

Computer Architecture

CH3 Computer Arithmetic (III) Floating Point

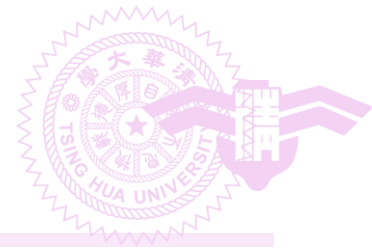
Prof. Ren-Shuo Liu
NTHU EE
Fall 2017





Outline

- Overview
- IEEE 754 standard
 - Single-precision
 - Double-precision
 - Special numbers
- Floating-point operations
 - Addition
 - Multiplication
 - Rounding



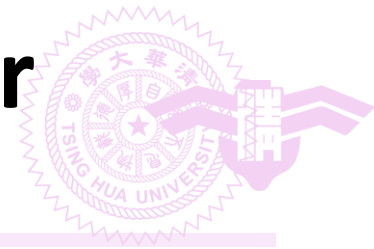
Outline

- Overview
 - RISC-V floating-point instructions
 - Fixed-point and floating-point representations
- IEEE 754 standard
- Floating-point operations



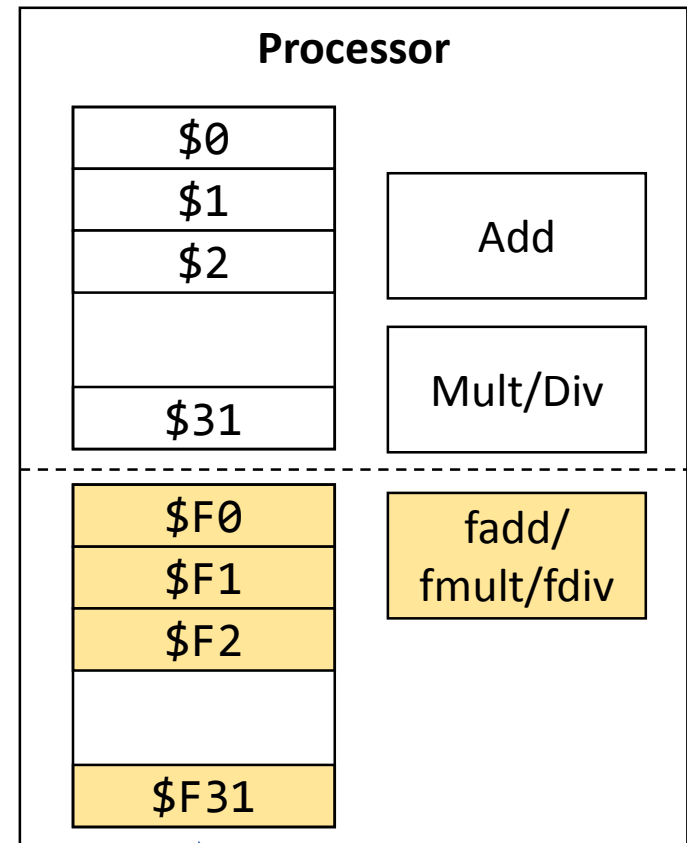
RISC-V Floating Instructions

- Arithmetic
 - fadd.s, fsub.s, fmul.s, fdiv.s # s means single-precision
 - fadd.d, fsub.d, fmul.d, fdiv.d # d means double-precision
- Comparisons
 - feq.s, feq.d # equal
 - flt.s, flt.d # less than
 - fle.s, fle.d # less than or equal
- Load / store
 - flw, fsw
 - fld, fsd



Floating Point Unit and Register Files

- Separate 32 registers for floating-point
 - Register pairs (e.g., \$F0 and \$F1) for double precision
 - \$F0 is not always zero
- Floating-point instructions can be optional
 - Many embedded systems do not utilize them



整數和浮點數無法直接做運算，但在寫 C code 的時候，Compiler 會幫你把整數轉成浮點數。



Fixed-Point

- Integers scaled by an **implicit** (隱含的) factor
- The **scaling factor** for each variable **does not change** (i.e., **fixed**) during the entire computation

