

Quotient
Shift left

4 bits

01001010

Remainder
01001010 Write

1

Control
test

shift with 0

8 bits

Restore the original remainder

# Example: 1st Iteration – 3

bring down next dividend bit to remainder
(Not equivalent in meaning,
but it has a similar meaning)

Quotient

Dividend

```
       01001
1000│01001010
      0100
      01001
    -  1000
         10
        101
       1010
      - 1000
         10
```

Divisor

Remainder

01000000

Divisor
Shift right

shift

8 bits

8-bit ALU

0000

Quotient
Shift left

4 bits

Remainder
01001010 Write

Control
test

8 bits

*In this example, we should check if it is the 5th repetition*

`loop len(divisor)+1 times`

33rd repetition?

No: < 33 repetitions

Yes: 33 repetitions

Done

# Example: 2nd Iteration – 1

remainder = remainder – divisor;

**Quotient**

**Dividend**

01001

1000 | 01001010

**Divisor**

0100

01001

− 1000

10

101

1010

- 1000

**Remainder** 10

01000000

Divisor
Shift right

8 bits

01001010   01000000

sub

8-bit ALU

00001010

Remainder
00001010 Write

8 bits

1

Control
test

0000

Quotient
Shift left

4 bits

Quotient

Dividend

Divisor

```
       01001
  1000 01001010
       0100
       01001
     -  1000
          10
         101
        1010
      -  1000
           10
```

Remainder

01000000

Remainder ≥ 0 | Test Remainder | Remainder < 0

if divisor ≤ remainder:

Divisor — Shift right

8 bits

0000

8-bit

Quotient Shift left

4 bits

Remainder
00001010 Write

Control test

8 bits

Remainder ≥ 0
(divisor < remainder)

0 (positive)

# Example: 2nd Iteration – 2a

quotient bit = 1;

Quotient

Dividend

0**1**001

1000 | 01001010

Divisor

0100

01001

- 1000

10

101

1010

- 1000

Remainder 10

0**1000**000

Divisor
Shift right

8 bits

8-bit ALU

Remainder
**00001010** Write

8 bits

000**1**

Quotient
Shift left

4 bits

Control
test

shift with 1

# Example: 2nd Iteration – 3

bring down next dividend bit to remainder
(Not equivalent in meaning, but it has a similar meaning)

00**1000**00

**Quotient**

**Dividend**

**Divisor**

```
          01001
1000 | 01001010
       0100
       01001
     -  1000
          10
         101
        1010
      - 1000
          10
```

**Remainder**

shift



Divisor
Shift right

8 bits

8-bit ALU

00**01**

Quotient
Shift left

4 bits

Remainder
**00001010** Write

8 bits

Control test

# Division Hardware: Algorithm Summary

Start

1. Subtract the Divisor register from the Remainder register and place the result in the Remainder register

Test Remainder

Remainder ≥ 0

Remainder < 0

2a. Shift the Quotient register to the left, setting the new rightmost bit to 1

2b. Restore the original value by adding the Divisor register to the Remainder register and placing the sum in the Remainder register. Also shift the Quotient register to the left, setting the new least significant bit to 0

3. Shift the Divisor register right 1 bit

33rd repetition?

No: < 33 repetitions

Yes: 33 repetitions

Done

Divisor

Shift right

64 bits

64-bit ALU

Remainder

Write

64 bits

Quotient

Shift left

32 bits

Control test

*Can we optimize this division hardware?*

# Optimized Divider

# What has Changed?

- 64 => 32 bits
- No shift

64 => 32-bit ALU

Divisor

32 bits

32-bit ALU

Divisor

Shift right

64 bits

64-bit ALU

Quotient

Shift left

32 bits

Remainder    Quotient

Shift right
Shift left
Write

Remainder

64 bits

Control test

Remainder

Write

64 bits

Remainder register:
- Left-half is remainder
- Right-half is quotient

Shift left/right

# Optimized Divider – Performance Gained

- Looks a lot like an optimized multiplier!
- The hardware is optimized to halve the width of the ALU and registers (64 bits ⇒ 32 bits, Clock cycle time ↓)
- Perform steps in parallel: add/shift (# of clock cycle ↓)

# Optimized Divider: Algorithm



Divisor

32 bits

32-bit ALU

Remainder

64 bits

Shift right
Shift left
Write

Control test

Start

1. Shift the Remainder register left 1 bit

2. Subtract the Divisor register from the left half of the Remainder register and place the result in the left half of the Remainder register

