

HPY 201– Ψηφιακοί Υπολογιστές

Γ. Παπαευσταθίου

Floating Point Numbers
Examples

Exponential Notation

- The following are equivalent representations of 1,234

$$123,400.0 \times 10^{-2}$$

$$12,340.0 \times 10^{-1}$$

$$1,234.0 \times 10^0$$

$$123.4 \times 10^1$$

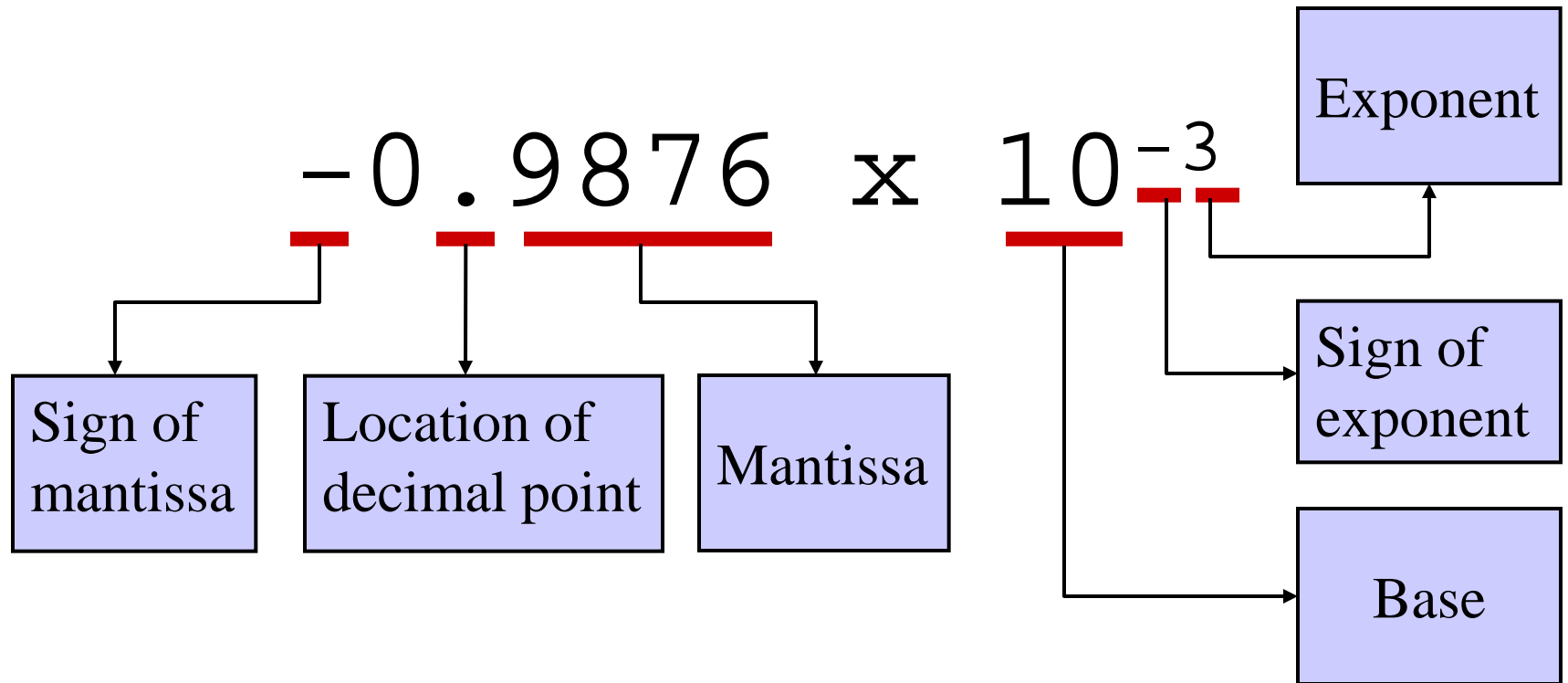
$$12.34 \times 10^2$$

$$1.234 \times 10^3$$

$$0.1234 \times 10^4$$

The representations differ in that the decimal place – the “point” -- “floats” to the left or right (with the appropriate adjustment in the exponent).