- Examples

- 3.14 is represented as
 - 314 (scaling factor = $1/100$)
 - 3140 (scaling factor = $1/1000$)
- 5,000,000 is represented as
 - 5 (scaling factor = 1,000,000)
 - 50 (scaling factor = 100,000)

把小數表示成一個
整數跟一個scaling
factor

計算時小數點位置不變，
隱含著計算的每個數值都
有相同的scaling factor



Floating-Point ~ = Scientific Notation (科學記號表示法)

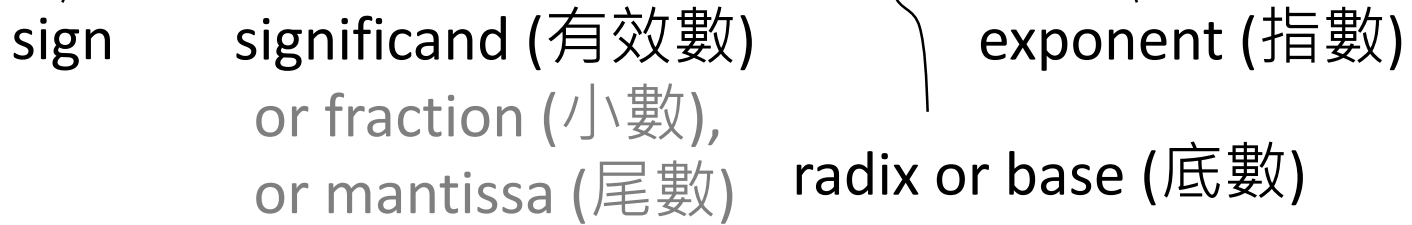
- 光速 $+ 2.99792458 \times 10^8$ (m/s)
 - 電子電量 $- 1.60217733 \times 10^{-19}$ (C)
 - 0.5莫耳碳原子 $+ 6.00000000 \times 10^{-3}$ (kg)
- sign significand (有效數)
 or fraction (小數),
 or mantissa (尾數) radix or base (底數)
- exponent (指數)

- **Normalized form**
 - Exactly one non-zero significant digit to the left of the point
 - 29.9×10^7 and 0.299×10^9 are not normalized forms



Floating-Point ~ = Scientific Notation (科學記號表示法)

- 光速 $+ 2.99792458 \times 10^8$ (m/s)
- 電子電量 $- 1.60217733 \times 10^{-19}$ (C)
- 0.5莫耳碳原子 $+ 6.00000000 \times 10^{-3}$ (kg)



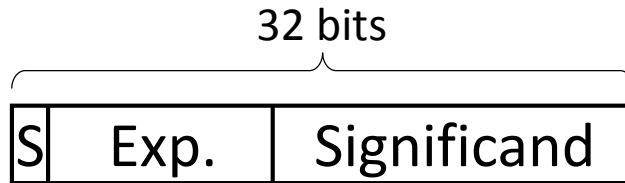
- **Scaling factor** (the exponent) is **explicit**
- "Floating"
 - Scaling factor can change during computation



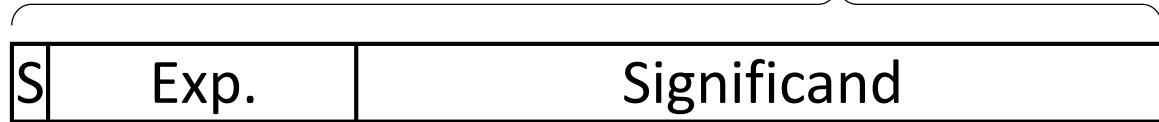
Floating Point Number

- IEEE 754 standard

- Single-precision



- Double-precision



只存小數部分，因為Normalized form規定只有一個數能在小數點左邊，且此為二進位(故直接寫1)

