# Finite representation of real numbers Floating-point numbers

Dr. Ing. Rodrigo Gonzalez

rodralez@ingenieria.uncu.edu.ar

Control y Sistemas

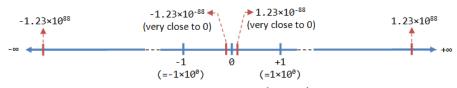
Facultad de Ingeniería, Universidad Nacional de Cuyo

## Summary

- Floating-point Representation
  - Floating-point Representation
  - Standards
  - Normalized Form
- Ploating-point Examples
- De-normalized Form
- Special values
- Rounding schemes
- Opposition of the property of the property
- Precision
- Precision problems

## Floating-point Representation

A floating-point number can represent a very large or a very small value, positive and negative.



Floating-point Numbers (Decimal)

A floating-point number is typically expressed in the scientific notation in the form of

$$(-1)^{\mathcal{S}} \times F \times r^{\mathcal{E}}$$
,

where,

- S, sign bit.
- F, fraction.
- E, exponent.
- r, certain radix. r = 2 for binary; r = 10 for decimal.

#### Standards

Modern computers adopt IEEE 754-2008 standard for representing floating-point numbers.

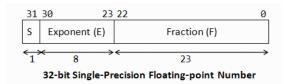
						IEE	E Sta	ndar	d P7	54 Fo	rmat					
Bit	31	30	29	28	27	26	25	24	23	22	21	20		2	1	0
	S	27	26	25	24	23	22	21	20	2 1	2 2	2 3		2 21	2 22	2-23
Sign	Sign (s) $\leftarrow$ Exponent (c) $\rightarrow$						$\leftarrow$ Fraction $(f) \rightarrow$									
							I	BM I	Form	at						
Bit	31	30	29	28	27	26	25	24	23	22	21	20		2	1	()
	S	26	<b>2</b> <sup>5</sup>	24	23	2 <sup>2</sup>	21	20	2-1	2 2	2-3	2 4		2 22	2 23	2 24
Sign					← Fraction $(f)$ →											
					DEC	(Dig	ital I	qui	omer	nt Co	rp.) F	orma	t			
Bit	31	30	29	28	27	26	25	24	23	22	21	20		2	1	0
	S	27	2 <sup>6</sup>	25	24	<b>2</b> <sup>3</sup>	2 <sup>2</sup>	21	20	2-2	2 -3	2 4		2 22	2 -23	2 24
Sigr	Sign (s) $\leftarrow$ Exponent (e) $\rightarrow$						$\leftarrow$ Fraction $(f) \rightarrow$									
						M	IL-S	TD 1	750	\ For	mat					
Bit	31	30	29		11	10	9	8	7	6	5	4	3	2	1	0
	20	2 1	2 2		2-20	2 21	2 -22	2 -23	27	26	25	24	23	22	21	20
	$\leftarrow$ Fraction $(f) \rightarrow$							$\leftarrow$ Exponent (e) $\rightarrow$								

#### IEEE 754-2008 standard

#### IEEE 754-2008 standard defines several formats.

	Binary form	nats (B = 2)	Decimal formats $(B = 10)$				
Parameter	Binary 16	Binary 32	Binary 64	Binary 128	Decimal 132	Decimal 164	Decimal 128
p, digits	10 + 1	23 + 1	52 + 1	112 + 1	7	16	34
$e_{max}$	+15	+127	+1023	+16383	+96	+384	+16,383
$e_{min}$	-14	-126	-1022	-16382	-95	-383	-16,382
Common name	Half precision	Single precision	Double precision	Quadruple precision			

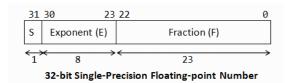
## IEEE-754 32-bit Single-Precision



$$(-1)^S \times F \times r^{(E-bias)}$$

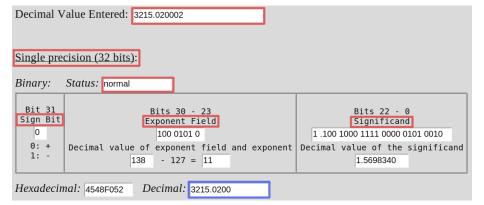
- S, sign bit. 0 for positive numbers and 1 for negative numbers.
- E, 8-bits exponent.
- We need to represent both positive and negative exponents.
- E = [1, 254], bias = 127;  $-126 \le E bias \le 127$ .
- E = 0 and E = 255 are reserved.
- F, 23-bits fraction.