Parts of a Floating Point Number



IEEE 754 Standard

- Most common standard for representing floating point numbers
- Single precision: 32 bits, consisting of...
 - Sign bit (1 bit)
 - Exponent (8 bits)
 - Mantissa (23 bits)
- Double precision: 64 bits, consisting of...
 - Sign bit (1 bit)
 - Exponent (11 bits)
 - Mantissa (52 bits)

Representation

- Sign + exponent + coefficient

- | | | |
|---|-----|-------------|
| S | Exp | Coefficient |
|---|-----|-------------|

 - $1 + 8 + 23 = 32$ bits
 - $1 + 11 + 52 = 64$ bits (double precision)

The sign bit

- 0 indicates a positive number
- 1 a negative number

The exponent (I)

- 8 bits for single precision
- 11 bits for double precision
- With 8 bits, we can represent exponents between -126 and + 127
 - All-zeroes value is reserved for the zeroes and denormalized numbers
 - All-ones value are reserved for the infinities and NaNs (Not a Number)

The exponent (II)

- Exponents are represented using a biased notation
 - **Stored value = actual exponent + bias**
- For 8 bit exponents, bias is 127
 - Stored value of 1 corresponds to -126
 - Stored value of 254 corresponds to $+127$

0 and 255 are reserved for special values

The exponent (III)

- Biased notation simplifies comparisons:
 - If two normalized floating point numbers have different exponents, the one with the bigger exponent is the bigger of the two

Special values (I)

- *Signed zeroes:*
 - IEEE 754 distinguishes between $+0$ and -0
 - Represented by
 - Sign bit: 0 or 1
 - Biased exponent: all zeroes
 - Coefficient: all zeroes

Special values (II)

- *Denormalized numbers:*
 - Numbers whose coefficient cannot be normalized
 - *Smaller than 2^{-126}*
 - Will have a coefficient with leading zeroes and exponent field equal to zero
 - Reduces the number of significant digits
 - Lowers accuracy

Special values (III)

- *Infinities:*
 - $+\infty$ and $-\infty$
 - Represented by
 - Sign bit: 0 or 1
 - Biased exponent: all ones
 - Coefficient: all zeroes

Special values (IV)

- *NaN*:
 - For *Not a Number*
 - Often result from illegal divisions:
0/0, ∞/∞ , $\infty/-\infty$, $-\infty/\infty$, **and** $-\infty/-\infty$
 - Represented by
 - Sign bit: 0 or 1
 - Biased exponent: all ones
 - Coefficient: non zero

The coefficient

- Also known as *fraction* or *significand*
- *Most significant bit is always one*
 - Implicit and not *represented*
- Biased exponent is 127_{ten}
- True coefficient is *implicit one* followed by all zeroes

0	01...1	00000000000000000000000000000000
---	--------	----------------------------------

Decoding a floating point number

- Sign indicated by first bit
- Subtract 127 from biased exponent to obtain power of two:
 $\text{<be>} - 127$
- Use coefficient to construct a normalized binary value with a binary point:
 $1.\text{<coefficient>}$
- Number being represented is
 $1.\text{<coefficient>} \times 2^{\text{<be>} - 127}$

First example

[illegible]

- Sign bit is zero:
Number is positive
- Biased exponent is 127
Power of two is zero
- Normalized binary value is
1.0000000
- Number is **$1 \times 2^0 = 1$**

Second example

[illegible]

- Sign bit is zero:
Number is positive
- Biased exponent is 128
Power of two is 1
- Normalized binary value is
1.1000000
- Number is **$1.1 \times 2^1 = 11 = 3_{\text{ten}}$**

Third example

1	01...10	11000000000000000000000000000000
---	---------	----------------------------------

- Sign bit is 1:
Number is negative
- Biased exponent is 126
Power of two is -1
- Normalized binary value is
1.1100000
- Number is $-1.11 \times 2^{-1} = -0.111 = -7/8_{\text{ten}}$

Can we do it now?