用Excess-K表示出小的數字

- Represented value: $(-1)^S \times \underline{1} \cdot \text{Significand}_{\text{two}} \times 2^{(\text{Exponent} - \text{Bias})}$

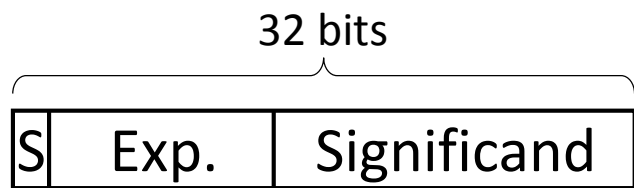
	Sign bit	Exp. bits	Significand bits	Bias
Single precision (<i>float</i> in C/C++)	1	0-255 8 unsigned	23	127
Double precision (<i>double</i> in C/C++)	1	11	52	1023

Single Precision: Exponent-Bias = -127->128



Special Floating Point Numbers

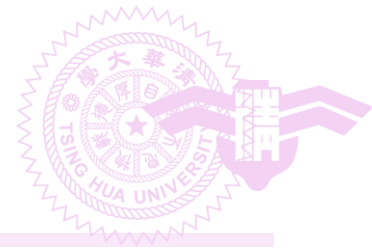
- 1. 32bit 的浮點數可以比32bit 的整數有更大的數值range (因為有Exponent)
- 2. 32bit 浮點數和32bit 整數一樣有 2^{32} 種組合。但是浮點數的0有重複(+0)且正負無限大和NaN都不是數值。所以32bit 整數能表示的數值個數比浮點數多



	S	Exp.	Significand
zero	S	0...00	0...00
Denormalized value	S	0...00	non-zero
$+\infty$	0	1...11	0...00
$-\infty$	1	1...11	0...00
Not a number (NaN)	d	1...11	non-zero

用於表示很小的數字 (0.000000X)

i.e., the maximal and minimal exponent values are reserved for special floating-point numbers



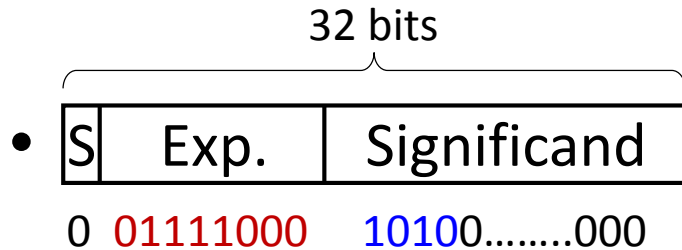
Denormalized Value

- | | | |
|---|--------|-------------|
| S | 0...00 | Significant |
|---|--------|-------------|

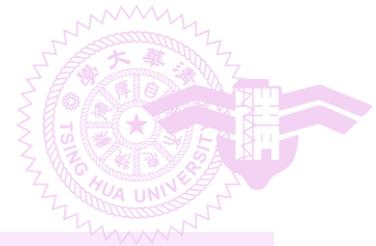
 are denormalized values
 - $(-1)^S \times 0.\text{Significant} \times 2^{(1-\text{bias})}$
 - No leading one to the left of the point
- Objective
 - Represent very small value
 - Gradual underflow



Floating Point Examples



$$\begin{aligned} &= 1.1010\dots\dots000_{\text{two}} \times 2^{(120-127)} \\ &= 1.1010_{\text{two}} \times 2^{-7} \\ &= 1.625_{\text{ten}} \times 2^{-7} \\ &= 0.0126953125_{\text{ten}} \end{aligned}$$



Floating Point Examples

- Convert -3.14 to 32-bit float

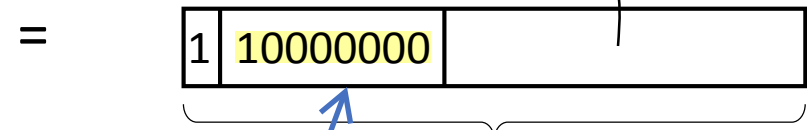
- $3 = 11_{\text{two}}$

- 3.14

$$= 11.0010_0011_1101_0111_0000_1010\dots_{\text{two}}$$

$$= 1.1001_0001_1110_1011_1000_010 \times 2^1$$

23-bit significand (assume not rounded)



