# ECS404: Computer Systems and Networks 2016 Week 5

Signed Integers, Floating Point, Character Sets

# This week

- Number representation

  - Brief recap on Unsigned

  - Signed

  - Floating Point

  - Rounding errors: Why you should be careful

- Character Sets

  - Common character sets and their structure.

# Learning Objectives: Signed Integers

- Use of 2's complement to represent signed integers.

- Basic structure of 2's complement representation: which numbers can be represented, and the geometrical relationship between numbers and representations.

- Why we use it rather than sign and magnitude: how operations are implemented.

- How to encode and decode numbers into 2's complement: algorithms

# Learning Objectives: Floating Point

- Floating point real numbers

- Scientific representation and its structure

- IEEE binary floating point numbers

- Rounding errors and problems with fixed width floating point systems.

# Learning Objectives: Character Sets

- Text representation:

- ASCII, ISO-xxxx and Unicode-based character sets: numbers of bits used and basic facts about representations.

# Number systems and recap

# Three types of numbers

- Unsigned: we can think of these as positive integers (whole numbers)

- Signed: regular integers, can be positive or negative

- Floating point: real numbers, not necessarily whole numbers.

# Three types of numbers

| class | examples | representation | Java |
|---|---|---|---|
| unsigned | 0,1,2,3.. | unsigned binary | |
| signed | ..-3,-2,-1,0,1,2,3.. | 2's complement | byte, short, int, long |
| floating point | -2.80,1.00,3.14 | IEEE floating point | float, double |

# Java has 8 primitive datatypes

- byte: 8 bit signed 2's complement

- short: 16 bit signed 2's complement

- int: 32 bit signed 2's complement

- long: 64 bit signed 2's complement

- float: 32 bit IEEE floating point

- double: 64 bit IEEE floating point

# C has a subtly different collection

- unsigned - unsigned integers

- int, long - signed integers

- float, double - floating point

All of these have fixed size (typically 32 or 64 bits), but can differ between implementations.

# Examples: unsigned

- 32 bit: 4 8-byte blocks, numbers from 0 to $2^{32}$-1

- Standard binary representation

- 00000000 00000000 00000000 00001001 = 9

# Operations: unsigned

- Can use standard long addition to add (and subtraction to subtract, though we did not see that)

- Can use standard long multiplication to multiply… but binary version is easier because we only have to copy (shift) multiplicand.

# Operations: unsigned

- Equality testing is easy: equal if and only if bit patterns the same

- Checking for < is easy: scan from left and first difference gives the order.

# Operations: unsigned

- Checking for < is easy: scan from left and first difference gives the order.

00010101011110011101011

00010101101010001101011

First difference shows top larger than bottom.

# Signed Integers
# (2's complement)

# Negative numbers

- Of course we also need negative numbers.

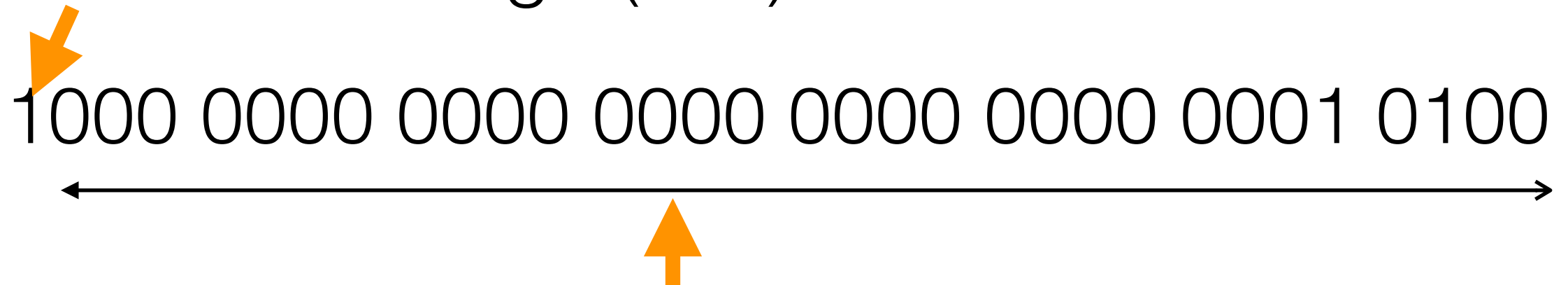- Let's think about how we might do that.

# Obvious answer

- Do what we do…

- Use one bit to give the sign, and the rest to give the magnitude.

# Obvious answer

First bit is sign (1=-)

1000 0000 0000 0000 0000 0000 0001 0100

Remaining 31 bits are magnitude (20)

So this is -20

# Problem 1

- We have two zeroes

- This makes some of our most common operations harder: testing for 0 and testing for equality

+0 = 0000 0000 0000 0000 0000 0000 0000 0000

-0 = 1000 0000 0000 0000 0000 0000 0000 0000

# Problem 2

- The algorithm for adding two positive or two negative numbers is basically addition.

- The algorithm for adding a positive number to a negative one is basically subtraction.

- So our circuitry is more complicated, and probably slower.

# Example

- Suppose we are using **8 bits**, and we use the first bit as a sign 0 is positive, and 1 is negative, say.

- 00001001 represents 9

- 10001001 represents -9

- 00000101 represents 5

- 10000101 represents -5

# Example

- To compute 00001001 + 00000101 we add the last seven bits: 0001001 + 0000101

- Similarly to compute 10001001 + 10000101 we add the last seven bits.

- To compute 00001001 + 10000101 we subtract the last seven bits: 0001001 - 0000101 (and must check we get the sign right).

- To compute 10001001 + 00000101 we subtract the last seven bits in the other order: 0000101 - 0001001 (and again must check we get the sign right)… though we could use the above and just negate at the last minute.

- In any case, this is quite complicated.

- there is a better way: 2's complement

# 2's complement

- 32 bit sign and magnitude uses unsigned $0...2^{31}-1$ to represent signed $0...2^{31}-1$, and the unsigned $2^{31}...2^{32}-1$ to represent signed $0...-(2^{31}-1)$ (decreasing)

- 32 bit 2's complement uses unsigned $0...2^{31}-1$ to represent signed $0...2^{31}-1$, and the unsigned $2^{31}...2^{32}-1$ to represent signed $2^{31}..-1$ (increasing)

**Unsigned**

$0$ $\qquad$ $2^{31}\text{-}1$ $\quad$ $2^{31}$ $\qquad\qquad$ $2^{32}\text{-}1$

bit pattern: 000...000

bit pattern: 011...111

bit pattern: 100...000

bit pattern: 111...111

**Unsigned**

0    $2^{31}-1$  $2^{31}$    $2^{32}-1$

bit pattern: 000...000

bit pattern: 011...111

bit pattern: 100...000

bit pattern: 111...111

$-(2^{31}-1)$    0    $2^{31}-1$

**Sign and magnitude**

**Unsigned**

$0$            $2^{31}-1$   $2^{31}$            $2^{32}-1$

bit pattern: 000...000

bit pattern: 011...111

bit pattern: 100...000

bit pattern: 111...111

$-2^{31}$            $-1\ 0$            $2^{31}-1$

**2's complement**

# 2's complement

- We have not duplicated any representations.