0	129 _{ten}	10100000000000000000000000000000
---	--------------------	----------------------------------

- Sign bit is 0:
Number is _____
- Biased exponent is 129
Power of two is _____
- Normalized binary value is
1._____
- Number is _____

Encoding a floating point number

- Use sign to pick sign bit
- Normalize the number:
Convert it to form $1.\langle\textit{more bits}\rangle \times 2^{\langle\textit{exp}\rangle}$
- Add **127** to exponent $\langle\textit{exp}\rangle$ to obtain
biased exponent $\langle\textit{be}\rangle$
- Coefficient $\langle\textit{coeff}\rangle$ is equal to fractional part $\langle\textit{more bits}\rangle$ of number

First example

- Represent 7:
 - Convert to binary: 111
 - Normalize: 1.11×2^2
 - Sign bit is 0
 - Biased exponent is $127 + 2 = 10000001_{\text{two}}$
 - Coefficient is 1100...0

0	10...01	11000000000000000000000000000000
---	---------	----------------------------------

Second example

- Represent $1/2$
 - Convert to binary: 0.1
 - Normalize: 1.0×2^{-1}
 - Sign bit is 0
 - Biased exponent is $127 - 1 = 01111110_{\text{two}}$
 - Coefficient is $00\dots0$

0	01...10	00000000000000000000000000000000
---	---------	----------------------------------

Third example

- Represent **-2**
 - **Convert to binary:** **10**
 - **Normalize:** **1.0×2^1**
 - **Sign bit is 1**
 - **Biased exponent is $127 + 1 = 10000000$** _{two}
 - **Coefficient is 00...0**

1	10...00	00000000000000000000000000000000
---	---------	----------------------------------

Fourth example

- Represent **9/4**
 - **Convert to binary:** 1001×2^{-2}
 - **Normalize:** 1.001×2^1
 - **Sign bit is 0**
 - **Biased exponent is $127 + 1 = 10000000_{\text{two}}$**
 - **Coefficient is 0010...0**

1	10...00	00100000000000000000000000000000
---	---------	----------------------------------

Can we do it now?

- Represent **6.25**:

- **Convert to binary:**

- **Normalize:**

$$\frac{\quad}{1.\quad} \times 2^{\quad}$$

- **Sign bit is** \quad

- **Biased exponent is** $127 + \quad = \quad_{\text{ten}}$

- **Coefficient is** \quad



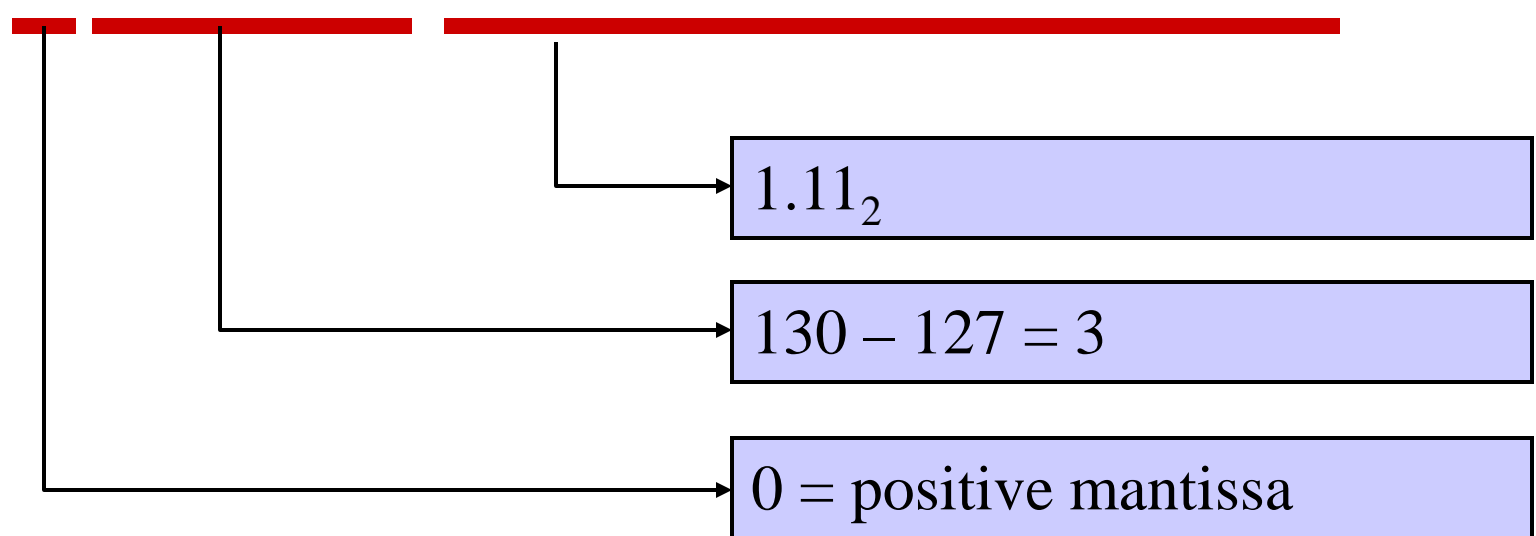
Range

- Can represent numbers between $1.00\dots0 \times 2^{-126}$ and $1.11\dots1 \times 2^{127}$
 - Say between 2^{-126} and 2^{128}
- **Observing that $2^{10} \approx 10^3$**
we divide the exponents by 10 and multiply them by 3 to obtain the interval expressed in powers of 10
 - Approximate range is 10^{-38} to 10^{38}

More Examples

- Single precision

0 10000010 110000000000000000000000



$$+1.11_2 \times 2^3 = 1110.0_2 = 14.0_{10}$$

Hexadecimal

- It is convenient and common to represent the original floating point number in hexadecimal
- The preceding example...

0	100	00001	0	110	00000	00000	00000	00000	00000
4	1	6	0	0	0	0	0	0	

Go all the way again

- E.g., What decimal value is represented by the following 32-bit floating point number?

$C17B0000_{16}$

- Express in binary and find S, E, and M

$$\text{C17B0000}_{16} =$$

1 10000010 111101100000000000000000₂

S E M

1 = negative
0 = positive

- Step 2
 - Find “real” exponent, n
 - $n = E - 127$
 $= 10000010_2 - 127$
 $= 130 - 127$
 $= 3$

- Step 3
 - Put S, M, and n together to form binary result
 - (Don't forget the implied “1.” on the left of the mantissa.)

$$-1.1111011_2 \times 2^n =$$

$$-1.1111011_2 \times 2^3 =$$

$$-1111.1011_2$$

- Step 4
 - Express result in decimal

$$\begin{array}{r} \underline{-1111} . \underline{1011}_2 \\ -15 \end{array}$$

$2^{-1} = 0.5$
 $2^{-3} = 0.125$
 $2^{-4} = \underline{0.0625}$
 0.6875

Answer: -15.6875

Converting to Floating Point

- E.g., Express 36.5625_{10} as a 32-bit floating point number (in hexadecimal)

- Step 1
 - Express original value in binary

$$36.5625_{10} =$$

$$100100.1001_2$$

- Step 2
 - Normalize

$$100100.1001_2 =$$

$$1.001001001_2 \times 2^5$$

- Step 3
 - Determine S, E, and M

$$\underbrace{+1}_S . \underbrace{001001001}_M_2 \times 2^{\underbrace{5}_n}$$

$$\begin{aligned} E &= n + 127 \\ &= 5 + 127 \\ &= 132 \\ &= 10000100_2 \end{aligned}$$

$S = 0$ (because the value is positive)

- Put S, E, and M together to form 32-bit binary result

<u>0</u>	<u>10000100</u>	<u>00100100100000000000000000000000</u> ₂
S	E	M

- Step 5
 - Express in hexadecimal

$$\begin{array}{cccccccc}
 0 & 10000100 & 00100100 & 10000000 & 00000000 & 00000000 & 00000000 & 00000000_2 = \\
 0100 & 0010 & 0001 & 0010 & 0100 & 0000 & 0000 & 0000_2 = \\
 4 & 2 & 1 & 2 & 4 & 0 & 0 & 0_{16}
 \end{array}$$

Answer: 42124000₁₆

Double precision arithmetic (I)

- Use 64-bit double words
- Allows us to have
 - One bit for sign
 - Eleven bits for exponent
 - 2,048 possible values
 - Fifty-two bits for coefficient
 - Plus the implicit leading bit

Double precision arithmetic (II)

- Exponents are still represented using a biased notation
 - **Stored value = actual exponent + bias**
- For 11-bit exponents, bias is **1023**
 - Stored value of 1 corresponds to **-1,022**
 - Stored value of 2,046 corresponds to **+1,023**
 - Stored values of 0 and 2,047 are reserved for special cases

Double precision arithmetic (III)

- Can now represent numbers between $1.00\dots0 \times 2^{-1,022}$ and $1.11\dots1 \times 2^{1,203}$
 - Say between $2^{-1,022}$ and $2^{1,204}$
 - Approximate range is 10^{-307} to 10^{307}
 - In reality, more like 10^{-308} to 10^{308}

Double precision arithmetic (IV)

- We now have 53 significant bits
 - Theoretical precision of $1/2^{53}$. that is, roughly $1/10^{16}$
- Can now add correctly billions or trillions

If that is now enough, ...

