# Tutorial 04: Rounding and Normalization

*Computer Science Department*
*CS2208: Introduction to Computer Organization and Architecture*
*Fall 2022-2023*
*Instructor: Mahmoud R. El-Sakka*
*Office: MC-419*
*Email: elsakka@csd.uwo.ca*
*Phone: 519-661-2111 x86996*

# Rounding

❑ The rounding mechanisms include

○ *Truncation* (i.e., *dropping unwanted bits*) by *rounding towards zero*; a.k.a., *rounding down*

○ *Rounding towards positive or negative infinity*: the *nearest valid floating-point number* in the direction of positive infinity (for positive values) or negative infinity (for negative values) is chosen to decide the rounding; a.k.a., *rounding up*.

○ *Rounding to nearest*: the *closest valid floating-point number* to the actual value is used.

# Rounding

❑ *Example 1*: *Round to the nearest* the following numbers to 8 digits after the binary point.

```
0.110101011001000 ==> 0.11010101      0.110101001001000 ==> 0.11010100
                    + 0.00000001                          + 0.00000001
                    = 0.11010110                          = 0.11010101
```

> **1001000 > 1000000**

If it is == case, and this bit = 1, you round up.

If it is == case, and this bit = 0, you round down.

```
0.110101011000000 ==> 0.11010101      0.110101001000000 ==> 0.11010100
                    + 0.00000001                          + 0.00000000
                    = 0.11010110                          = 0.11010100
```

> **1000000 == 1000000**

**Mid-way ➜ round to even significand**

> **1000000 == 1000000**

```
0.1101010100xxxxxx ==> 0.11010101     0.1101010000xxxxxx ==> 0.11010100
                    + 0.00000000                          + 0.00000000
                    = 0.11010101                          = 0.11010100
```

> **0xxxxxx < 1000000**

> **0xxxxxx < 1000000**

3    *CS 2208: Introduction to Computer Organization and Architecture*

# Normalization

- ***Example 2***: Convert the unsigned value AB.BA$_{16}$ to binary. ***Normalize*** *your answer.*

AB.BA$_{16}$

➔ $10101011.10111010_2$

*After normalization,*

➔ $1.0101011101110102 \times 2^{+7}$

In base b, a normalized number will have the form

$\pm b_0 . b_1 b_2 b_3... \times b^n$

where $b_0 \neq 0$, and $b_0, b_1, b_2, b_3 ...$ are integers between 0 and b -1

| | |
|---|---|
| 0 = 0000 | |
| 1 = 0001 | |
| 2 = 0010 | |
| 3 = 0011 | |
| 4 = 0100 | |
| 5 = 0101 | |
| 6 = 0110 | |
| 7 = 0111 | |
| 8 = 1000 | |
| 9 = 1001 | |
| A = 1010 | |
| B = 1011 | |
| C = 1100 | |
| D = 1101 | |
| E = 1110 | |
| F = 1111 | |

# Normalization and Rounding

■ *__Example 3__*: Consider the unsigned normalized binary value $1.010101110111010_2 \times 2^{+7}$

  □ *limit it (using __truncation / rounding down__) to 6 bits (1 + 5 bits) in total*
  □ *limit it (using __rounding up__) to 6 bits (1 + 5 bits) in total*
  □ *limit it (using __rounding to the nearest__) to 6 bits (1 + 5 bits) in total*

  □ *limit it (using __truncation / rounding down__) to 9 bits (1 + 8 bits) in total*
  □ *limit it (using __rounding up__) to 9 bits (1 + 8 bits) in total*
  □ *limit it (using __rounding to the nearest__) to 9 bits (1 + 8 bits) in total*

  □ *limit it (using __truncation / rounding down__) to 14 bits (1 + 13 bits) in total*
  □ *limit it (using __rounding up__) to 14 bits (1 + 13 bits) in total*
  □ *limit it (using __rounding to the nearest__) to 14 bits (1 + 13 bits) in total*

Calculate the rounding error in each case.

