
Floating Point Numbers

Philipp Koehn

HW4 - due Monday
Oct 28th

~~7 November 2016~~
23 Oct 2019



Numbers



- So far, we only dealt with integers
- But there are other types of numbers

Numbers



- So far, we only dealt with integers
- But there are other types of numbers
- Rational numbers (from ratio \simeq fraction)
 - $3/4 = 0.75$
 - $10/3 = 3.33333333\dots$

} use 2 integers!

Numbers



- So far, we only dealt with integers
- But there are other types of numbers
- Rational numbers (from ratio \simeq fraction)
 - $3/4 = 0.75$
 - $10/3 = 3.33333333\dots$
- Real numbers
 - $\pi = 3.14159265\dots$
 - $e = 2.71828182\dots$

integers = "countable"

} by definition, approximate

Very Large Numbers



- Distance of sun and earth

150,000,000,000 meters

- Scientific notation

1.5×10^{11} meters

- Another example: number of atoms in 12 gram of carbon-12 (1 mol)

$6.022140857 \times 10^{23}$

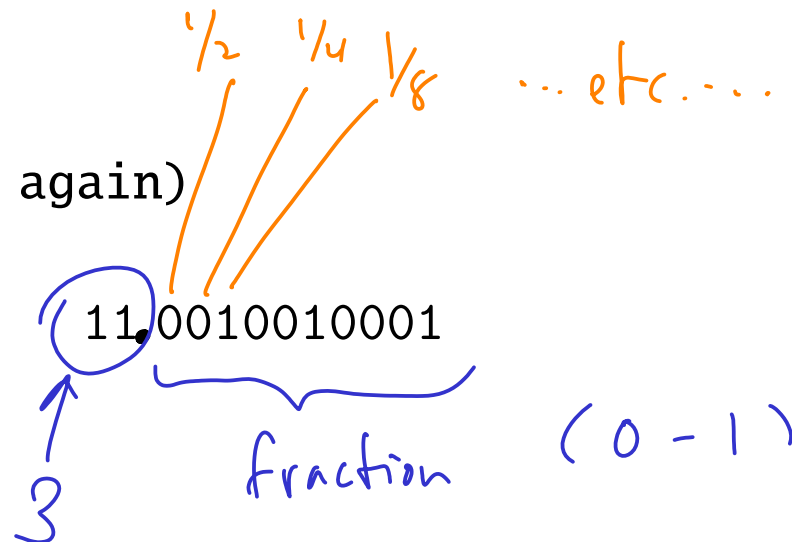


Binary Numbers in Scientific Notation

3



- Example binary number (π again)



- Scientific notation

$$1.10010010001 \times 2^1$$

- General form

$$1.x \times 2^y$$

Representation

- IEEE 754 floating point standard

- Uses 4 bytes (single precision) float

31	30	29	28	27	26	25	24	23	22	21	20	...	2	1	0
s	exponent								fraction						
1 bit	8 bits								23 bits						

1 = negative
0 = positive

- Exponent is offset with a bias of 127

e.g. $2^{-6} \rightarrow \text{exponent} = -6 + 127 = 121$

$1.x \cdot 2^y$

Conversion into Binary



- $\pi = 3.14159265$
- Number before period: $3_{10} = 11_2$
- Conversion of fraction .14159265

Conversion into Binary



- $\pi = 3.14159265$
- Number before period: $3_{10} = 11_2$
- Conversion of fraction .14159265

Digit Calculation

$$0.14159265 \times 2 \downarrow$$

Conversion into Binary



- $\pi = 3.14159265$
- Number before period: $3_{10} = 11_2$
- Conversion of fraction .14159265

Digit Calculation

	$0.14159265 \times 2 \downarrow$
0	0.2831853

Conversion into Binary



- $\pi = 3.14159265$
- Number before period: $3_{10} = 11_2$
- Conversion of fraction .14159265

Digit Calculation

	$0.14159265 \times 2 \downarrow$
0	$0.2831853 \times 2 \downarrow$
0	0.5663706

Conversion into Binary



- $\pi = 3.14159265$
- Number before period: $3_{10} = 11_2$
- Conversion of fraction .14159265

Digit Calculation

	$0.14159265 \times 2 \downarrow$
0	$0.2831853 \times 2 \downarrow$
0	$0.5663706 \times 2 \downarrow$
1	0.1327412

Conversion into Binary

- $\pi = 3.14159265$
- Number before period: $3_{10} = 11_2$
- Conversion of fraction .14159265