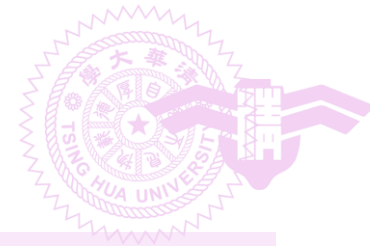
要考慮Bias(128)

32 bits

小數部分乘以2，取整數部分當作輸出Binary

fraction part x2

0.14	0.28	0.88
	0.56	1.76
	1.12	1.52
	0.24	1.04
	0.48	0.08
	0.96	0.16
	1.92	0.32
	1.84	0.64
	1.68	1.28
	1.36	0.56
	0.72	1.12
	1.44	0.24



IEEE 754 Online Converter

IEEE 754 CONVERTER

This page allows you to convert between the decimal representation of numbers (like "1.02") and the binary format used by all modern CPUs (IEEE 754 floating point). The conversion is limited to single precision numbers (32 Bit). The purpose of this webpage is to help you understand floating point numbers.

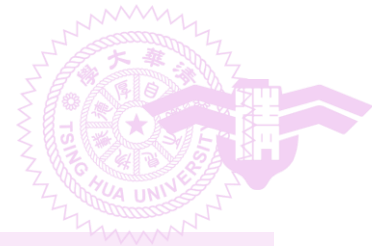
IEEE 754 Converter (JavaScript), V0.13

Note: This JavaScript-based version is still under development, please report errors [here](#).

	Sign	Exponent	Mantissa
Value:	+1	2^{-7}	1.625
Encoded as:	0	120	5242880
Binary:	<input type="checkbox"/>	<input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	<input checked="" type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
		<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
	Decimal Representation		<input type="text" value="0.0126953125"/>
	Binary Representation		<input type="text" value="00111100010100000000000000000000"/>
	Hexadecimal Representation		<input type="text" value="0x3c500000"/>
	After casting to double precision		<input type="text" value="0.0126953125"/>

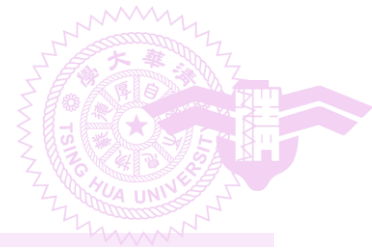
<https://www.h-schmidt.net/FloatConverter/IEEE754.html>