Test Remainder

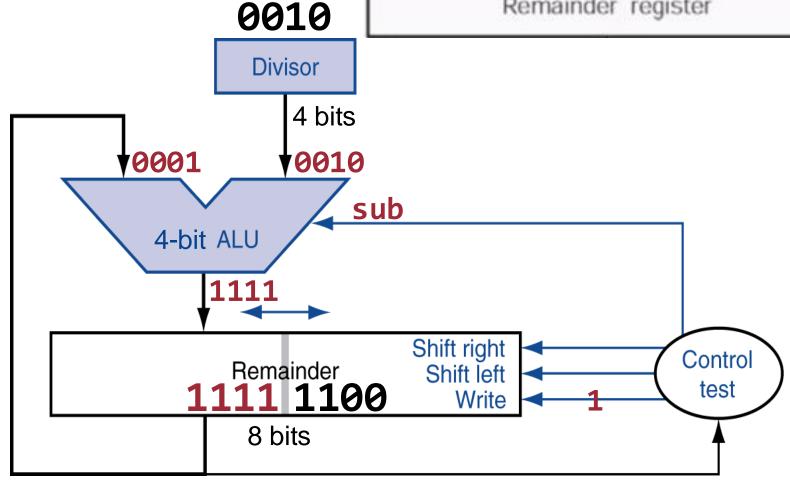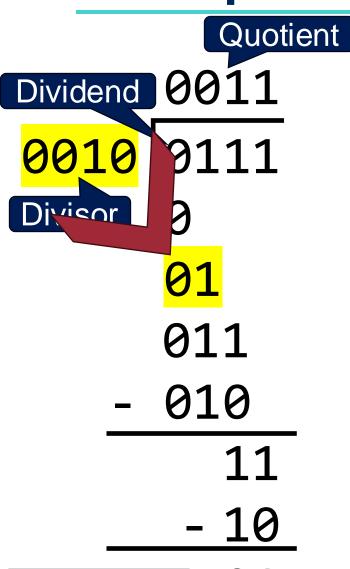Remainder ≥ 0          Remainder < 0

3a. Shift the Remainder register to the left, setting the new rightmost bit to 1

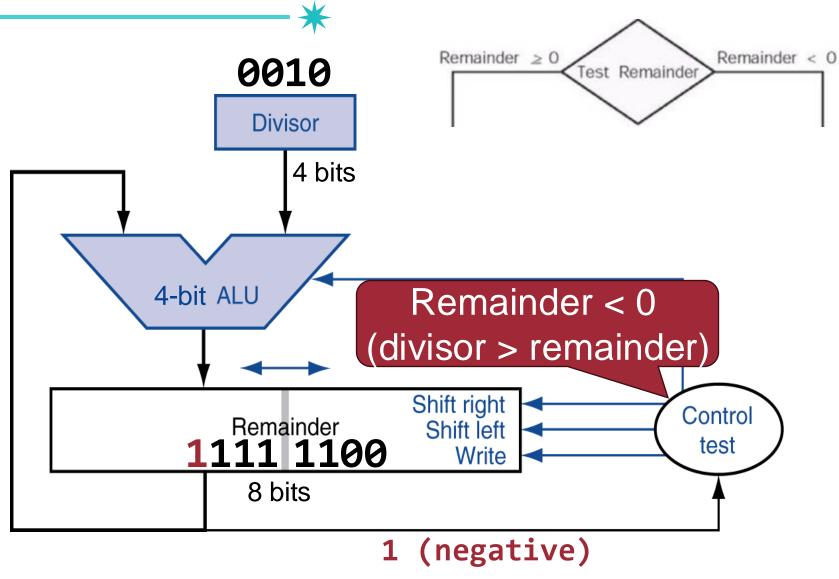3b. Restore the original value by adding the Divisor register to the left half of the Remainder register and place the sum in the left half of the Remainder register. Also shift the Remainder register to the left, setting the new rightmost bit to 0

32nd repetition?          No:  < 32 repetitions

Yes:  32 repetitions

Done. Shift left half of Remainder right 1 bit

# Optimized Divider: Example

Quotient

Dividend

```
        0011
0010 | 0111
        0
        01
        011
     -  010
     _____
        11
      - 10
     _____
        01
```

Divisor

Remainder



Divisor

4 bits

4-bit ALU

Remainder | Shift right / Shift left / Write

8 bits

Control test

# Example: Initial State

Quotient

Dividend

```
      0011
0010 |0111
       0
       01
       011
     - 010
     ───────
        11
      - 10
     ───────
```