Digit	Calculation	Digit	Calculation
	$0.14159265 \times 2 \downarrow$	1	$0.9817472 \times 2 \downarrow$
0	$0.2831853 \times 2 \downarrow$	1	$0.9634944 \times 2 \downarrow$
0	$0.5663706 \times 2 \downarrow$	1	$0.9269888 \times 2 \downarrow$
1	$0.1327412 \times 2 \downarrow$	1	$0.8539776 \times 2 \downarrow$
0	$0.2654824 \times 2 \downarrow$	1	$0.7079552 \times 2 \downarrow$
0	$0.5309648 \times 2 \downarrow$	1	$0.4159104 \times 2 \downarrow$
1	$0.0619296 \times 2 \downarrow$	0	$0.8318208 \times 2 \downarrow$
0	$0.1238592 \times 2 \downarrow$	1	$0.6636416 \times 2 \downarrow$
0	$0.2477184 \times 2 \downarrow$	1	$0.3272832 \times 2 \downarrow$
0	$0.4954368 \times 2 \downarrow$	0	$0.6545664 \times 2 \downarrow$
0	$0.9908736 \times 2 \rightarrow$	1	0.3091328×2

- Binary: 11.001001000011111101101

Encoding into Representation

- π

$$1.1001001000011111101101 \times 2^1$$

- Encoding

Sign	Exponent	Fraction
0	10000000	1001001000011111101101

$$128 = (1 + 127)$$

- Note: leading 1 in fraction is omitted

bias

Special Cases



- Zero

Special Cases



- Zero
- Infinity ($1/0$)
- Negative infinity ($-1/0$)

Special Cases



- Zero
- Infinity ($1/0$)
- Negative infinity ($-1/0$)
- Not a number ($0/0$ or $\infty - \infty$)

NaN

Encoding



Exponent	Fraction	Object
0	0	zero
0	>0	denormalized number
1-254	anything	floating point number
255	0	infinity
255	>0	NaN (not a number)

} could be negative

(denormalized number: $0.x \times 2^{-126}$)

Double Precision



- Single precision = 4 bytes
- Double precision = 8 bytes

double

Sign	Exponent	Fraction
1 bit	8 bits	23 bits
1 bit	11 bits	52 bits



addition

Addition with Scientific Notation



- Decimal example, with 4 significant digits in encoding
- Example

$$0.1610 + 99.99$$

- In scientific notation

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

Addition with Scientific Notation

- Decimal example, with 4 significant digits in encoding
- Example

$$0.1610 + 99.99$$

- In scientific notation

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

- Bring lower number on same exponent as higher number

$$0.01610 \times 10^1 + 9.999 \times 10^1$$

Addition with Scientific Notation



- Round to 4 significant digits

$$0.016 \times 10^1 + 9.999 \times 10^1$$

Addition with Scientific Notation

12



- Round to 4 significant digits

$$0.016 \times 10^1 + 9.999 \times 10^1$$

- Add fractions

$$0.016 + 9.999 = 10.015$$

Addition with Scientific Notation

- Round to 4 significant digits

$$0.016 \times 10^1 + 9.999 \times 10^1$$

- Add fractions

$$0.016 + 9.999 = 10.015$$

- Adjust exponent

$$10.015 \times 10^1 = 1.0015 \times 10^2$$

Addition with Scientific Notation

- Round to 4 significant digits

$$0.016 \times 10^1 + 9.999 \times 10^1$$

- Add fractions

$$0.016 + 9.999 = 10.015$$

- Adjust exponent

$$10.015 \times 10^1 = 1.0015 \times 10^2$$

- Round to 4 significant digits

$$1.002 \times 10^2$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10}$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10}$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10} = 0.1_2$$

Binary Floating Point Addition

13



- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$$

$$-0.4375_{10} = -\frac{7}{16}_{10}$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$$

$$-0.4375_{10} = -\frac{7}{16}_{10} = -\frac{7}{2^4}_{10} = 0.0111_2$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$$

$$-0.4375_{10} = -\frac{7}{16}_{10} = -\frac{7}{2^4}_{10} = 0.0111_2 = -1.110_2 \times 2^{-2}$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$$

$$-0.4375_{10} = -\frac{7}{16}_{10} = -\frac{7}{2^4}_{10} = 0.0111_2 = -1.110_2 \times 2^{-2}$$

- Bring lower number on same exponent as higher number

$$-1.110 \times 2^{-2} = -0.111 \times 2^{-1}$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$$

$$-0.4375_{10} = -\frac{7}{16}_{10} = -\frac{7}{2^4}_{10} = 0.0111_2 = -1.110_2 \times 2^{-2}$$