#### Normalized Form



- Representation of a floating point number may not be unique:  $11.01_2 = 1.101_2 \times 2^1 = 110.1_2 \times 2^{-1}$ .
- Therefore, the fractional part *F* is normalized.
- 1.F, implicit leading 1.

#### Represent 3215.020002<sub>10</sub>



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

Represent  $3215.020002_{10} \times 2 = 6430.040004_{10}$ 

Decimal Value Entered: 6430.040004

#### Single precision (32 bits):

```
Binary:
          Status: normal
```

```
Bit 31
                          Bits 30 - 23
                                                                    Bits 22 - 0
Sign Bit
                         Exponent Field
                                                                    Significand
  0
                            10001011
                                                            1 .10010001111000001010010
  0: +
         Decimal value of exponent field and exponent
                                                         Decimal value of the significand
  1: -
                            - 127 = 12
                      139
                                                                     1.5698340
```

Hexadecimal: 45C8F052

Decimal: 6430.0400

Represent  $3215.020002_{10}/4 = 803.7550005_{10}$ 

Decimal Value Entered: 803.7550005

#### Single precision (32 bits):

Binary: Status: normal

```
Bit 31 | Bits 30 - 23 | Bits 22 - 0 | Significand | 1.0001000 | Decimal value of exponent field and exponent | Decimal value of exponent field and exponent | 1.5698340 | 1.5698340 | 1.5698340
```

Hexadecimal: 4448F052 Decimal: 803.75500

Floating-point numbers are auto-scaled!

#### De-normalized Form

## Not all real numbers in the range are representable



#### Normalized floating-point numbers



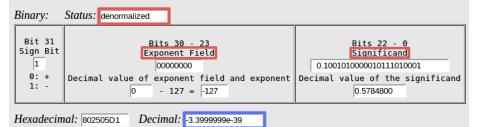
#### Denormalized floating-point numbers

- Normalized form has a serious problem.
- With an implicit leading 1 for the fraction, it cannot represent the number zero!
- De-normalized form is devised to represent zero and small numbers.
- $E = 0 \Rightarrow 0.F$ , implicit leading 0.

#### Represent -3.4E-39<sub>10</sub>

Decimal Value Entered: -3.4e-39

#### Single precision (32 bits):



## Special values

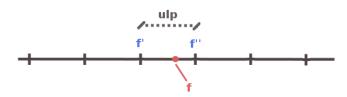
- **Zero**: E = 0, F = 0. Two representations: **+0** (S = 0) and **-0** (S = 1).
- Inf (Infinity): E = 0xFF, F = 0. Two representations: +Inf (S = 0) and -Inf (S = 1).
- NaN (Not a Number): E = 0xFF,  $F \neq 0$ . A value that cannot be represented as a real number (e.g. 0/0).

## **MATLAB**

- $0 \sim a = 1/0$
- 2 » ans = Inf

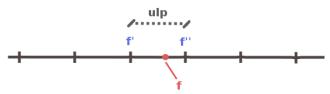
- $\bigcirc$  » c = 0/0

## Rounding schemes



- ulp (unit of least precision, eps ()).
- f, significant, f = 1.F.
- f' and f" being two successive multiples of ulp.
- Assume that f' < f < f'', f'' = f' + ulp,
- Then, the rounding function round(f) associates to f either f' or f'', according to some rounding strategy.

## Rounding schemes II



#### Rounding schemes are:

- Truncation (also called round toward 0 or chopping):
  - if f is positive, round(f) = f'.
  - if f is negative, round(-f) = f''.
- Round toward plus infinity: round(f) = f''.
- Round toward minus infinity: round(f) = f'.
- Round to nearest (default):
  - if f < f' + ulp/2, round(f) = f'.
  - if f > f' + ulp/2, round(f) = f''.

## Dynamic range

$$DR_{dB} \approx 6.02 \cdot 2^{b_E}$$

where  $b_E$  is the number of bits of E.