Divisor

Remainder    01

**0010**

Divisor

4 bits

4-bit ALU

Remainder
**0000 0111**

Shift right
Shift left
Write

8 bits

Control
test

# Example: 1

**Quotient**

**Dividend** 0011

0010 | 0111

**Divisor**

0

01

011

- 010

11

- 10

**Remainder** 01

0010

Divisor

4 bits

4-bit ALU

Remainder
0000 1110

Shift right
Shift left  **shift with 0**
Write

Control
test

8 bits

# Example: 1st Iteration – 2

**Quotient**

**Dividend** 0011

0010 | 0111

**Divisor** 0

01

011

- 010

11

- 10

**Remainder** 01

0010

**Divisor**

4 bits

0000    0010

sub

**4-bit ALU**

1110

Remainder
**1110** 1110

Shift right
Shift left
Write    1

8 bits

Control test

Quotient

Dividend
0011

0010  0111

Divisor  0

01

011

- 010

11

- 10

Remainder  01

0010

Divisor

4 bits

Remainder ≥ 0    Test Remainder    Remainder < 0

4-bit ALU

Remainder < 0
(divisor > remainder)

Shift right
Remainder    Shift left
**1110 1110**    Write

Control
test

8 bits

1 (negative)

# Example: 1st Iteration – 3b

**Quotient**

**0010**

Divisor

**Dividend**

0011

0010 | 0111

**Divisor**

0

01

011

- 010

11

- 10

**Remainder** 01

4 bits

**1110**  **0010**

**add**

4-bit ALU

**0000**

Remainder
**0000 1110**   Shift right / Shift left / Write

**1**

Control test

8 bits

# Example: 1st Iteration – 3b

Quotient

Dividend

Divisor

Remainder

```
     0011
0010)0111
      0

      01
     011
    - 010
    ─────
      11
    - 10
    ─────
      01
```

0010
Divisor

4 bits

4-bit ALU

Remainder
0001 1100

Shift right
Shift left
Write

**shift with 0**

Control test

8 bits

*In this example, we should check if it is the 4th repetition*

Quotient

Dividend 0011

0010

Divisor

0010

0

**32nd repetition?** → No: < 32 repetitions

Yes: 32 repetitions

0001 1100

8 bits

Write

11

Control test

10

Remainder 01

# Example: 2nd Iteration – 2

Quotient

**0011**

Dividend

**0010 | 0111**

Divisor

0

01

011

- 010

11

- 10

Remainder 01

**0010**

Divisor

4 bits

**0001**   **0010**

sub

4-bit ALU

**1111**

Remainder
**1111 1100**

Shift right
Shift left
Write   **1**

Control test

8 bits

**Quotient**

**Dividend**

**Divisor**

**Remainder**

```
      0011
0010 )0111
       0
       01
       011
     - 010
       11
     - 10
       01
```

0010

Divisor

4 bits

4-bit ALU

Remainder ≥ 0    Test Remainder    Remainder < 0

**Remainder < 0
(divisor > remainder)**

Remainder
**1**111 1100

Shift right
Shift left
Write

Control
test

8 bits

**1 (negative)**

# Example: 2nd Iteration – 3b

Quotient

Dividend

Divisor

Remainder

$$
\begin{array}{r}
0011 \\
0010 \overline{\smash{\big)}\,0111} \\
0 \\
01 \\
011 \\
-\ 010 \\
\hline
11 \\
-\ 10 \\
\hline
01 \\
\end{array}
$$

0010

Divisor

4 bits

1111    0010

add

4-bit ALU

0001

Remainder
0001 1100

Shift right
Shift left
Write

1

Control
test

8 bits

# Example: 2nd Iteration – 3b

Quotient

Dividend

Divisor

Remainder

```
      0011
     ------
0010 |0111
      0
      --
      01
      011
    - 010
    ------
      11
    - 10
    ------
      01
```

0010

Divisor

4 bits

4-bit ALU

Remainder

0011 1000

8 bits

Shift right
Shift left  **shift with 0**
Write

Control test

# Example: 3rd Iteration – 2

Quotient

Dividend  0011

Divisor  0010 | 0111

0

01

011
- 010

11

- 10

Remainder  01

0010

Divisor

4 bits

0011    0010

sub

4-bit ALU

0001

Remainder
0001 1000

Shift right
Shift left
Write

1

Control test

8 bits

Quotient

Dividend

Divisor

Remainder

```
       0011
      _____
0010 | 0111
       0
       __
       01
       011
      -010
      _____
       11
      - 10
      _____
       01
```

0010

Divisor

4 bits

4-bit ALU

Remainder ≥ 0
(divisor < remainder)

Remainder ≥ 0   Test Remainder   Remainder < 0

Remainder
0001 1000

Shift right
Shift left
Write

Control
test

8 bits

0 (positive)