- (So tests for 0 and equality will be simple).

- We have not reversed any line segments.

# 2's complement

- Positive signed 0 to $2^{31}-1$ are represented by themselves.

- Negative signed $-2^{31}$ to $-1$ are represented by (themselves plus $2^{32}$).

- So -3 is represented by $-3+2^{32} = 2^{32}-3$.

- Note that this is between 0 and $2^{32}$.

# 8 bit 2's complement

- For exercises we will use 8 bit 2's complement.

- 8 bit ($2^7$ = 128)

- 8 bit unsigned can represent numbers 0..255

- 8 bit 2's complement represents $-2^7$=-128 .. 127=$2^7$-1

unsigned

| 0 | .... | 127 | 128 | .... | 255 |
|---|------|-----|-----|------|-----|
| 0 | .... | 127 | -128 -127 | .... | -2 -1 |

signed

# 2's complement

- 32 bit

- 32 bit unsigned can represent numbers $0..2^{32}-1$

- 32 bit 2's complement represents $-2^{31} .. 2^{31}-1$

unsigned

| 0 | .... | $2^{31}-1$ | $2^{31}$ | .... | $2^{32}-1$ |
|---|------|------------|----------|------|------------|

| 0 | .... | $2^{31}-1$ | $-2^{31}$ | $-2^{31}+1$ | .... | $-2$ $-1$ |

signed

# 2's complement

- 32 bit

- On this bit signed = unsigned

unsigned

$0$ .... $2^{31}-1$ $2^{31}$ .... $2^{32}-1$

$0$ .... $2^{31}-1$ $-2^{31}$ $-2^{31}+1$ .... $-2$ $-1$

signed

# 2's complement

- 32 bit

- On this bit signed = unsigned - $2^{32}$

unsigned

$$0 \quad\quad .... \quad\quad 2^{31}-1 \quad 2^{31} \quad\quad .... \quad\quad 2^{32}-1$$

$$0 \quad\quad .... \quad\quad 2^{31}-1 \; -2^{31} \; -2^{31}+1 \quad .... \quad\quad -2 \; -1$$

signed

# 2's complement

- So it is always the case that if we take any 32-bit sequence, then the (unsigned value) - (signed value) is divisible by $2^{32}$

- In other words the signed value = unsigned value mod $2^{32}$

# Operations mod n

Mathematical Property:

If $a1 = a2 \bmod n$ and $b1 = b2 \bmod n$, then

- $a1 + b1 = a2 + b2 \bmod n$

- $a1 - b1 = a2 - b2 \bmod n$

- $a1 * b1 = a2 * b2 \bmod n$

# Practical consequence

- With 2's complement, we can use unsigned addition, subtraction, multiplication algorithms to implement signed addition, subtraction, multiplication.

# Translating between unsigned and 2's complement

- There is more than one way.

- You ALWAYS have to distinguish between:

- signed - distinguish between positive and negative (>= 0 and <0)

- unsigned - distinguish between $< 2^{n-1}$ and $>= 2^{n-1}$

# To get the 2's complement bit pattern representing a signed integer

- We will use **8 bit** 2's complement.

- Case 1: integer is positive

- Method: convert to unsigned binary. Pad with initial zeroes to correct length.

- Example: 20

- Convert to unsigned binary: 10100

- Pad with initial zeroes: 00010100

# To get the 2's complement bit pattern representing a signed integer

- We will use **8 bit** 2's complement.

- Case 2: integer is negative

- Goal: if integer is x, then representation will be unsigned binary translation of $2^n+x$

- Example: x=-20, then representation is $2^8+(-20)$ = 128 - 20

# To get the 2's complement bit pattern representing a signed integer

- **8 bit** 2's complement, negative input: -20

- Representation of x is unsigned binary translation of $2^n+x$

- Example: x=-20, then representation is $2^8+(-20)$ = 256 - 20

- There are many ways to calculate this.

# To get the 2's complement bit pattern representing a signed integer

- **8 bit** 2's complement, negative input: -20

- Method 1: Do the subtraction 256-20 in decimal. Convert result to binary.

- 256 - 20 = 236

- $236_{10} = 11101100_2$

- Result: 1110 1100

# To get the 2's complement bit pattern representing a signed integer

- **8 bit** 2's complement, negative input: -20

- Method 1: Do the subtraction 128-20 in decimal, then convert result to binary. Pad with zeroes.

- Problem: can leave you with a (very) large number to convert to binary.

# To get the 2's complement bit pattern representing a signed integer

- **8 bit** 2's complement, negative input: -20

- Method 2: Convert 20 to binary. Do the subtraction 256-20 in binary.

- $20_{10} = 10100_2$

- Result: 1 0000 0000 - 1 0100

# To get the 2's complement bit pattern representing a signed integer

- **8 bit** 2's complement, negative input: -20

- Method 2: Convert 20 to binary. Do the subtraction 256-20 in binary.

```
          1 0000 0000
 _____

               1 0100    -
          _____
 _____
```

# To get the 2's complement bit pattern representing a signed integer

First two digits OK

1 0000 0000

1 0100    -

00

# To get the 2's complement bit pattern representing a signed integer

Next needs borrow

```
          1 0000 0000
             1111 1100        Borrows
             1 0100      -
          _____
                  100
          _____
```

# To get the 2's complement bit pattern representing a signed integer

From then on, OK
Effectively flip the bit.

```
  1 0000 0000
    1111 1100        Borrows
      1 0100      -
  _____
    1110 1100
  _____
```

# To get the 2's complement bit pattern representing a signed integer

```
    1 0000 0000
      1111 1100        Borrows
       1 0100    -
      ─────────
      1110 1100
```

Result is: 1110 1100

# To get 2's complement representation of a signed integer

- If you look at the last algorithm, you see that you can write it as:

  - translate to binary

  - do something complicated with the bits (scan from right till first one, leave that unchanged, then flip all the bits to the left…

  - some of you may have been taught this

# To get 2's complement representation of a signed integer

A better method:

**To compute -n where n is in 2's complement binary:**

**1. flip the bits of n (ie change 0's to 1's and vice versa)**

Example: take 6 in 8-bit: 0000 0110 goes to 1111 1001

(Note: if we add these two together we get 1111 1111, so in unsigned terms this is equivalent to computing 255 - n = 256 - 1 - n = 256 - n -1. In general if we are using k bits it is $2^k$-n-1).