- Can use 128-bit quad words
- Allows us to have
 - One bit for sign
 - Fifteen bits for exponent
 - From **-16382** to **+16383**
 - One hundred twelve bits for coefficient
 - Plus the implicit leading bit

Decimal floating point addition (I)

- $5.25 \times 10^3 + 1.22 \times 10^2 = ?$
- Denormalize number with smaller exponent:
 $5.25 \times 10^3 + 0.122 \times 10^3$
- Add the numbers:
 $5.25 \times 10^3 + 0.122 \times 10^3 = 5.372 \times 10^3$
- Result is normalized

Decimal floating point addition (II)

- $9.25 \times 10^3 + 8.22 \times 10^2 = ?$
- Denormalize number with smaller exponent:
 $9.25 \times 10^3 + 0.822 \times 10^3$
- Add the numbers:
 $9.25 \times 10^3 + 0.822 \times 10^3 = 10.072 \times 10^3$
- Normalize the result:
 $10.072 \times 10^3 = 1.0072 \times 10^4$

Binary floating point addition (I)

- Say $1001 + 10$ or $1.001 \times 2^3 + 1.0 \times 2^1$
- Denormalize number with smaller exponent:
 $1.001 \times 2^3 + 0.01 \times 2^3$
- Add the numbers:
 $1.001 \times 2^3 + 0.01 \times 2^3 = 1.011 \times 2^3$
- Result is normalized

Binary floating point addition (II)

- Say $101 + 11$ or $1.01 \times 2^2 + 1.1 \times 2^1$
- Denormalize number with smaller exponent:
 $1.01 \times 2^2 + 0.11 \times 2^2$
- Add the numbers:
 $1.01 \times 2^2 + 0.11 \times 2^2 = 10.00 \times 2^2$
- Normalize the results
 $10.00 \times 2^2 = 1.000 \times 2^3$

Binary floating point subtraction

- Say $101 - 11$ or $1.01 \times 2^2 - 1.1 \times 2^1$
- Denormalize number with smaller exponent:
 $1.01 \times 2^2 - 0.11 \times 2^2$
- Perform the subtraction:
 $1.01 \times 2^2 - 0.11 \times 2^2 = 0.10 \times 2^2$
- Normalize the results
 $0.10 \times 2^2 = 1.0 \times 2^1$

Decimal floating point multiplication

- Exponent of product is the *sum* of the exponents of multiplicand and multiplier
- Coefficient of product is the product of the coefficients of multiplicand and multiplier
- Compute sign using usual rules of arithmetic
- May have to renormalize the product

Decimal floating point multiplication

- $6 \times 10^3 \times 2.5 \times 10^2 = ?$
- Exponent of product is:
 $3 + 2 = 5$
- Multiply the coefficients:
 $6 \times 2.5 = 15$
- **Result will be positive**
- Normalize the result:
 $15 \times 10^5 = 1.5 \times 10^6$

Binary floating point multiplication

- Exponent of product is the *sum* of the exponents of multiplicand and multiplier
- Coefficient of product is the product of the coefficients of multiplicand and multiplier
- Compute sign using usual rules of arithmetic
- May have to renormalize the product

Binary floating point multiplication

- Say 110×11 or $1.1 \times 2^2 \times 1.1 \times 2^1$
- Exponent of product is:
$$2 + 1 = 3$$
- Multiply the coefficients:
$$1.1 \times 1.1 = 10.01$$
- Result will be positive
- Normalize the result:
$$10.01 \times 2^3 = 1.001 \times 2^4$$

FP division

- Very tricky
- One good solution is to multiply the dividend by the inverse of the divisor

A trap

- Addition does not necessarily commute:

- $-9 \times 10^{37} + 9 \times 10^{37} + 4 \times 10^{-37}$

- Observe that

- $(-9 \times 10^{37} + 9 \times 10^{37}) + 4 \times 10^{-37} = 4 \times 10^{-37}$

while

- $-9 \times 10^{37} + (9 \times 10^{37} + 4 \times 10^{-37}) = 0$

due to the limited accuracy of FP numbers