- Bring lower number on same exponent as higher number

$$-1.110 \times 2^{-2} = -0.111 \times 2^{-1}$$

- Add the fractions

$$1.000_2 \times 2^{-1} + (-0.111 \times 2^{-1}) = 0.001 \times 2^{-1}$$

Binary Floating Point Addition

- Numbers

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$$

$$-0.4375_{10} = -\frac{7}{16}_{10} = -\frac{7}{2^4}_{10} = 0.0111_2 = -1.110_2 \times 2^{-2}$$

- Bring lower number on same exponent as higher number

$$-1.110 \times 2^{-2} = -0.111 \times 2^{-1}$$

- Add the fractions

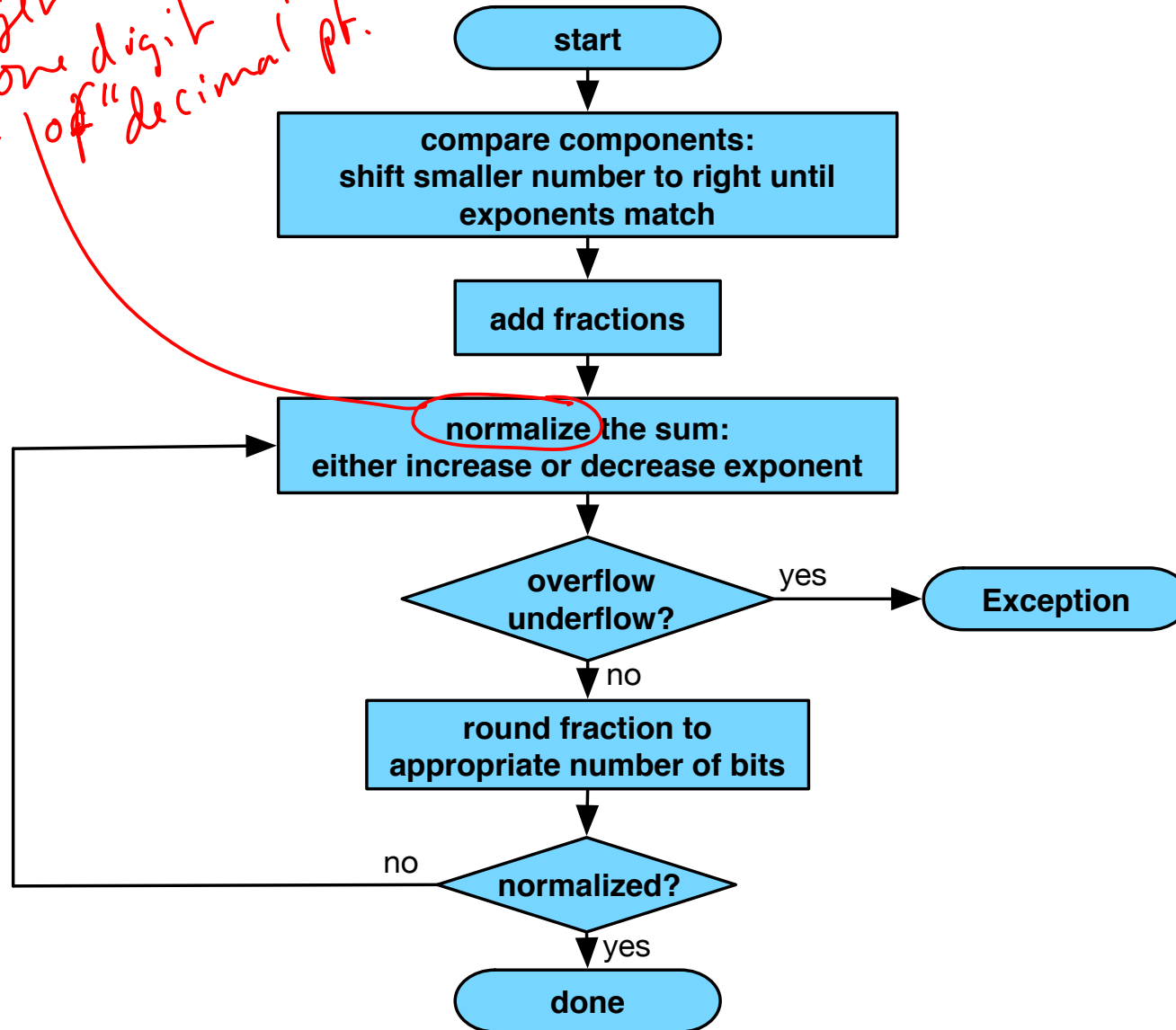
$$1.000_2 \times 2^{-1} + (-0.111 \times 2^{-1}) = 0.001 \times 2^{-1}$$