# To get 2's complement representation of a signed integer

**To compute -n where n is in 2's complement binary:**

**1. flip the bits of n (ie change 0's to 1's and vice versa)**

**2. add 1**

Example: take 6 in 8-bit: 0000 0110 goes to

1. 1111 1001

2. 1111 1010

(Note: in stage 1 we computed 256 - n -1, so this second stage gives 256 -n. The last 8 bits of this are the 2's complement representation of -n. In general if we are using k bits, we get $2^k$-n).

# To get 2's complement representation of a signed integer

**To compute -n where n is in 2's complement binary:**

**1. flip the bits of n (ie change 0's to 1's and vice versa)**

**2. add 1**

Note: both of these operations are basic binary operations that the cpu will have as standard instructions.

# Translation to and from 2's complement

- There are other ways of doing this.

- You will NOT be penalised for using them PROVIDED you explain your method clearly.

# Examples: signed

- 32 bit: 4 8-byte blocks, numbers from $-2^{31}$ to $2^{31}-1$

- Standard 2's complement representation

- 00000000 00000000 00000000 00000101 = 5

- 11111111 11111111 11111111 11111011 = -5

# Floating Point

# Floating point

- Computers use the IEEE floating point standard to represent real numbers (real is in the mathematical sense of numbers that are not integers).

- This is based on standard scientific notation, but using a normalised form, and expressed in binary.

- Most common is 64 bit and we will use that in examples.

# Scientific Notation

- Scientists write very large or very small numbers in **scientific notation.**

- Number is written as:

  - $2.9979 \times 10^{8}$ (speed of light)

  - or $6.626 \times 10^{-34}$ (Planck's constant)

# Scientific Notation

Parts of the number:

- 2.9979 x $10^8$ (speed of light)

**significand
or mantissa**

**exponent**

- or 6.626 x $10^{-34}$ (Planck's constant)

# Scientific Notation

Parts of the number:

- 2.9979 x $10^8$ (speed of light)

**significand
is between 1 and 10**

**exponent is either
positive or negative**

- or 6.626 x $10^{-34}$ (Planck's constant)

# Scientific notation

$$-1.23456 \times 10^{12}$$

Sign

Exponent

Significand
(formerly mantissa)

# Multiplying two numbers in scientific notation

$$(2.4 \times 10^{20}) * (1.5 \times 10^{-5})$$

Multiply significands          Add exponents

$$= 3.6 \times 10^{15}$$

# Adding two numbers in scientific notation

$(2.4 \times 10^{20}) + (1.2 \times 10^{19})$
$= 2.52 \times 10^{20}$

Translate to same exponent:
$1.2 \times 10^{19} = 0.12 \times 10^{20}$

Add significands:
$2.4 + 0.12 = 2.52$

# Computers

- Use a similar scheme to represent real numbers.

- It is called floating point.

- It uses a fixed size representation.

- But first we think about fixed size decimal scientific notation.

# Fixed width decimal

- Use a similar scheme to represent real numbers.

- It is called floating point.

- It uses a fixed size representation (we refer to it as **fixed width**).

- But first we think about fixed width decimal scientific notation.

# Fixed width decimal

$$-1.23456 \times 10^{12}$$

Sign

Significand:
limited number of
significant figures

Exponent:
limited range

# Fixed width decimal: example

$$-1.23456 \times 10^{12}$$

Sign

Significand:
limited to 6
significant figures

Exponent:
limited to -99 … 99

# Problems with fixed width representations

1. You can't represent every number accurately.

- Some numbers are too big: $10^{102}$

- Some numbers are just too small: $10^{-102}$

- Some numbers have too many decimal places: 10/3 = 3.33333333….

- All of these numbers have to be approximated.

- This introduces **rounding errors.**

# Rounding errors

- Given a number we **round** it to the nearest number we can represent.

  - $10^{-102}$ is too small to represent and would be rounded to 0.

  - $10/3 = 3.33333333\ldots$ would be rounded to $3.33333 \times 10^{0}$

  - $10^{102}$ would probably not be rounded (it is 100 times larger than the closest number we can represent). We would just say we had a number too big to represent.

  - The difference between the accurate number and the nearest representable number is the **rounding error**.

# Problems with fixed width representations

2. You get rounding errors when you add or multiply two representable numbers together.

- $1 \times 10^7 + 4 \times 10^0 = 1.000004 \times 10^7$ which gets rounded to $1.00000 \times 10^7$

- $1.00001 \times 10^0 * 1.1 \times 10^0 = 1.100011 \times 10^0$ which gets rounded to $1.10001 \times 10^0$

# Problems with fixed width representations

3. If you use two different mathematically equivalent methods to compute a value, you may get slightly different results.

- it is easy (and sobering) to run a test and find out for how many whole numbers N, N*(1/N) is *not* equal to 1.

# Problems with fixed width floating-point

- Rounding: you can't represent numbers accurately

- Operations: even when two numbers are representable accurately, it is usually not the case that their sum/difference/ product/quotient is. So operations of addition, subtraction, multiplication, division introduce further errors.

- Because of rounding, if you try to compute a number in two different ways, you will often get different results (technically this is statistically usually, but not always for simple examples).

- Therefore you should NEVER use equality testing on floating point reals.

# Computers

- use standard forms of fixed with reals

- very similar to the fixed width decimal we just saw

- formalised by the IEEE.

# IEEE binary floating point

- The way this should be implemented is set out in IEEE standards (IEEE=Institute of Electrical and Electronics Engineers).

- The relevant one is IEEE 754 as revised in 2008: IEEE 754-2008, on QMPlus.

# Basic structure ( IEEE p8)

— Signed zero and non-zero floating-point numbers of the form $(-1)^s \times b^e \times m$, where

  — $s$ is 0 or 1.

  — $e$ is any integer $emin \leq e \leq emax..$

  — $m$ is a number represented by a digit string of the form

  $d_0 \bullet d_1 d_2 ... d_{p-1}$ where $d_i$ is an integer digit $0 \leq d_i < b$ (therefore $0 \leq m < b$).

**Table 3.2—Parameters defining basic format floating-point numbers**

| parameter | Binary format $(b=2)$ | | | Decimal format $(b=10)$ | |
|---|---|---|---|---|---|
| | binary32 | binary64 | binary128 | decimal64 | decimal128 |
| $p$, digits | 24 | 53 | 113 | 16 | 34 |
| $emax$ | +127 | +1023 | +16383 | +384 | +6144 |

$emin$ shall be $1-emax$ for all formats.

# 64-bit reals

| 1 bit | MSB | $w$ bits | LSB | MSB | $t = p-1$ bits | LSB |
|-------|-----|----------|-----|-----|----------------|-----|
| $S$ (sign) | | $E$ (biased exponent) | | | $T$ (trailing significand field) | |

$E_0$.....................$E_{w-1}$ $d_1$..........................................................................................................$d_{p-1}$

**Figure 3.1—Binary interchange floating-point format**

- 1 bit for the sign

- 11 bits for the exponent

- 53 bit precision

- This makes 65 bits.