<http://babbage.cs.qc.cuny.edu/IEEE-754.old/Decimal.html>



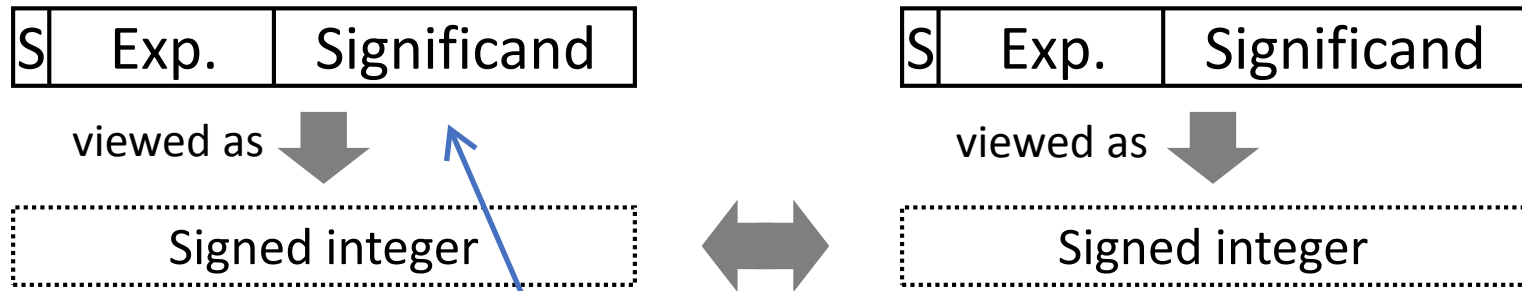
Floating Point Operations

- Comparisons
- Addition
- Multiplication



Comparisons

- Almost identical to **signed** integer comparison



- Rationales

- Positive > negative
- Between two **positive** floating point numbers
 - One with **larger** {exponent, significand} is greater
- Between two **negative** floating point numbers
 - One with **smaller** {exponent, significand} is greater

直接把他看成一個Signed Integer:
(1) 先看正負號
(2) 看Exponent
(3) 最後看Signifi cand



Comparisons (Cont'd)

- Cases directly supported by unsigned comparisons
 - $+\infty == +\infty$
 - $-\infty == -\infty$
 - $-\infty < \text{all numbers} < \infty$
- Special cases that signed comparisons do not directly support
 - **!=** involving any **NaN** yields **true**
 - **All other comparisons** involving NaN yield **false**
 - $\text{NaN} < 10 ? \rightarrow \text{false}$
 - $\text{NaN} > \text{NaN} ? \rightarrow \text{false}$
 - $\text{NaN} == \text{NaN} ? \rightarrow \text{false}$
 - $0 == -0$



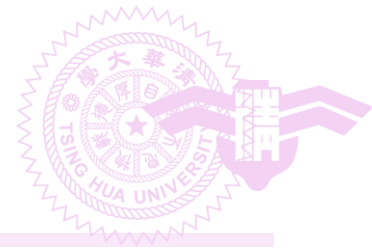
Addition

- Steps

1. Align (shift smaller number right)
2. Perform addition
3. Normalize
4. Round
5. Re-normalize

- Examples

- $9.999_{\text{ten}} \times 10^1 + 1.610_{\text{ten}} \times 10^{-1}$
- $1.101_{\text{two}} \times 2^9 + 1.110_{\text{two}} \times 2^{12}$
- Assume the four-digit significands



Decimal Example

Align
小數字往大數字對齊

Add

Normalize

Round

Renormalize

$$\begin{aligned} & 9.999_{\text{ten}} \times 10^1 & + 1.610_{\text{ten}} \times 10^{-1} \\ & = 9.999_{\text{ten}} \times 10^1 & + 0.01610_{\text{ten}} \times 10^1 \\ & = 10.01500_{\text{ten}} \times 10^1 \\ & = \underline{1.001500}_{\text{ten}} \times 10^2 \\ & = 1.002_{\text{ten}} \times 10^2 \\ & = 1.002_{\text{ten}} \times 10^2 \text{ (no change)} \end{aligned}$$