- Adjust exponent

$$0.001 \times 2^{-1} = \underline{1.000} \times 2^{-4}$$

Flowchart

try to get exactly one digit to left of "decimal" pt.





multiplication

Multiplication with Scientific Notation¹⁶



- Example: multiply 1.110×10^{10} and 9.200×10^{-5}

Multiplication with Scientific Notation¹⁶



- Example: multiply 1.110×10^{10} and 9.200×10^{-5}

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

Multiplication with Scientific Notation¹⁶



- Example: multiply 1.110×10^{10} and 9.200×10^{-5}

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

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

Multiplication with Scientific Notation¹⁶



- Example: multiply 1.110×10^{10} and 9.200×10^{-5}

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

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

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

Multiplication with Scientific Notation¹⁶



- Example: multiply 1.110×10^{10} and 9.200×10^{-5}

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

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

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

- Add exponents

$$-5 + 10 = 5$$

Multiplication with Scientific Notation¹⁶



- Example: multiply 1.110×10^{10} and 9.200×10^{-5}

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

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

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

- Add exponents

$$-5 + 10 = 5$$

- Multiply fractions

$$1.110 \times 9.200 = 10.212$$

Multiplication with Scientific Notation¹⁶



- Example: multiply 1.110×10^{10} and 9.200×10^{-5}

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

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

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

- Add exponents

$$-5 + 10 = 5$$

- Multiply fractions

$$1.110 \times 9.200 = 10.212$$

- Adjust exponent

$$\underline{10.212} \times 10^5 = 1.0212 \times 10^6$$

*need to
round to 4
significant
digits*

Binary Floating Point Multiplication

17



- Example

$$1.000 \times 2^{-1} \times -1.110 \times 2^{-2}$$

Binary Floating Point Multiplication

17



- Example

$$1.000 \times 2^{-1} \times -1.110 \times 2^{-2}$$

- Add exponents

$$-1 + (-2) = -3$$

Binary Floating Point Multiplication

17



- Example