- Lose one bit by using normalised form (if you know the leading bit is 1, you don't need to record it).

# 64-bit reals: sign

| 1 bit | MSB   $w$ bits   LSB | MSB   $t = p-1$ bits   LSB |
|-------|----------------------|----------------------------|
| S (sign) | E (biased exponent) | T (trailing significand field) |

$E_0$..................$E_{w-1}$  $d_1$.........................................................$d_{p-1}$

**Figure 3.1—Binary interchange floating-point format**

- 1 bit for the sign

- 0 is positive

- 1 is negative

- This is because $(-1)^0 = 1$, and $(-1)^1 = -1$

# 64-bit reals: exponent

| 1 bit | MSB | w bits | LSB MSB | t = p − 1 bits | LSB |
|---|---|---|---|---|---|
| S (sign) | | E (biased exponent) | | T (trailing significand field) | |
| | $E_0$................$E_{w-1}$ | | $d_1$................................................................$d_{p-1}$ | | |

**Figure 3.1—Binary interchange floating-point format**

- 11 bits for the expon

- emax is 1023

- emin is 1 - emax = 1 - 1023 = -1022

- biased means "start at emin"

- 100 0000 0000 = 1024 represents 1

- 011 1111 1111 = 1023 represents 0

- 100 0000 0001 = 1025 represents 2

# 64-bit reals: exponent



| 1 bit | MSB | $w$ bits | LSB MSB | $t = p-1$ bits | LSB |
|---|---|---|---|---|---|
| S (sign) | | E (biased exponent) | | T (trailing significand field) | |

$E_0$...................$E_{w-1}$  $d_1$...........................................................................................$d_{p-1}$

**Figure 3.1—Binary interchange floating-point format**

- 11 bits for the exponent

- emax is 1023

- emin is 1 - emax = 1 - 1023 = -1022

- biased means "start at emin"

- 100 0000 0000 = 1024 represents 1

- 011 1111 1111 = 1023 represents 0

- 100 0000 0001 = 1025 represents 2

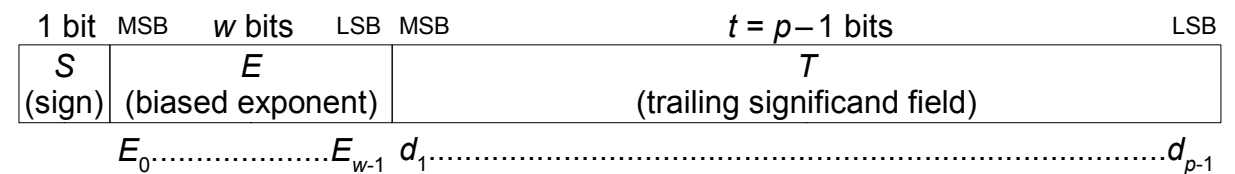- So we get exponent by subtracting 1023 from the unsigned value.

# 64-bit reals: exponent

| 1 bit | MSB        $w$ bits        LSB | MSB                    $t = p - 1$ bits                    LSB |
|-------|-------------------------------|--------------------------------------------------------------|
| $S$ (sign) | $E$ (biased exponent) | $T$ (trailing significand field) |

$E_0$···················$E_{w-1}$  $d_1$·····················································································$d_{p-1}$

**Figure 3.1—Binary interchange floating-point format**

- 11 bits for the exponent

- emax is 1023

- emin is 1 - emax = 1 - 1023 = -1022

- biased means "start at emin"

- We get exponent by subtracting 1023 from the unsigned value.

- This means 000 0000 0000 = 0 represents exponent 0 -1023 = -1023, which is out of range.

- And 111 1111 1111 = 2047 represents exponent 2047 -1023 = 1024, which is also out of range.

- That gives us some room to represent other things: e.g. 0, overflow, infinities, NaN's

# 64-bit reals: exponent

| 1 bit | MSB | w bits | LSB | MSB | t = p−1 bits | LSB |
|---|---|---|---|---|---|---|
| S (sign) | | E (biased exponent) | | | T (trailing significand field) | |

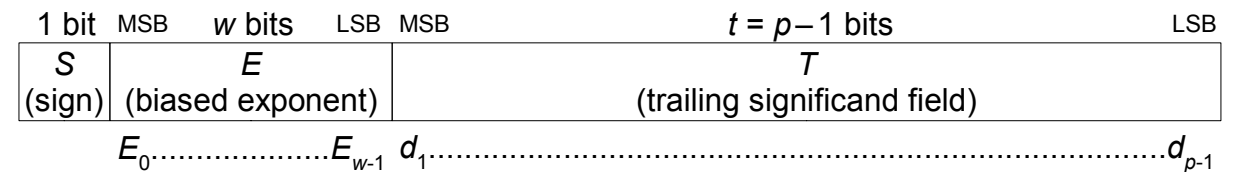$E_0$...................$E_{w-1}$  $d_1$.....................................................................$d_{p-1}$

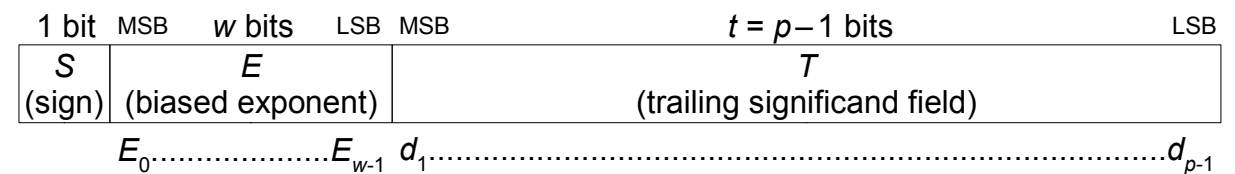**Figure 3.1—Binary interchange floating-point format**

- 11 bits for the exponent

- emax is 1023

- emin is 1 - emax = 1 - 1023 = -1022

- biased means "start at emin"

- This means 000 0000 0000 = 0 represents exponent 0 -1023 = -1023, which is out of range.

- Note that in our experiment 0 was represented by 0000 0000 0000 0000 …. 0000

- This uses that out of range exponent.

Within each format, the following floating-point data shall be represented:

— Signed zero and non-zero floating-point numbers of the form $(-1)^s \times b^e \times m$, where

    — $s$ is 0 or 1.

    — $e$ is any integer $emin \leq e \leq emax..$

    — $m$ is a number represented by a digit string of the form

        $d_0 \bullet d_1 d_2 ... d_{p-1}$ where $d_i$ is an integer digit $0 \leq d_i < b$ (therefore $0 \leq m < b$).