# Example: 3rd Iteration – 3a

Quotient

**0011**

Dividend

0010 | 0111

Divisor

0
01
011
-  010
_____
**11**
-  10
_____

Remainder    01

0010

Divisor

4 bits

4-bit ALU

8 bits

Remainder

**0011 0001**

Shift right
Shift left
Write

**shift with 1**

Control test

# Example: 4th Iteration – 2

Quotient

0011

Dividend

0010 | 0111

Divisor

0

01

011

- 010

11

- 10

Remainder 01

**0010**

Divisor

4 bits

**0011**   **0010**

**sub**

4-bit ALU

**0001**

Remainder
**0001** **0001**

Shift right
Shift left
Write

**1**

Control test

8 bits

Quotient

Dividend

Divisor

0011

0010 | 0111

0

011

- 010

**11**

- 10

Remainder 01

0010
Divisor

4 bits

4-bit ALU

Remainder ≥ 0
(divisor < remainder)

Remainder
**0001 0001**

8 bits

Shift right
Shift left
Write

Control
test

0 (positive)

Remainder ≥ 0    Test Remainder    Remainder < 0

# Example: 4th Iteration – 3a

Quotient

Dividend

**0011**

0010 | 0111

Divisor

0

01

011

- 010

11

- 10

Remainder  01

0010

Divisor

4 bits

4-bit ALU

Remainder

**0010 0011**

8 bits

Shift right
Shift left
Write

shift with 1

Control test

# Example: Final Step

**Quotient**

**Dividend**

```
      0011
0010 |0111
      0
      01
      011
    - 010
      ─────
      11
    - 10
      ────
      01
```

**Divisor**

**Remainder**

0010

Divisor

4 bits

4-bit ALU

Shift left-half of remainder

Remainder

0001 0011

8 bits

Shift right
Shift left
Write

**shift**

Control test

# Division of Signed Number

- The *magnitude* of quotient/remainder depends on the *magnitude* of dividend/divisor

- The *sign* of quotient/remainder depends on the *sign* of dividend/divisor

| Example | Quotient | Remainder |
|---------|----------|-----------|
| +7 / +2 | +3 | +1 |
| -7 / +2 | -3 | -1 |
| +7 / -2 | -3 | +1 |
| -7 / -2 | +3 | -1 |

# Division of Signed Number

- The **magnitude** of <u>quotient/remainder</u> depends on the **magnitude** of <u>dividend/divisor</u>

- The **sign** of quotient/remainder depends on

**After calculating with positive values, convert based on the sign at the end**

| +7  /  +2 | +3 | +1 |
|-----------|----|----|
| -7  /  +2 | -3 | -1 |
| +7  /  -2 | -3 | +1 |
| -7  /  -2 | +3 | -1 |

# Faster Division?

- Can't use parallel hardware as in multiplier
  - Subtraction is conditional on sign of remainder

# MIPS Division

- **Use `HI`/`LO` registers for result**
  - HI: 32-bit remainder
  - LO: 32-bit quotient

- **Instructions**
  - `div rs, rt` / `divu rs, rt`
  - No overflow or divide-by-0 checking
    - Software must perform checks if required
  - Use `mfhi`, `mflo` to access result

# Floating-point Number Arithmetic

# Background: Number System

- Decimal number system

$$\dots \quad 10^2 \quad 10^1 \quad 10^0 \quad 10^{-1} \quad 10^{-2} \quad \dots$$

x10     x10     /10     /10

- Binary number system

$$\dots \quad 2^2 \quad 2^1 \quad 2^0 \quad 2^{-1} \quad 2^{-2} \quad \dots$$

x2     x2     /2     /2

# Floating-point Number: Motivation

**We need a way to represent ...**

- *Infinite decimal (e.g., 3.1415926535…)*

- *Very small numbers*

- *Very large numbers*

***from computer!***

**Solution: Floating-point Number Representation**

# Floating-point Number: Motivation

## We need a way to represent ...

- *Infinite decimal (e.g., 3.1415926535…)* → *Approximate value* → $3.1415$

- *Very small numbers* → *Floating decimal point* → $0.001 \times 10^{-20}$

- *Very large numbers* → *Floating decimal point* → $3.15576 \times 10^{19}$

**Can be represented with a limited number of bits!**

**Solution: Floating-point Number Representation**

# Floating-point Number: Notations

- **Scientific notation**: renders numbers with a <mark>single digit</mark> to the *left* of the point
  - Example: $7.15576 \times 10^4$, $0.314 \times 10^1$

*Normalized*

- **Normalized scientific notation**: scientific notation that has *no leading 0s*
  - Example: $7.15576 \times 10^4$, $3.14 \times 10^0$

# Floating-point Number: Notations

- **Scientific notation**: renders numbers with a <u>single digit</u> to the *left* of the point

$$1 \leq \textit{Significand} < 10$$

- **Normalized scientific notation**: <u>scientific notation</u> that has *no leading 0s*
  - Example: 7.15576 x $10^4$, 3.14 x $10^0$

Significand     Base     Exponent

$$(sign) \times significand \times base^{exponent}$$

# Floating-point Number: Binary

- **Scientific notation**: renders numbers with a <u>single digit</u> to the *left* of the point