For single precision (32-bits):

$$DR_{dB} \approx 6.02 \cdot 2^8 \approx 1541 \, dB$$

For fixed-point Q31 (32-bits):

$$DR_{dB} \approx 6.02 \cdot 31 \approx 186 \, dB$$

#### Precision

- Precision is not constant throughout floating point numbers' range.
- As the numbers get larger, the precision gets worse.

#### **MATLAB**

```
    w u = linspace(-15,15,1000);

    w q = quantizer([6 4],'float'); % [wordlength exponentlength]

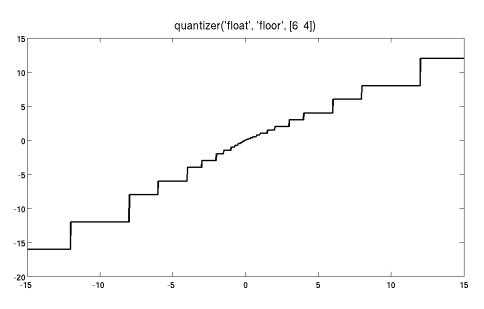
    w y1 = quantize(q,u);

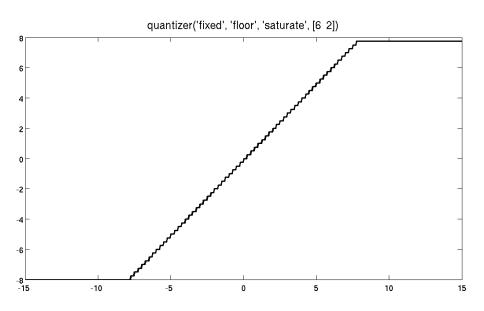
    w plot(u,y1); title(tostring(q))

    w q = quantizer('fixed',[6 2]); % [wordlength fractionlength]

    w y2 = quantize(q,u);

    w plot(u,y2); title(tostring(q))
```





#### Precision IV

eps(x) returns the positive distance from abs(x) to the next larger in magnitude floating point number of the same precision.

#### **MATLAB**

- $\bigcirc$  » e1 = eps(single(1))
- 2 » e1 = 1.1920929e-07
- 3 » e2 = eps(single(1e1))
- $\bigcirc$  » e3 = eps(single(1e10))
- $\bullet$  » e3 = 1024
- 0 > t = single(1e10) + single(1300)

## Sum of two floating-point positive numbers

Perform 0.5 + (-0.4375) using 4 bits for the mantissa.

$$0.5 = 0.1 \times 2^0 = 1.000 \times 2^{-1} \text{ (normalised)}$$
 
$$-0.4375 = -0.0111 \times 2^0 = -1.110 \times 2^{-2} \text{ (normalised)}$$

Matches with the exponent of the larger number:

$$-1.110 \times 2^{-2} = -0.1110 \times 2^{-1}$$

Add the mantissas:

$$1.000 \times 2^{-1} + -0.1110 \times 2^{-1} = 0.001 \times 2^{-1}$$

Normalise the sum, checking for overflow/underflow:

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

$$-126 \le -4 \le 127$$
 No overflow or underflow

Round the sum:

The sum fits in 4 bits so rounding is not required

## Sum of two floating-point positive numbers, II

Perform 1e10 + 1300 using IEEE-754 single precision.

Matches with the exponent of the larger number:

- Add the mantissas:

Normalise the sum, checking for overflow/underflow:

$$1.001010100000010111111001 \times r^{(160-127)}$$
  $-126 <= (160-127) <= 127$  No overflow or underflow

Round the sum:

The sum fits in 23 bits so rounding is not required

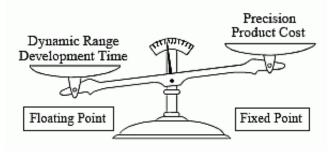
## More examples

 When calculations involve large and small numbers at the same time, the loss of precision affects the small number and the result.

## MATLAB

- 2 » ans = 0
- $\bigcirc$  » t = tan(x) sin(x)/cos(x)
- $\mathbf{6}$  » t = 2.2204e-16 % eps(1)

## Fixed-point vs floating-point



## Bibliography

1 Jean-Pierre Deschamps, Gustavo D. Sutter, and Enrique Cantó. Guide to FPGA Implementation of Arithmetic Functions, Chapter 12.