— Two infinities, $+\infty$ and $-\infty$.

— Two NaNs, qNaN (quiet) and sNaN (signaling).

These are the only floating-point data represented.

# 64-bit reals: significand

| 1 bit | MSB | $w$ bits | LSB | MSB | $t = p-1$ bits | LSB |
|-------|-----|----------|-----|-----|----------------|-----|
| S (sign) | | E (biased exponent) | | | T (trailing significand field) | |

$E_0$.....................$E_{w-1}$ $d_1$.........................................................................$d_{p-1}$

**Figure 3.1—Binary interchange floating-point format**

- 52 bits for the significand.

- Unless a binary number is 0, its most significant digit will be 1.

- This contrasts with decimal, where the most significant digit can be any of 1..9.

- This means that a normalised binary real in scientific notation will always look like:

- +/- $1.d_1d_2d_3\ldots * 2^e$

- In IEEE floating point, the significand bits give us $d_1d_2d_3\ldots$

# 64-bit reals: significand

| 1 bit | MSB | w bits | LSB | MSB | t = p − 1 bits | LSB |
|---|---|---|---|---|---|---|
| S (sign) | | E (biased exponent) | | | T (trailing significand field) | |

$E_0$..................$E_{w-1}$  $d_1$..................................................................................$d_{p-1}$

**Figure 3.1—Binary interchange floating-point format**

- 52 bits for the significand.

- This means that a normalised binary real in scientific notation will always look like:

- +/- $1.d_1d_2d_3\ldots * 2^e$

- In IEEE floating point, the significand bits give us $d_1d_2d_3\ldots$

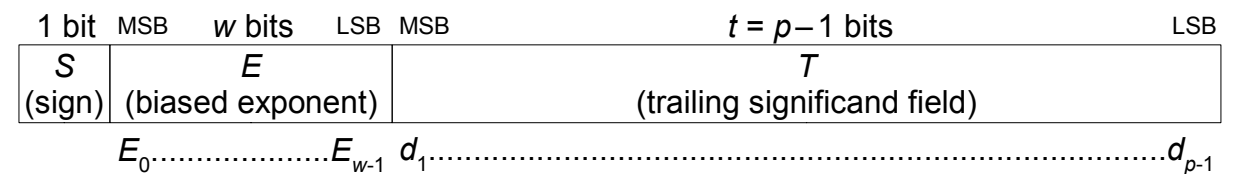- So 0000 0000 … 0000 means we have 1.000… as our significand (hence a power of 2: 1/4,1/2,1,2,4,…)

- Similarly 0100 0000 … 0000 means we have 1.01 times a power of 2, and 1.01 is 5/4 = 1.25 ($101_2 = 5_{10}$, and shift two places corresponds to division by $2^2 = 4$).

# Big Endian and Little Endian

- Remember: data is organised into 8-bit bytes.

- There is a choice about which order to store the bytes in.

- We are used to seeing high-order bytes first. This is Big Endian.

- Intel and others use Little Endian (low-order bytes first).

- This is used for both integers and reals.

- It can be a little confusing when looking at bit dumps.

# Floating Point

- One theme of this course is that many of the techniques used in computing correspond very closely to things we use or do already in real life.

- The computer representation of reals corresponds to scientific notation.

- More specifically it corresponds to scientific notation to a fixed number of places, and with a fixed range of possible powers of 10.

# Floating Point: additional material

# Compare what we have just seen with

- Speed of light: $2.997930 * 10^8$ ms$^{-1}$

- Mass of hydrogen atom: $1.67 \times 10^{-24}$ g

- Avogadro's constant: $6.02214179 \times 10^{23}$

# Rounding 1

- If we take an arbitrary real number, then it probably isn't one of our floating-point reals.

- Example: the closest we get to 1/3 is 3.33 e -1.

- This is out by 0.00033333....

- So we get errors when we put numbers into the system.

- These errors can build up as the program runs.

# Rounding 1

- and we certainly can't represent numbers like pi, e, or the square root of 2 exactly.

# Rounding 1

- Moreover, how accurately we can represent a number depends on how big it is.

- Near 0 we can get within $10^{-9}$ of a number. Smallest non-zero number we can represent is $1.00*10^9$ if we only allow significands beginning with 1-9. It's $0.01*10^{-9} = 1.00*10^{-11}$ if we allow ones beginning with 0.

- That's 0.000000001 in the first case and 0.00000000001 in the second.

# Rounding 1

- Up near 1,000,000 we have:

- $1,000,000 = 1.00*10^6$

- so the next number we can represent is $1.01*10^6$

- $1.01*10^6 = 1,010,000$

- So the gap between numbers we can represent has jumped to 10,000.

- In other words, we cannot represent 1,005,000 accurately.

# Class question

- What is the smallest whole number we can't represent?

- Ans: 1001

- It would be $1.001*10^3$ which requires 4 digits in the exponent.

# Rounding 2

- We basically use two algorithms, one for addition and one for subtraction.

- Let's concentrate on addition.

- ASSUME: significand is between 1.000… and 9.999…

- IE it has a single digit to the left of the decimal point, and that digit is not 0.

# Addition

- To add $x*10^a$ and $y*10^{b:}$

- Take the greater of a and b (without loss of generality, suppose it's a).

- Shift y right by a-b, to get y' $(y*10^b = y'*10^a)$

- Add x and y' essentially as integers to get z

- Round z to the correct number of places to get z'

- If the result is less than 10, answer is $z'*10^a$

- If the result is bigger than 10, shift right by 1 to get z", answer is z" $*10^{a+1}$.

# Addition (Example)

- To add $x*10^a = 2.01*10^3$ and $y*10^b = 3.24*10^{2:}$

- Take the greater of a and b (a=3).

- Shift y right by a-b=1, to get y'=0.324  ($y*10^b = y'*10^a$)

- Add x and y' essentially as integers (2010 and 0324) to get z=2.334

- Round z to the correct number of places to get z'=2.33

- If the result is less than 10, answer is $z'*10^a = 2.33*10^3$

# Addition (Example)

- To add $x*10^a = 9.80 * 10^4$ and $y*10^b = 4.67 * 10^3$

- Take the greater of a and b (a=4).

- Shift y right by a-b=1, to get y'=0.467  ($y*10^b = y'*10^a$)

- Add x and y' essentially as integers to get z = 10.267

- Round z to the correct number of places to get z' = 10.3

- If the result is less than 10, answer is $z'*10^a$

- If the result is bigger than 10, shift right by 1 to get z"= 1.03, answer is z" * $10^{a+1}$ = 1.03 * $10^5$.

# Rounding 2

- When we add, multiply or subtract two floating-point numbers we don't necessarily get another one.

