## Floating Point Numbers

## IEEE 754 Representation

Single Precision: 32 bits

Double Precision: 64 bits

## Single Precision

```
 1    8      23        ← No. of bits
┌───┬─────┬──────────┐
│ S │ exp │ mantissa │
└───┴─────┴──────────┘
```

Sign
0 = +
1 = −

Exponent
Excess 127
127 = 7F

significand
mantissa       diff. names for the
fraction       same thing

Normalized representation
$1 < x < 10$ scientific notation

$A2B = A.2B \times 16^2$

$101.01 = 1.0101 \times 2^2$

$0.001 = 1.0 \times 2^{(-3)}$

mantissa = after .

before . digit implied
not stored

## example   Dec 80.5 in Single Precision

$\dfrac{80}{16} = 5$     Remainder = 0

$\dfrac{5}{16} = 0$      Remainder: 5

write number in this order

$80 = 50_{16}$

### Fractional Part

$0.5$
$\times 16$
$8.0$      when 0 stop else continue

$0.5 = 0.8_{16}$

$80.5 = 50.8_{16}$

$0101\ 0000\ .\ 1000$

$101\ 0000.10 \Rightarrow 1.\overbrace{01\ 0000\ 010}^{mantissa} \times 2^6$

mantissa        not stored

| | 127 | 7F |
|---|---|---|
| | 6 | 6 |
| | 133 | $85_{16}$ |

↑ exponent

S
0 1000 0101 01000010000 . 0
  4    2    A    1        ← we need to have 32 bit

+ve number

```
┌──────────────┐
│ 4 2 A1 00 00 │
└──────────────┘
          4
```

FP Nos / 1 – Oct 20, '20 (5 pages)

**Example 1**

$1 \quad 11 \quad 52 \quad \leftarrow$ size in bits

| | exp. | mantissa |
|---|---|---|

excess $1023 = 3FF_{16}$

Exponent is stored with an excess value of $1023$ $(3FF_{16})$

$80.5 \Rightarrow 50.8_{16}$

$\underbrace{101\ 0000}.\underbrace{1000} \Rightarrow \overbrace{1.01000010}^{\text{Mantissa}} \times 2^{6}$

Normalized
Scientific Representation
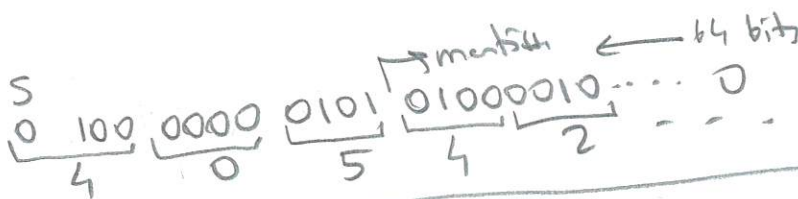
Excess $1023$ representation
for the exponent $6$

$$\begin{array}{cc} 1023 & 3FF \\ +\ 6 & +\ 6 \\ \hline 1029 & 405_{16} \end{array}$$

For example for decimal
number
$$0.072 \Rightarrow 7.2 \times 10^{-2}$$
$$1234.5 \Rightarrow 1.2345 \times 10^{3}$$
Normalized or Scientific
Representation

$\overset{S}{\underset{\phantom{x}}{0}} \underbrace{100}_{4} \underbrace{0000}_{0} \underbrace{0101}_{5} \overbrace{\underbrace{0100}_{4}\underbrace{0010}_{2} \cdots \cdots 0}^{\text{mantissa} \leftarrow 64 \text{ bits}}$

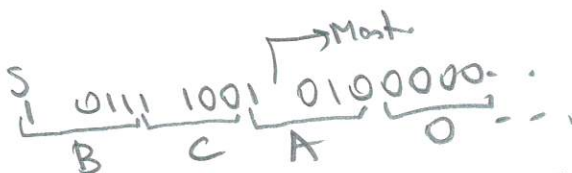| 40 54 20 00 00 00 00 00 00$_{16}$ |
|---|

**Example 2**

$-0.05_{16}$ show in single precision

$0.0000\underbrace{0101} \Rightarrow 1.01 \times 2^{-6}$
(not stored)