Renormalize 至多
做一次



Binary Example

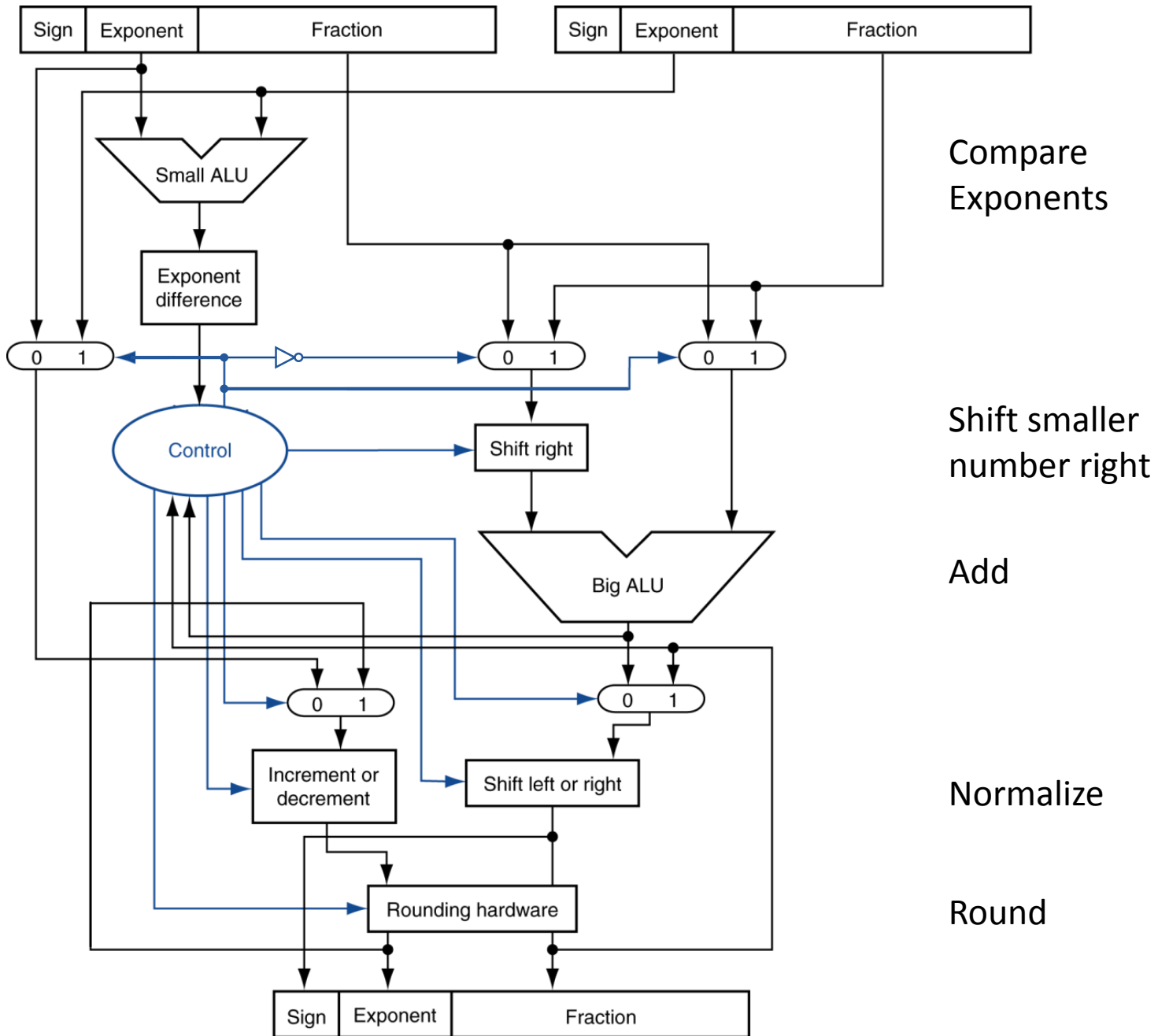
Align $1.101_{\text{two}} \times 2^9$ + $1.110_{\text{two}} \times 2^{12}$
 $= 0.001101_{\text{two}} \times 2^{12}$ + $1.110_{\text{two}} \times 2^{12}$

Add $= 1.111101_{\text{two}} \times 2^{12}$

Normalize $= \underline{1.111}101_{\text{two}} \times 2^{12}$ (no change)

Round $= 10.000_{\text{two}} \times 2^{12}$ 101大於一半(四捨五入)

Renormalize $= 1.000_{\text{two}} \times 2^{13}$





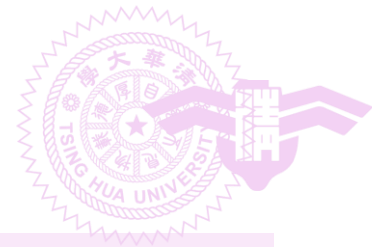
Multiplication

- Steps

1. Add exponents without bias
2. Multiply the significands (with sign determined)
3. Normalize (and check over/underflow)
4. Round
5. Re-normalize (and re-check over/underflow)