- Example: 1.23e1 + 4.56e-1

  $$= 1.23e1 + 0.0456\ e1$$

  $$= 1.2756e1$$

  - Now we have to decide which number to use as the value of the sum.: it could be either 1.27e1 or 1.28e1

  - We could decide to: round to $+\infty$, round to $-\infty$, round to nearest, round toward zero, round away from 0.

  - In fact, IEEE does not allow round away from 0, and there are two variants of round to nearest: round ties to even and round ties away from 0.

  - These errors can build up too.

# Rounding 2

- Just to ram this point home: take a couple of extreme examples.

- $1.00 + 0.001 = 1.0\ e0 + 1.0e-3 = 1.0\ e0$ !!!

- $1000 + 1 = 1.0e3 + 1.0e0 = 1.0\ e3 = 1000$ !!!

- This kind of thing really happens on your computer (not with these numbers), and we will be doing some experiments to see more detail.

- Notice that we are getting this problem because the numbers are apart by a factor of $10^3$, which corresponds to the number of significant figures we have in the representation.

# Experiment week4.6.c

- Live experiment

# Another test: week4.7.c

- Actually, we don't need the constants to show this. We can find values.

```c
#include <stdio.h>

main()
{
  double a = 1.0;
  double e = 1.0;
  int n = 0;
  while ((a+e)!=a) {
    e = e/2;
    n++;
    }
  printf("For a = %e: n = %i; e = %e\n",a,n,e);
}
```

```
Edmunds-MacBook:C-code edmundr$ ./week4.7
For a = 1.000000e+00: n = 53; e = 1.110223e-16
Edmunds-MacBook:C-code edmundr$
```

# Operations

**5.4.1 Arithmetic operations**

Implementations shall provide the following *formatOf* general-computational operations, for destinations of all supported arithmetic formats, and, for each destination format, for operands of all supported arithmetic formats with the same radix as the destination format. These operations shall not propagate non-canonical results.

- *formatOf-**addition**(source1, source2)*

  The operation **addition**$(x, y)$ computes $x+y$.

  The preferred exponent is $\min(Q(x), Q(y))$.

- *formatOf-**subtraction**(source1, source2)*

  The operation **subtraction**$(x, y)$ computes $x-y$.

  The preferred exponent is $\min(Q(x), Q(y))$.

- *formatOf-**multiplication**(source1, source2)*

  The operation **multiplication**$(x, y)$ computes $x \times y$.

  The preferred exponent is $Q(x)+Q(y)$.

- *formatOf-**division**(source1, source2)*

  The operation **division**$(x, y)$ computes $x/y$.

  The preferred exponent is $Q(x) - Q(y)$.

- *formatOf-**squareRoot**(source1)*

  The operation **squareRoot**$(x)$ computes $\sqrt{x}$. It has a positive sign for all operands $\geq 0$, except that **squareRoot**$(-0)$ shall be $-0$.

  The preferred exponent is $\mathrm{floor}(Q(x)/2)$.

- *formatOf-**fusedMultiplyAdd**(source1, source2, source3)*

  The operation **fusedMultiplyAdd**$(x, y, z)$ computes $(x \times y)+z$ as if with unbounded range and precision, rounding only once to the destination format. No underflow, overflow, or inexact exception (see 7) can arise due to the multiplication, but only due to the addition; and so fusedMultiplyAdd differs from a multiplication operation followed by an addition operation.

  The preferred exponent is $\min(Q(x)+Q(y), Q(z))$.

- *formatOf-**convertFromInt**(int)*

  It shall be possible to convert from all supported signed and unsigned integer formats to all supported arithmetic formats. Integral values are converted exactly from integer formats to floating-point formats whenever the value is representable in both formats. If the converted value is not exactly representable in the destination format, the result is determined according to the applicable rounding-direction attribute, and an inexact or floating-point overflow exception arises as specified in Clause 7, just as with arithmetic operations. The signs of integer zeros are preserved. Integer zeros without signs are converted to $+0$.

  The preferred exponent is 0.

- The standard also specifies arithmetical operations...

# More on rounding

- The inexactness in floating-point means that you need to be careful how you do things.

- There are some basic rules:

  - never base a test on whether two floating-point numbers are equal (see experiment)

  - never simply use raw values (almost always better to work with a base and offset, x+e).

- For more (much more) see "What every computer scientist should know about floating-point arithmetic"

# Reading

- What every computer scientist should know about floating-point arithmetic.

- IEEE 754-2008

# Floating Point: end of additional material

# Character Sets

# Text: the simplest example

- Computers need to deal with text.

- Text is made up of individual characters.

- Each character is represented as a number.

- Exactly how it is represented as a number depends on the encoding and the "character set" being used.

# From a mail message header

```
From: <concurrency-request@listserver.tue.nl>
Subject: Concurrency Digest, Vol 8, Issue 90
To: <concurrency@listserver.tue.nl>
Reply-To: <concurrency@listserver.tue.nl>
Date: Sat, 14 Jun 2014 19:54:09 +0200
Message-ID: <mailman.11.1402768449.7567.concurrency@listserver.tue.nl>
Content-Type: text/plain; charset="us-ascii"
Content-Transfer-Encoding: 7bit
X-BeenThere: concurrency@listserver.tue.nl
X-Mailman-Version: 2.1.12
```

# Another

```
X-MS-Exchange-Organization-AuthMechanism: 03
X-MS-Exchange-Organization-AuthSource: DB4PR07MB332.eurprd07.prod.outlook.com
X-MS-Has-Attach:
X-MS-Exchange-Organization-SCL: -1
X-MS-TNEF-Correlator:
Content-Type: text/plain; charset="iso-8859-2"
Content-Transfer-Encoding: quoted-printable
MIME-Version: 1.0

Dear Bill,

In the "strong links with business" part of the letter I'd add "Philips, Ho=
neywell and BBC".
```

# QMPlus Web Page

```html
<!DOCTYPE html>
<html  dir="ltr" lang="en" xml:lang="en">
<head>
    <title>Course: ECS404U - Computer Systems and Networks - 2014/15</title>
    <link rel="shortcut icon" href="http://qmplus.qmul.ac.uk/theme/image.php/qmul science/the
    <meta http-equiv="Content-Type" content="text/html; charset=utf-8" /><script type="text/j
<meta name="keywords" content="moodle, Course: ECS404U - Computer Systems and Networks - 2014
<meta http-equiv="pragma" content="no-cache" />
<meta http-equiv="expires" content="0" />
<script type="text/javascript">
```

# Characters: ASCII

- The oldest and still the most famous representation is ASCII (American Standard Code for Information Interchange).