$$1.000 \times 2^{-1} \times -1.110 \times 2^{-2}$$

- Add exponents

$$-1 + (-2) = -3$$

- Multiply fractions

$$1.000 \times -1.110 = -1.110$$

Binary Floating Point Multiplication

17



- Example

$$1.000 \times 2^{-1} \times -1.110 \times 2^{-2}$$

- Add exponents

$$-1 + (-2) = -3$$

- Multiply fractions

$$1.000 \times -1.110 = -1.110$$

$$1000 \times 1110 = 1110000$$

Binary Floating Point Multiplication

17



- Example

$$1.000 \times 2^{-1} \times -1.110 \times 2^{-2}$$

- Add exponents

$$-1 + (-2) = -3$$

- Multiply fractions

$$1.000 \times -1.110 = -1.110$$

$$1000 \times 1110 = 1110000$$

$$-1.110000$$

Binary Floating Point Multiplication

17



- Example

$$1.000 \times 2^{-1} \times -1.110 \times 2^{-2}$$

- Add exponents

$$-1 + (-2) = -3$$

- Multiply fractions

$$1.000 \times -1.110 = -1.110$$

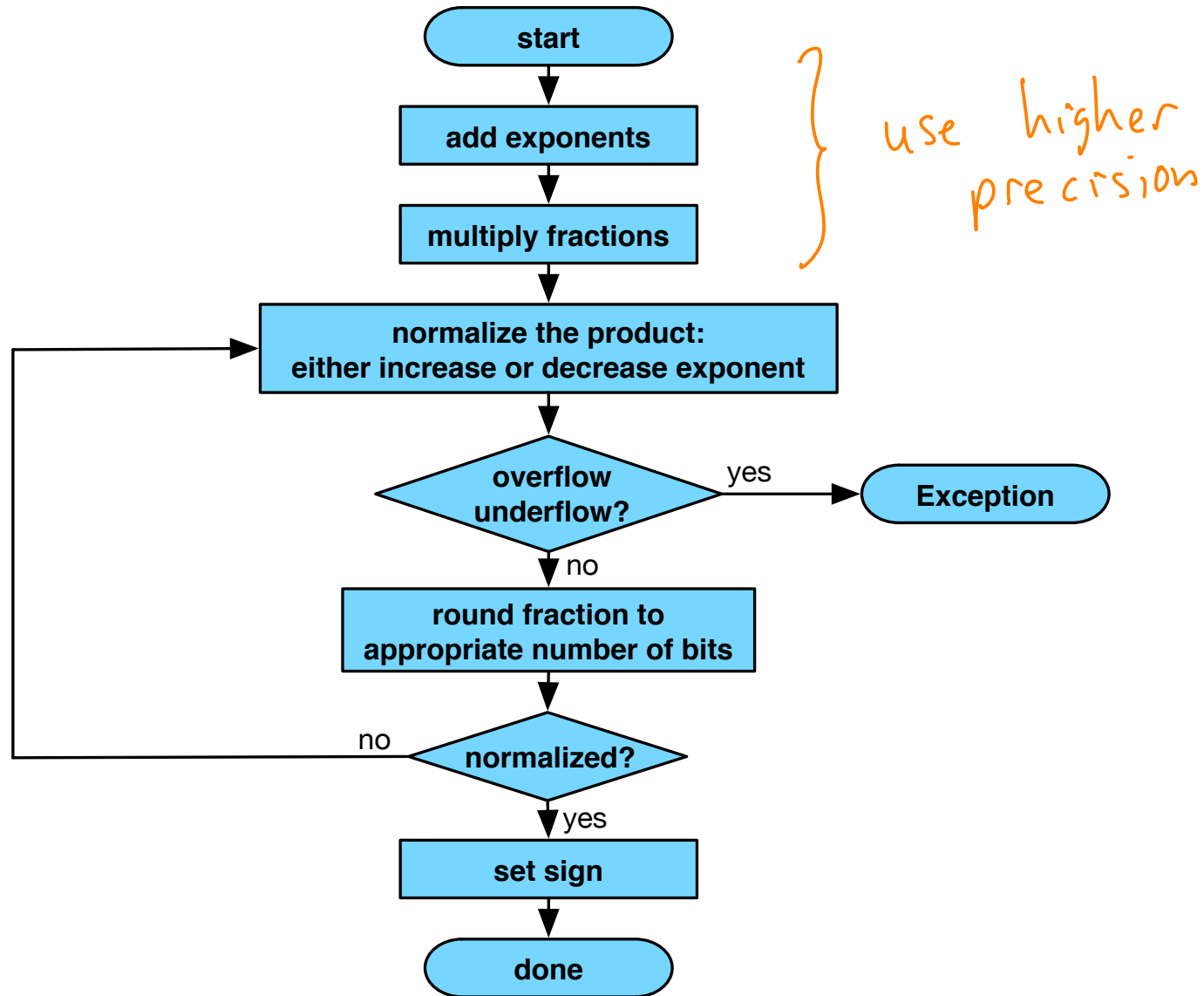
$$1000 \times 1110 = 1110000$$

$$\underline{-1.110000}$$

- Adjust exponent (not needed)

$$-1.110 \times 2^{-3}$$

Flowchart





Important idea
order / type of operations is significant

$a + b - b$ might not yield a

mips instructions

Instructions

- Both single precision (s) and double precision (d)
- Addition (add.s / add.d)
- Subtraction (sub.s / sub.d)
- Multiplication (mul.s / mul.d)
- Division (div.s / div.d)
- Comparison (c.x.s / c.x.d)
 - equality (x = eq), inequality (x = neq)
 - less than (x = lt), less than or equal (x = le)
 - greater than (x = gt), greater than or equal (x = ge)
- Floating point branch on true (bclt) or fals (bclf)

} sets a status bit!

Floating Point Registers

- MIPS has a separate set of registers for floating point numbers
- Little overhead, since used for different instructions
 - no need to specify in add, subtract, etc. instruction codes
 - different wiring for floating point / integer registers
 - much more limited use for floating point registers
(e.g., never an address)
- Double precision = 2 registers used

Example

- Conversion Fahrenheit to Celsius ($5.0/9.0 \times (x - 32.0)$)
- Input value x stored in register $\$f12$, constant in offsets to $\$gp$
- Code

lwcl

```
lwcl $f16, const5($gp) ; load 5.0
lwcl $f18, const9($gp) ; load 9.0
div.s $f16, $f16, $f18 ; $f16 = 5.0/9.0
lwcl $f18, const32($gp) ; load 32.0
sub.s $f18, $f12, $f18 ; $f18 = x-32.0
mul.s $f0, $f16, $f18 ; $f0 = result
```