  **_1 ≤ Significand < 2_**

- **Normalized scientific notation**: <u>scientific notation</u> that has *no leading 0s (= Always 1 to the left of the point)*
  - Example: **$1.1 \times 2^{-1}$**  ($0.75_{ten} \to 0.5 + 0.25$
    $\to 0.11_{two} \to 1.1 \times 2^{-1}$)

**2**

$$\text{(sign) x significand x base}^{\text{exponent}}$$

# Floating-point Number: **Binary**

- **Scientific notation**: renders numbers with a single digit to the *left* of the point

$$1 \leq Significand < 2$$

## How are floating-point numbers represented in computers?

$\rightarrow 0.11_{two} \rightarrow 1.1 \times 2^{-1})$

$$\text{(sign)} \times \text{significand} \times \text{base}^{exponent}$$

2

# IEEE 754 Floating-point Standard

- Developed in response to *divergence of representations*

**Divergence of representations**

$$0.11_{two} = 1.1 \times 2^{-1}{}_{two} = 11 \times 2^{-2}{}_{two}$$

**Normalized representation**

$$1.1 \times 2^{-1}{}_{two}$$

**IEEE 756 representation**

| S | Exponent | Fraction |
|---|----------|----------|

# IEEE 754 Floating-point Standard

- Two representations
  - **Single precision (32-bit)**: type `float` in C
  - **Double precision (64-bit)**: type `double` in C

**Sign**
- Single: 1 bit
- Double: 1 bit

**Exponent**
- Single: 8 bits
- Double: 11 bits

**Fractional part of significand**
- Single: 23 bits
- Double: 52 bits

| S | Exponent | Fraction |
|---|----------|----------|

$$+1.1 \times 2^{-1}{}_{two}$$

# IEEE 754 Floating-point Standard

- Two representations
  - **Single precision (32-bit)**: type `float` in C
  - **Double precision (64-bit)**: type `double` in C

**Fractional part of significand**
- Single: 23 bits
- Double: 52 bits

| S | Exponent | Fraction |
|---|----------|----------|

$$+1.1 \times 2^{-1}{}_{two}$$

Leading 1 bit is implicit

# IEEE 754 Floating-point Standard

- Two representations
  - **Single precision (32-bit)**: type `float` in C
  - **Double precision (64-bit)**: type `double` in C

**Exponent**
- Single: 8 bits
- Double: 11 bits

**Fractional part of significand**
- Single: 23 bits
- Double: 52 bits

| S | Exponent | Fraction |
|---|----------|----------|

**-1+bias**

$+1.1 \times 2^{-1}{}_{two}$

Leading 1 bit is implicit

# Exponent: Why Biased?

- To make sorting easier
- **Bias**: *127 for single precision* and *1023 for double precision*

| S | Exponent | Fraction |
|---|----------|----------|

$$-1 + bias$$

$$+1.1 \times 2^{-1}{}_{two}$$

# Exponent: Why Biased?

- To make sorting easier
- **Bias**: *127 for single precision* and *1023 for double precision*

| Exponent (8 bits) | Decimal (Unsigned) | Biased by 127 |
|---|---|---|
| 00000000 | 0 | 0-127 = -127 |
| 00000001 | 1 | 1-127 = -126 |
| 00000010 | 2 | 2-127 = -125 |
| ... | ... | ... |
| 11111111 | 255 | 255-127 = 128 |

| S | Exponent | |
|---|---|---|

**-1+bias**

$+1.1 \times 2^{-1}{}_{two}$

# Exponent: Why Biased?

- To make sorting easier
- **Bias**: *127 for single precision* and *1023 for double precision*

| Exponent (8 bits) | Decimal (Unsigned) | Biased by 127 |
|---|---|---|
| 00000000 | 0 | 0-127 = -127 |
| 00000001 | 1 | 1-127 = -126 |
| 00000010 | 2 | 2-127 = -125 |
| ... | ... | ... |
| 11111111 | 255 | 255-127 = 128 |

| S | Exponent |
|---|---|

**-1+bias**

$+1.1 \times 2^{-1}$ two

Easy sorting

# Tradeoff between Precision and Range

- Increasing *the size of the* fraction enhances ***the precision***
  - Shorter length: $1.1 \times 2^{-1}$
  - Longer length: $1.10110 \times 2^{-1}$

- Increasing *the size of the* exponent increases ***the range***
  - Shorter length: $1.1 \times 2^{2}$
  - Longer length : $1.1 \times 2^{23}$