- It uses 7 bits, and so numbers $0\ldots2^7-1 = 127$

- It is based on the idea of an old-style line-printer.

| dec | hex | chr | | dec | hex | chr | | dec | hex | chr | dec | hex | chr | dec | hex | chr | dec | hex | chr | dec | hex | chr | dec | hex | chr |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 000d | 00h | ` | (nul) | 016d | 10h | ► | (dle) | 032d | 20h | ␣ | 048d | 30h | 0 | 064d | 40h | @ | 080d | 50h | P | 096d | 60h | ' | 112d | 70h | p |
| 001d | 01h | ☺ | (soh) | 017d | 11h | ◄ | (dc1) | 033d | 21h | ! | 049d | 31h | 1 | 065d | 41h | A | 081d | 51h | Q | 097d | 61h | a | 113d | 71h | q |
| 002d | 02h | ● | (stx) | 018d | 12h | ↕ | (dc2) | 034d | 22h | " | 050d | 32h | 2 | 066d | 42h | B | 082d | 52h | R | 098d | 62h | b | 114d | 72h | r |
| 003d | 03h | ♥ | (etx) | 019d | 13h | ‼ | (dc3) | 035d | 23h | # | 051d | 33h | 3 | 067d | 43h | C | 083d | 53h | S | 099d | 63h | c | 115d | 73h | s |
| 004d | 04h | ♦ | (eot) | 020d | 14h | ¶ | (dc4) | 036d | 24h | $ | 052d | 34h | 4 | 068d | 44h | D | 084d | 54h | T | 100d | 64h | d | 116d | 74h | t |
| 005d | 05h | ♠ | (enq) | 021d | 15h | § | (nak) | 037d | 25h | % | 053d | 35h | 5 | 069d | 45h | E | 085d | 55h | U | 101d | 65h | e | 117d | 75h | u |
| 006d | 06h | ♠ | (ack) | 022d | 16h | ▬ | (syn) | 038d | 26h | & | 054d | 36h | 6 | 070d | 46h | F | 086d | 56h | V | 102d | 66h | f | 118d | 76h | v |
| 007d | 07h | · | (bel) | 023d | 17h | ↨ | (etb) | 039d | 27h | ' | 055d | 37h | 7 | 071d | 47h | G | 087d | 57h | W | 103d | 67h | g | 119d | 77h | w |
| 008d | 08h | ▫ | (bs) | 024d | 18h | ↑ | (can) | 040d | 28h | ( | 056d | 38h | 8 | 072d | 48h | H | 088d | 58h | X | 104d | 68h | h | 120d | 78h | x |
| 009d | 09h | | (tab) | 025d | 19h | ↓ | (em) | 041d | 29h | ) | 057d | 39h | 9 | 073d | 49h | I | 089d | 59h | Y | 105d | 69h | i | 121d | 79h | y |
| 010d | 0Ah | ▪ | (lf) | 026d | 1Ah | | (eof) | 042d | 2Ah | * | 058d | 3Ah | : | 074d | 4Ah | J | 090d | 5Ah | Z | 106d | 6Ah | j | 122d | 7Ah | z |
| 011d | 0Bh | ♂ | (vt) | 027d | 1Bh | ← | (esc) | 043d | 2Bh | + | 059d | 3Bh | ; | 075d | 4Bh | K | 091d | 5Bh | [ | 107d | 6Bh | k | 123d | 7Bh | { |
| 012d | 0Ch | | (np) | 028d | 1Ch | ∟ | (fs) | 044d | 2Ch | ' | 060d | 3Ch | < | 076d | 4Ch | L | 092d | 5Ch | \ | 108d | 6Ch | l | 124d | 7Ch | | |
| 013d | 0Dh | ♪ | (cr) | 029d | 1Dh | ↔ | (gs) | 045d | 2Dh | - | 061d | 3Dh | = | 077d | 4Dh | M | 093d | 5Dh | ] | 109d | 6Dh | m | 125d | 7Dh | } |
| 014d | 0Eh | ♫ | (so) | 030d | 1Eh | ▲ | (rs) | 046d | 2Eh | . | 062d | 3Eh | > | 078d | 4Eh | N | 094d | 5Eh | ^ | 110d | 6Eh | n | 126d | 7Eh | ~ |
| 015d | 0Fh | ☼ | (si) | 031d | 1Fh | ▼ | (us) | 047d | 2Fh | / | 063d | 3Fh | ? | 079d | 4Fh | O | 095d | 5Fh | _ | 111d | 6Fh | o | 127d | 7Fh | ⌂ |

# Some things to notice

- 7 bits: abcdefg

- A-Z occupy: 65-90: 1000001-1011010

- a-z occupy: 97-122: 1100001-1111010

- they differ by exactly 32 (hence in one bit), and capitals start at binary 100000.

- 0-9 occupies: 48-57 (binary 0110000 - 0111001)

# Some things to notice

- There's a whole bunch of strange stuff, eg carriage return (moves the print head back to the start of the line) and line-feed (moves the paper up a line).

- There are no pound signs (only dollar=36), and no funny accents.

# More modern character sets

- Computers work in units of $2^n$ bits. So more modern character sets would have 8 bits not 7, and 256 characters (0..255) not 128 (0..127).

- Most modern character sets are standardised by ISO (International Standards Organisation).

# ISO-8859-1

- Extension of ASCII (to 8 bits)

- (Was) default character set in many browsers (now UTF-8)

- One of a number of ISO character sets designed to allow 8-bit encoding of other national alphabets.

- This one is Latin, and includes most European symbols (eg accents).

# ISO-8859-1 (Latin)

| Char | Code | Name | Description |
|------|------|------|-------------|
| à | 224 | agrave | a grave |
| á | 225 | aacute | a acute |
| â | 226 | acirc | a circumflex |
| ã | 227 | atilde | a tilde |
| ä | 228 | auml | a umlaut |
| å | 229 | aring | a ring |
| æ | 230 | aelig | ae ligature |
| ç | 231 | ccedil | c cedilla |
| è | 232 | egrave | e grave |
| é | 233 | eacute | e acute |
| ê | 234 | ecirc | e circumflex |
| ë | 235 | euml | e umlaut |
| ì | 236 | igrave | i grave |
| í | 237 | iacute | i acute |
| î | 238 | icirc | i circumflex |
| ï | 239 | iuml | i umlaut |

| Char | Code | Name | Description |
|------|------|------|-------------|
| ð | 240 | eth | eth |
| ñ | 241 | ntilde | n tilde |
| ò | 242 | ograve | o grave |
| ó | 243 | oacute | o acute |
| ô | 244 | ocirc | o circumflex |
| õ | 245 | otilde | o tilde |
| ö | 246 | ouml | o umlaut |
| ÷ | 247 | divide | division sign |
| ø | 248 | oslash | o slash |
| ù | 249 | ugrave | u grave |
| ú | 250 | uacute | u acute |
| û | 251 | ucirc | u circumflex |
| ü | 252 | uuml | u umlaut |
| ý | 253 | yacute | y acute |
| þ | 254 | thorn | thorn |
| ÿ | 255 | yuml | y umlaut |