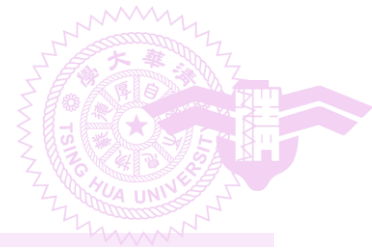
- Examples

- $1.110_{\text{ten}} \times 10^{10} \times 9.200_{\text{ten}} \times 10^{-5}$
- $1.000_{\text{two}} \times 2^{-1} \times (-1.110_{\text{two}}) \times 2^{-2}$
- Inputs and outputs have four-digit significand



Decimal Example

	$1.110_{\text{ten}} \times 10^{10}$	$\times 9.200_{\text{ten}} \times 10^{-5}$
Exponent	$10 + (-5) = 5$	
Multiply	$1.110_{\text{ten}} \times 9.200_{\text{ten}} = 10.212_{\text{ten}}$	
Normalize	$= \underline{1.0212}_{\text{ten}} \times 10^6$	
Round	$= 1.021_{\text{ten}} \times 10^6$	
Renormalize	$= 1.021_{\text{ten}} \times 10^6$ (no change)	



Binary Example

	$1.000_{\text{two}} \times 2^{-1}$	$\times (-1.110_{\text{two}}) \times 2^{-2}$
Exponent	$(-1) + (-2) = (-3)$	
Multiply	$1.000_{\text{two}} \times (-1.110_{\text{two}}) = (-1.110000_{\text{two}})$	
Normalize	$= (-\underline{1.110}000_{\text{two}}) \times 2^{-3}$ (no change)	
Round	$= (1.110_{\text{two}}) \times 2^{-3}$ (no change)	
Renormalize	$= (1.110_{\text{two}}) \times 2^{-3}$ (no change)	



Internal Format with Extra Bits

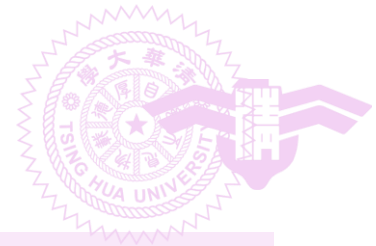
- **Extra bits** are needed during arithmetic operations to increase the arithmetic accuracy
- e.g., $1.101_{\text{two}} \times 2^9 + 1.110_{\text{two}} \times 2^{12}$

without extra bits

$$\begin{array}{r} 0.001_{\text{two}} \times 2^{12} \\ + 1.110_{\text{two}} \times 2^{12} \\ \hline = \underline{1.111}_{\text{two}} \times 2^{12} \end{array}$$

with extra bits

$$\begin{array}{r} 0.001101_{\text{two}} \times 2^{12} \\ + 1.110000_{\text{two}} \times 2^{12} \\ \hline = \underline{1.111101}_{\text{two}} \times 2^{12} \\ = 10.00_{\text{two}} \times 2^{12} \\ = 1.000_{\text{two}} \times 2^{13} \end{array}$$