*Good design demands good compromise!*

| S | Exponent | Fraction |
|---|----------|----------|

# IEEE 754 Floating-point Standard

- Two representations
  - **Single precision (32-bit)**: type `float` in C
  - **Double precision (64-bit)**: type `double` in C

| **Exponent** <br> • Single: 8 bits <br> • Double: 11 bits | **Fractional part of significand** <br> • Single: 23 bits <br> • Double: 52 bits |

| S | Exponent | Fraction |

**-1+bias**

$$+ 1.1 \times 2^{-1}{}_{two}$$

Leading 1 bit is implicit

# IEEE 754 Floating-point Standard

- Two representations
  - **Single precision (32-bit)**: type `float` in C
  - **Double precision (64-bit)**: type `double` in C

**Sign**
- Single: 1 bit
- Double: 1 bit

**Exponent**
- Single: 8 bits
- Double: 11 bits

**Fractional part of significand**
- Single: 23 bits
- Double: 52 bits

| S | Exponent | Fraction |

**-1+bias**

- 0: positive (+)
- 1: negative (-)

$$+1.1 \times 2^{-1}{}_{two}$$

Leading 1 bit is implicit

# IEEE 754 Floating-point Standard: Summary

$$(-1)^{sign} \times (1+fraction) \times base^{exponent-bias}$$

| S | Exponent | Fraction |
|---|----------|----------|

# Single-precision Range

- Exponents 00000000 and 11111111 reserved

- **Smallest value**
  - −Exponent: 00000001
    $\Rightarrow$ actual exponent = 1 − 127 = −126
  - −Fraction: 000…00 $\Rightarrow$ significand = 1.0
  - −$\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$

- **Largest value**
  - −Exponent: 11111110
    $\Rightarrow$ actual exponent = 254 − 127 = +127
  - −Fraction: 111…11 $\Rightarrow$ significand $\approx$ 2.0
  - −$\pm 2.0 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$

Used for special cases

# Double-precision Range

- Exponents 0000…00 and 1111…11 reserved

Used for special cases

- **Smallest value**
  - Exponent: 00000000001
    $\Rightarrow$ actual exponent = $1 - 1023 = -1022$
  - Fraction: 000…00 $\Rightarrow$ significand = 1.0
  - $\pm 1.0 \times 2^{-1022} \approx \pm 2.2 \times 10^{-308}$

- **Largest value**
  - Exponent: 11111111110
    $\Rightarrow$ actual exponent = $2046 - 1023 = +1023$
  - Fraction: 111…11 $\Rightarrow$ significand $\approx 2.0$
  - $\pm 2.0 \times 2^{+1023} \approx \pm 1.8 \times 10^{+308}$

# Overflow and Underflow

The representable range using normalized notation

$-\infty \longleftarrow$ _____ $\longrightarrow +\infty$

$-2.0 \times 2^{127}$      $-1.0 \times 2^{-126}$   $+1.0 \times 2^{-126}$      $+2.0 \times 2^{127}$

0

# Overflow and Underflow

- **Overflow**: occurs when a result has a magnitude **too big** to be represented

The representable range using normalized notation

Floating-point overflow

$-\infty \leftarrow$ ——— $\rightarrow +\infty$

$-2.0 \times 2^{127}$  $-1.0 \times 2^{-126}$  $0$  $+1.0 \times 2^{-126}$  $+2.0 \times 2^{127}$

# Overflow and Underflow

- **Overflow**: occurs when a result has a magnitude *too big* to be represented

- **Underflow**: occurs when a result has a magnitude *too small* to be represented



The representable range using normalized notation

Floating-point underflow

Floating-point overflow

$-\infty$ ← $\quad$ → $+\infty$

$-2.0 \times 2^{127}$ $\quad$ $-1.0 \times 2^{-126}$ $\quad$ $+1.0 \times 2^{-126}$ $\quad$ $+2.0 \times 2^{127}$

0

# IEEE 754: Special Cases

The representable range using normalized notation

$-\infty$ $\longleftarrow$ $+\infty$

$-2.0 \times 2^{127}$      $-1.0 \times 2^{-126}$   $+1.0 \times 2^{-126}$      $+2.0 \times 2^{127}$

0

# IEEE 754: Special Cases

- Exponent = 00…0, Fraction = 00…0
  - → Not $1.0 \times 2^{-127}$ but **0**

The representable range using normalized notation

How to represent 0?