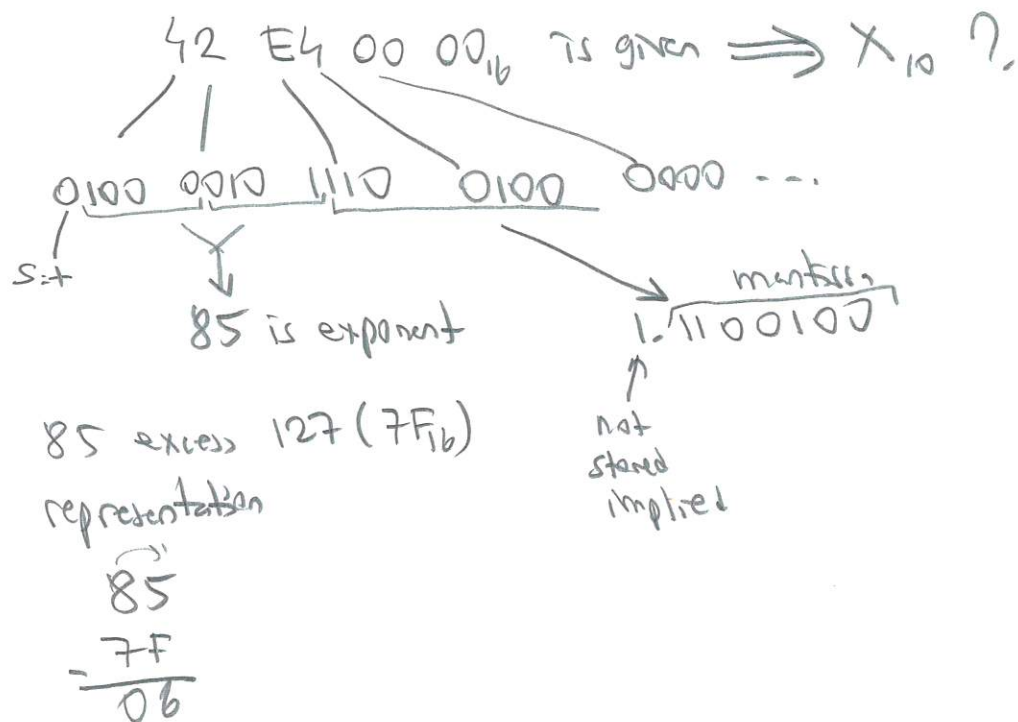$$\begin{array}{c} 7F \\ -6 \\ \hline 79 \end{array}$$

$\overset{S}{\underset{\phantom{x}}{1}} \underbrace{0111}_{B} \underbrace{1001}_{C} \overbrace{\underbrace{0100000}_{A}\underbrace{0}_{0} \cdots \cdots}^{\text{Mant.}}$

| B C A 0 00 00$_{16}$ |
|---|

**Example 3**

Try double precision: Is it | BF 94 00 00 00 00 00 00 00$_{16}$ | ?

FP Numbers /2 — Oct. 20, '20

## To Decimal Conversion

$42 \; E4 \; 00 \; 00_{16}$ is given $\Rightarrow X_{10}$ ?

0100 0010 1110 0100 0000 ...

S:t

85 is exponent

mantissa,
1. 1100100

not
stored
implied

85 excess 127 ($7F_{16}$)
representation

$$85$$
$$-7F$$
$$\overline{06}$$

1. 1100100 × $2^6$ $\Rightarrow$ 1110 0100

$$\underset{7}{\underleftarrow{1110}} \; \underset{2}{\underleftarrow{0010.0}} \Rightarrow 72_{16}$$

$72_{16} = 7 \times 16 + 2 = 112 + 2 \Rightarrow \boxed{114_{10}}$

## Problem to be solved

Given $A1 \; 49 \; 00 \; 00_{16}$

Interpret this number as a

1. Floating point number
2. Interpret as an integer
3. Interpret as a sign magnitude number

S magnitude

1 = -ve   0 = +ve

0100 $\Rightarrow$ +4
1100 $\Rightarrow$ -4
0000 $\Rightarrow$ +0
1000 $\Rightarrow$ -0

FP No/3 - Oct. 20, '20

==Decimal example==

$$9.995 \times 10^1 + 1.610 \times 10^{-1}$$

Step 1: Align the decimal point of the numbers with the larger exponent

$$9.999 \times 10^1 \qquad 1.610 \times 10^{-1} \Rightarrow 0.161 \times 10^0$$
$$\Rightarrow 0.016 \times 10^1$$

↑ lost a digit

Step 2    Add the new form of the numbers

$$9.999 \times 10^1$$
$$0.016 \times 10^1$$
$$\overline{10.015 \quad \times 10^1}$$

Step 3    Normalize the result

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

Assuming that we are allowed to keep three digits after decimal point

$$\Rightarrow 1.001 \times 10^2 \quad \underline{lost\ precision}$$

Example    $1.000_2 \times 2^{-1}$    $- 1.110 \times 2^{-2}$

called ↑ binary point

Step 1: Align the binary point of the numbers

$1.000_2 \times 2^{-1}$    $1.110 \times 2^{-2} \Rightarrow 0.111 \times 2^{-1}$

Step 2   Add numbers

$$1.000 \times 2^{-1}$$
$$-0.111 \times 2^{-1}$$
$$\overline{0.001 \times 2^{-1}}$$

borrow 1 comes as 10 -> decimal
borrow 1 comes as 2 -> binary

Step 4   Normalize result

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

---

Example

$1.011 \times 2^{-1} + 1.011 \times 2^{-6}$

$1.011 \times 2^{-6} \Rightarrow 0.1011 \times 2^{-5}$
$0.01011 \times 2^{-4}$
$0.001011 \times 2^{-3}$
$0.0001011 \times 2^{-2}$
$\underline{0.00001011 \times 2^{-1}}$

$$1.011 \times 2^{-1}$$
$$0.000 \times 2^{-1} \quad \longleftarrow lost$$
$$\overline{1.011 \times 2^{-1}}$$

FP No. 5 – Oct 20, '20