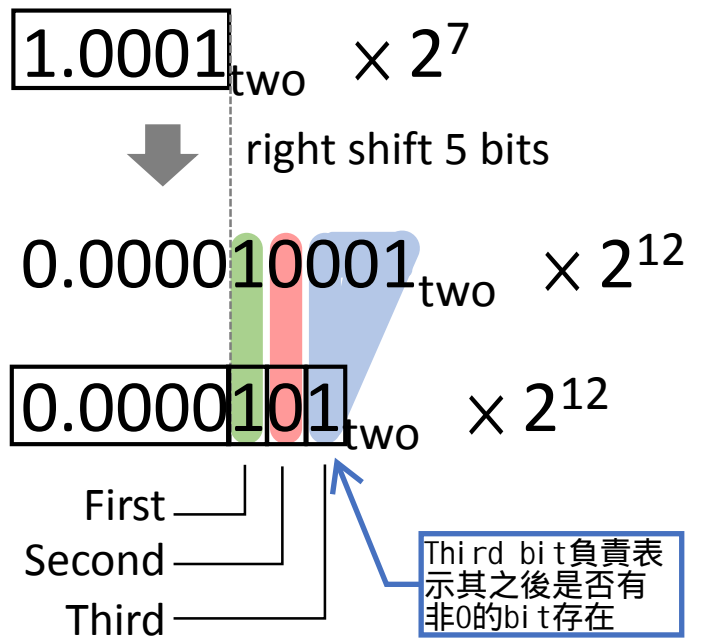


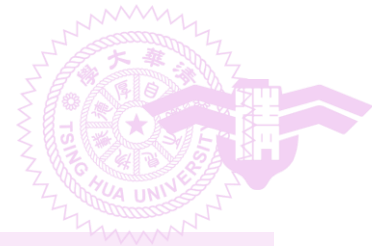
IEEE 754 Internal Format

- Three extra bits
 - The 3rd one represents any remaining nonzero bits to the right



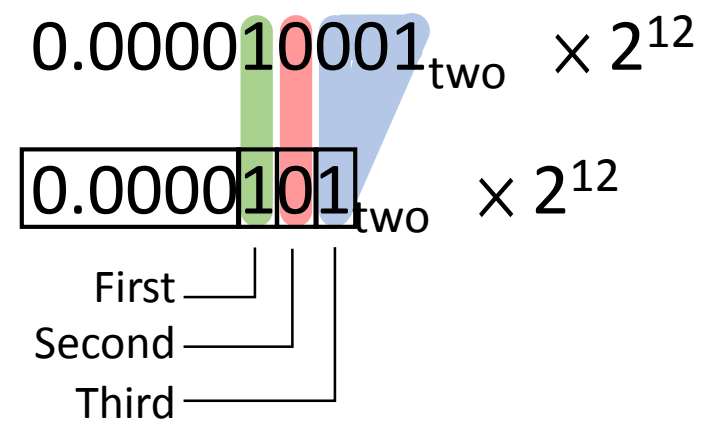
Internal format





IEEE 754 Internal Format

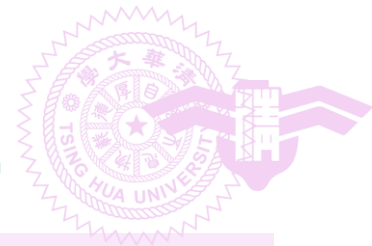
- Roles/names of the three extra bits
 - First: Guard
 - Second: Round
 - Third: Sticky





IEEE 754 Rounding Mode

- Four modes can be chosen by programmers
 - Toward 0 (also called **truncation**) 無條件捨去
 - Toward $+\infty$ 正數: 無條件進位、負數: 無條件捨去 (往正方向)
 - Toward $-\infty$ 負數: 無條件進位、正數: 無條件捨去 (往負方向)
 - Toward nearest even (**default mode**)
 - Choose the **even** one if there are two equally nearest values



Round Toward Nearest Even

- Binary examples

剛好一半的Case，看前面的數經過
運算後是否為偶數來決定

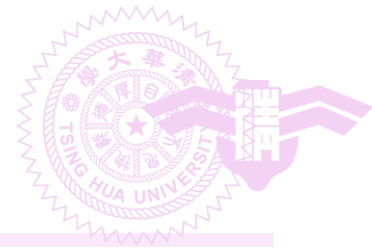
四捨六入

捨去後為偶數

進位後為偶數

	10.00 011	10.00 [↖] 101	10.10 100	10.11 [↖] 100
Results	10.00	10.01	10.10	11.00

- Reduce the statistical biases of rounding noises



Outline

- Overview
- IEEE 754 standard
 - Single-precision
 - Double-precision
 - Special numbers
- Floating-point operations
 - Addition
 - Multiplication
 - Rounding