$-2.0 \times 2^{127}$   $-1.0 \times 2^{-126}$   $+1.0 \times 2^{-126}$   $+2.0 \times 2^{127}$

$-\infty$   $0$   $+\infty$

# IEEE 754: Special Cases

- Exponent = 00…0, Fraction = 00…0
  - $\rightarrow$ Not $1.0 \times 2^{-127}$ but **0**

- Exponent = 00…0, Fraction $\neq$ 00…0
  - $\rightarrow$ Not $(1 + \text{fraction}) \times 2^{-127}$ but **$(0 + \text{fraction}) \times 2^{-126}$**
  - $\rightarrow$ Denormalized real numbers (to represent very small numbers)

The representable range using normalized notation

How to represent very small numbers?
$\rightarrow$ *Denormalized numbers*

$-\infty \longleftarrow$ $\longrightarrow +\infty$

$-2.0 \times 2^{127}$ $\quad$ $-1.0 \times 2^{-126}$ $\quad$ $+1.0 \times 2^{-126}$ $\quad$ $+2.0 \times 2^{127}$

$0$

# IEEE 754: Special Cases

- Exponent = 00…0, Fraction = 00…0
  - $\rightarrow$ Not $1.0 \times 2^{-127}$ but **0**

- Exponent = 00…0, Fraction ≠ 00…0
  - $\rightarrow$ Not $(1 + \text{fraction}) \times 2^{-127}$ but **$(0 + \text{fraction}) \times 2^{-126}$**
  - $\rightarrow$ Denormalized real numbers (to represent very small numbers)

- Exponent = 11…1, Fraction = 00…0
  - $\rightarrow$ **±infinity**

How to represent infinity?



$-\infty$ $\leftarrow$ ... $\rightarrow$ $+\infty$

$-2.0 \times 2^{127}$   $-1.0 \times 2^{-126}$   $+1.0 \times 2^{-126}$   $+2.0 \times 2^{127}$

0

# IEEE 754: Special Cases

- Exponent = 00…0, Fraction = 00…0
  - → Not $1.0 \times 2^{-127}$ but **0**

- Exponent = 00…0, Fraction ≠ 00…0
  - → Not $(1 + \text{fraction}) \times 2^{-127}$ but **$(0 + \text{fraction}) \times 2^{-126}$**
  - → Denormalized real numbers (to represent very small numbers)

- Exponent = 11…1, Fraction = 00…0
  - → **±infinity**

- Exponent = 11…1, Fraction ≠ 00…0
  - → **Not-a-Number (*NaN*)**
  - → Indicates illegal or undefined result

How to represent the result of invalid operations (e.g., 0/0)?

# IEEE 754 Encoding of Floating-point Numbers

| 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) |

# Floating-point Addition: Decimal

- Consider a 4-digit decimal example

$9.999 \times 10^1$ + $1.610 \times 10^{-1}$

## 1. Align decimal points

- Shift number with *smaller exponent*
- $9.999 \times 10^1$ + $0.016 \times 10^1$

# Floating-point Addition: Decimal

- Consider a 4-digit decimal example

  $9.999 \times 10^1 + 1.610 \times 10^{-1}$

## 1. Align decimal points

- − Shift number with *smaller exponent*
- − $9.999 \times 10^1 + 0.016 \times 10^1$

## 2. Add significands

- − $10.015 \times 10^1$

# Floating-point Addition: Decimal

- Consider a 4-digit decimal example

    $9.999 \times 10^1 + 1.610 \times 10^{-1}$

## 1. Align decimal points

- Shift number with *smaller exponent*
- $9.999 \times 10^1 + 0.016 \times 10^1$

## 2. Add significands

- $10.015 \times 10^1$

## 3. Normalize result & check for over/underflow

- $1.0015 \times 10^2$

## 4. Round

- $1.002 \times 10^2$

# Floating-point Addition: Binary

- Now consider a 4-digit binary example

$1.000_2 \times 2^{-1} + -1.110_2 \times 2^{-2}$ (0.5 + −0.4375)

# Floating-point Addition: Binary

- Now consider a 4-digit binary example

  $1.000_2 \times 2^{-1}$ + –$1.110_2 \times 2^{-2}$ (0.5 + –0.4375)

**1. Align binary points**
  - Shift number with *smaller exponent*
  - $1.000_2 \times 2^{-1}$ + –$0.111_2 \times 2^{-1}$

# Floating-point Addition: Binary

- Now consider a 4-digit binary example

  $1.000_2 \times 2^{-1} + {-}1.110_2 \times 2^{-2}$ $(0.5 + {-}0.4375)$

## 1. Align binary points

- − Shift number with *smaller exponent*
- − $1.000_2 \times 2^{-1} + {-}0.111_2 \times 2^{-1}$

## 2. Add significands

- − $0.001_2 \times 2^{-1}$

# Floating-point Addition: Binary

- Now consider a 4-digit binary example

    $1.000_2 \times 2^{-1} + {-}1.110_2 \times 2^{-2}$ (0.5 + –0.4375)

## 1. Align binary points

- Shift number with *smaller exponent*
- $1.000_2 \times 2^{-1} + {-}0.111_2 \times 2^{-1}$