Note that: The binary value $1.010101110111010_2 \times 2^{+7}$
$= 10101011.10111010_2 = AB.BA_{16}$

| | |
|---|---|
| 0 = 0000 | |
| 1 = 0001 | |
| 2 = 0010 | |
| 3 = 0011 | |
| 4 = 0100 | |
| 5 = 0101 | |
| 6 = 0110 | |
| 7 = 0111 | |
| 8 = 1000 | |
| 9 = 1001 | |
| A = 1010 | |
| B = 1011 | |
| C = 1100 | |
| D = 1101 | |
| E = 1110 | |
| F = 1111 | |

# Normalization and Rounding

- *Limiting the answer to 6 bits (1 + 5) in total,*
- ➔$1.01010\underline{1110111010}_2 \times 2^{+7}$

- ➔$1.01010_2 \times 2^{+7}$ *(using **truncation / rounding down**)*
  ➔$10101000_2$ ➔ $A8_{16}$
- ***Truncation** error = $AB.BA_{16} - A8_{16} = 3.BA_{16}$*

- ➔$1.01011_2 \times 2^{+7}$ *(using **rounding up**)*
  ➔$10101100_2$ ➔ $AC_{16}$
- ***Rounding up** error = $AB.BA_{16} - AC_{16} = -0.46_{16}$*

- As $1110111010_2 > 1000000000_2$
  ➔$1.01011_2 \times 2^{+7}$ *(using **rounding to the nearest**)*
  ➔$10101100_2$ ➔ $AC_{16}$
- ***Rounding to the nearest** error = $AB.BA_{16} - AC_{16} = -0.46_{16}$*

| | |
|---|---|
| 0 = 0000 | |
| 1 = 0001 | |
| 2 = 0010 | |
| 3 = 0011 | |
| 4 = 0100 | |
| 5 = 0101 | |
| 6 = 0110 | |
| 7 = 0111 | |
| 8 = 1000 | |
| 9 = 1001 | |
| A = 1010 | |
| B = 1011 | |
| C = 1100 | |
| D = 1101 | |
| E = 1110 | |
| F = 1111 | |

# Normalization and Rounding

- *Limiting the answer to 9 bits (1 + 8) in total,*
- ➔ $1.010101110111010_2 \times 2^{+7}$

- ➔ $1.010101111_2 \times 2^{+7}$ (*using **truncation / rounding down***)
  ➔ $10101011.1_2$ ➔ $AB.8_{16}$
- ***Truncation*** error = $AB.BA_{16} - AB.8_{16} = 0.3A_{16}$

- ➔ $1.01011000_2 \times 2^{+7}$ (*using **rounding up***)
  ➔ $10101100.0_2$ ➔ $AC_{16}$
- ***Rounding up*** error = $AB.BA_{16} - AC_{16} = -0.46_{16}$

- As $0111010_2 < 1000000_2$
  ➔ $1.010101111_2 \times 2^{+7}$ (*using **rounding to the nearest***)
  ➔ $10101011.1_2$ ➔ $AB.8_{16}$
- ***Rounding to the nearest*** error = $AB.BA_{16} - AB.8_{16} = 0.3A_{16}$

| | |
|---|---|
| 0 | = 0000 |
| 1 | = 0001 |
| 2 | = 0010 |
| 3 | = 0011 |
| 4 | = 0100 |
| 5 | = 0101 |
| 6 | = 0110 |
| 7 | = 0111 |
| 8 | = 1000 |
| 9 | = 1001 |
| A | = 1010 |
| B | = 1011 |
| C | = 1100 |
| D | = 1101 |
| E | = 1110 |
| F | = 1111 |

# Normalization and Rounding

- *Limiting the answer to 14 bits (1 + 13) in total,*
- ➔$1.010101110111010_2 \times 2^{+7}$

- ➔$1.01010111101110_2 \times 2^{+7}$ *(using **truncation / rounding down**)*
- ➔$10101011.101110_2$ ➔ $AB.B8_{16}$
- ***Truncation*** error $= AB.BA_{16} - AB.B8_{16} = 0.02_{16}$

- ➔$1.01010111101111_2 \times 2^{+7}$ *(using **rounding up**)*
- ➔$10101011.101111_2$ ➔ $AB.BC_{16}$
- ***Rounding up*** error $= AB.BA_{16} - AB.BC_{16} = -0.02_{16}$

- As $10_2 == 10_2$
- ➔$1.01010111101100_2 \times 2^{+7}$ *(using **rounding to the nearest**)*
- ➔$10101011.110110_2$ ➔ $AB.B8_{16}$
- ***Rounding to the nearest*** error $= AB.BA_{16} - AB.B8_{16} = 0.02_{16}$

  *round to even .*

- *Which **rounding mechanism** produces less error?*

  *nearest .*

| | |
|---|---|
| 0 = | 0000 |
| 1 = | 0001 |
| 2 = | 0010 |
| 3 = | 0011 |
| 4 = | 0100 |
| 5 = | 0101 |
| 6 = | 0110 |
| 7 = | 0111 |
| 8 = | 1000 |
| 9 = | 1001 |
| A = | 1010 |
| B = | 1011 |
| C = | 1100 |
| D = | 1101 |
| E = | 1110 |
| F = | 1111 |