## 2. Add significands

- $0.001_2 \times 2^{-1}$

## 3. Normalize result & check for over/underflow

- $1.000_2 \times 2^{-4}$, with no over/underflow

## 4. Round

- $1.000_2 \times 2^{-4} = 0.0625$

> Check *-126 ≤ -4 ≤ +127*
> in case of a single precision

# Floating-point Adder Hardware



Skip the details

# Exercise

$$1.0110_2 \times 2^3 + 1.1000_2 \times 2^2$$

## 1. Align binary points
- $1.0110_2 \times 2^3 + 0.1100_2 \times 2^3$

## 2. Add significands
- $10.0010_2 \times 2^3$

## 3. Normalize result & check for over/underflow
- $1.0001_2 \times 2^4$, with no over/underflow

## 4. Round
- $1.0001_2 \times 2^4$

# Floating-point Multiplication

- Consider a 4-digit decimal example

$1.110 \times 10^{10} \times 9.200 \times 10^{-5}$

## 1. Add exponents

- New exponent = 10 + –5 = 5

## 2. Multiply significands

- $1.110 \times 9.200 = 10.212 \implies 10.212 \times 10^{5}$

## 3. Normalize result & check for over/underflow

- $1.0212 \times 10^{6}$

## 4. Round

- $1.021 \times 10^{6}$

## 5. Determine sign of result from signs of operands

- $+1.021 \times 10^{6}$

# Floating-point Multiplication

- Now consider a 4-digit binary example
  $1.000_2 \times 2^{-1} \times -1.110_2 \times 2^{-2}$ ($0.5 \times -0.4375$)

**1. Add exponents**
- Unbiased: $-1 + -2 = -3$
- Biased: $(-1 + 127) + (-2 + 127) - 127 = -3 + 127$

> For biased exponents, subtract bias from sum

**2. Multiply significands**
- $1.000_2 \times 1.110_2 = 1.1102 \Rightarrow 1.110_2 \times 2^{-3}$

**3. Normalize result & check for over/underflow**
- $1.110_2 \times 2^{-3}$ (no change) with no over/underflow

**4. Round**
- $1.110_2 \times 2^{-3}$ (no change)

**5. Determine sign: + sign × – sign $\Rightarrow$ – sign**
- $-1.110_2 \times 2^{-3} = -0.21875$

# Floating-point Instructions in MIPS

- Separate floating-point registers
    - 32 single-precision: `$f0, $f1, … $f31`
    - Paired for double-precision: `$f0/$f1, $f2/$f3, …`
- FP instructions operate only on FP registers

| F1 | F0 |
|---|---|
| … | … |
| F31 | F30 |

| Category | Instruction | Example | Meaning | Comments |
|---|---|---|---|---|
| Arithmetic | FP add single | `add.s    $f2,$f4,$f6` | `$f2 = $f4 + $f6` | FP add (single precision) |
| | FP subtract single | `sub.s    $f2,$f4,$f6` | `$f2 = $f4 - $f6` | FP sub (single precision) |
| | FP multiply single | `mul.s    $f2,$f4,$f6` | `$f2 = $f4 × $f6` | FP multiply (single precision) |
| | FP divide single | `div.s    $f2,$f4,$f6` | `$f2 = $f4 / $f6` | FP divide (single precision) |
| | FP add double | `add.d    $f2,$f4,$f6` | `$f2 = $f4 + $f6` | FP add (double precision) |
| | FP subtract double | `sub.d    $f2,$f4,$f6` | `$f2 = $f4 - $f6` | FP sub (double precision) |
| | FP multiply double | `mul.d    $f2,$f4,$f6` | `$f2 = $f4 × $f6` | FP multiply (double precision) |
| | FP divide double | `div.d    $f2,$f4,$f6` | `$f2 = $f4 / $f6` | FP divide (double precision) |
| Data transfer | load word copr. 1 | `lwc1     $f1,100($s2)` | `$f1 = Memory[$s2 + 100]` | 32-bit data to FP register |
| | store word copr. 1 | `swc1     $f1,100($s2)` | `Memory[$s2 + 100] = $f1` | 32-bit data to memory |
| Condi-tional branch | branch on FP true | `bc1t     25` | if (cond == 1) go to PC + 4 + 100 | PC-relative branch if FP cond. |
| | branch on FP false | `bc1f     25` | if (cond == 0) go to PC + 4 + 100 | PC-relative branch if not cond. |
| | FP compare single (eq,ne,lt,le,gt,ge) | `c.lt.s $f2,$f4` | if ($f2 < $f4)     cond = 1; else cond = 0 | FP compare less than single precision |
| | FP compare double (eq,ne,lt,le,gt,ge) | `c.lt.d $f2,$f4` | if ($f2 < $f4)     cond = 1; else cond = 0 | FP compare less than double precision |

# Question?