| Char | Code | Name | Description | Char | Code | Name | Description | Char | Code | Name | Description | Char | Code | Name | Description |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 32 | - | Normal space | 0 | 48 | - | Digit 0 | @ | 64 | - | At sign | P | 80 | - | P |
| ! | 33 | - | Exclamation | 1 | 49 | - | Digit 1 | A | 65 | - | A | Q | 81 | - | Q |
| " | 34 | quot | Double quote | 2 | 50 | - | Digit 2 | B | 66 | - | B | R | 82 | - | R |
| # | 35 | - | Hash or pound | 3 | 51 | - | Digit 3 | C | 67 | - | C | S | 83 | - | S |
| $ | 36 | - | Dollar | 4 | 52 | - | Digit 4 | D | 68 | - | D | T | 84 | - | T |
| % | 37 | - | Percent | 5 | 53 | - | Digit 5 | E | 69 | - | E | U | 85 | - | U |
| & | 38 | - | Ampersand | 6 | 54 | - | Digit 6 | F | 70 | - | F | V | 86 | - | V |
| ' | 39 | - | Apostrophe | 7 | 55 | - | Digit 7 | G | 71 | - | G | W | 87 | - | W |
| ( | 40 | - | Open bracket | 8 | 56 | - | Digit 8 | H | 72 | - | H | X | 88 | - | X |
| ) | 41 | - | Close bracket | 9 | 57 | - | Digit 9 | I | 73 | - | I | Y | 89 | - | Y |
| * | 42 | - | Asterisk | : | 58 | - | Colon | J | 74 | - | J | Z | 90 | - | Z |
| + | 43 | - | Plus sign | ; | 59 | - | Semicolon | K | 75 | - | K | [ | 91 | - | Open square bracket |
| , | 44 | - | Comma | < | 60 | lt | Less than | L | 76 | - | L | \ | 92 | - | Backslash |
| - | 45 | - | Minus sign | = | 61 | - | Equals | M | 77 | - | M | ] | 93 | - | Close square bracket |
| . | 46 | - | Period | > | 62 | gt | Greater than | N | 78 | - | N | ^ | 94 | - | Pointer |
| / | 47 | - | Forward slash | ? | 63 | - | Question mark | O | 79 | - | O | _ | 95 | - | Underscore |

| Char | Code | Name | Description | Char | Code | Name | Description | Char | Code | Name | Description | Char | Code | Name | Description |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ` | 96 | - | Grave accent | p | 112 | - | p |   | 160 | nbsp | Non-breaking space | ° | 176 | deg | Degree sign |
| a | 97 | - | a | q | 113 | - | q | ¡ | 161 | iexcl | Inverted exclamation | ± | 177 | plusmn | Plus-minus sign |
| b | 98 | - | b | r | 114 | - | r | ¢ | 162 | cent | Cent sign | ² | 178 | sup2 | Superscript 2 |
| c | 99 | - | c | s | 115 | - | s | £ | 163 | pound | Pound sign | ³ | 179 | sup3 | Superscript 3 |
| d | 100 | - | d | t | 116 | - | t | ¤ | 164 | curren | Currency sign | ´ | 180 | acute | Spacing acute |
| e | 101 | - | e | u | 117 | - | u | ¥ | 165 | yen | Yen sign | µ | 181 | micro | Micro sign |
| f | 102 | - | f | v | 118 | - | v | ¦ | 166 | brvbar | Broken bar | ¶ | 182 | para | Paragraph sign |
| g | 103 | - | g | w | 119 | - | w | § | 167 | sect | Section sign | · | 183 | middot | Middle dot |
| h | 104 | - | h | x | 120 | - | x | ¨ | 168 | uml | Umlaut or diaeresis | ¸ | 184 | cedil | Spacing cedilla |
| i | 105 | - | i | y | 121 | - | y | © | 169 | copy | Copyright sign | ¹ | 185 | sup1 | Superscript 1 |
| j | 106 | - | j | z | 122 | - | z | ª | 170 | ordf | Feminine ordinal | º | 186 | ordm | Masculine ordinal |
| k | 107 | - | k | { | 123 | - | Left brace | « | 171 | laquo | Left angle quotes | » | 187 | raquo | Right angle quotes |
| l | 108 | - | l | \| | 124 | - | Vertical bar | ¬ | 172 | not | Logical not sign | ¼ | 188 | frac14 | One quarter |
| m | 109 | - | m | } | 125 | - | Right brace | | 173 | shy | Soft hyphen | ½ | 189 | frac12 | One half |
| n | 110 | - | n | ~ | 126 | - | Tilde | ® | 174 | reg | Registered trademark | ¾ | 190 | frac34 | Three quarters |
| o | 111 | - | o | ✗ | 127 | - | (Unused) | ¯ | 175 | macr | Spacing macron | ¿ | 191 | iquest | Inverted question mark |

| Char | Code | Name | Description | Char | Code | Name | Description | Char | Code | Name | Description | Char | Code | Name | Description |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| À | 192 | Agrave | A grave | Ð | 208 | ETH | ETH | à | 224 | agrave | a grave | ð | 240 | eth | eth |
| Á | 193 | Aacute | A acute | Ñ | 209 | Ntilde | N tilde | á | 225 | aacute | a acute | ñ | 241 | ntilde | n tilde |
| Â | 194 | Acirc | A circumflex | Ò | 210 | Ograve | O grave | â | 226 | acirc | a circumflex | ò | 242 | ograve | o grave |
| Ã | 195 | Atilde | A tilde | Ó | 211 | Oacute | O acute | ã | 227 | atilde | a tilde | ó | 243 | oacute | o acute |
| Ä | 196 | Auml | A umlaut | Ô | 212 | Ocirc | O circumflex | ä | 228 | auml | a umlaut | ô | 244 | ocirc | o circumflex |
| Å | 197 | Aring | A ring | Õ | 213 | Otilde | O tilde | å | 229 | aring | a ring | õ | 245 | otilde | o tilde |
| Æ | 198 | AElig | AE ligature | Ö | 214 | Ouml | O umlaut | æ | 230 | aelig | ae ligature | ö | 246 | ouml | o umlaut |
| Ç | 199 | Ccedil | C cedilla | × | 215 | times | Multiplication sign | ç | 231 | ccedil | c cedilla | ÷ | 247 | divide | division sign |
| È | 200 | Egrave | E grave | Ø | 216 | Oslash | O slash | è | 232 | egrave | e grave | ø | 248 | oslash | o slash |
| É | 201 | Eacute | E acute | Ù | 217 | Ugrave | U grave | é | 233 | eacute | e acute | ù | 249 | ugrave | u grave |
| Ê | 202 | Ecirc | E circumflex | Ú | 218 | Uacute | U acute | ê | 234 | ecirc | e circumflex | ú | 250 | uacute | u acute |
| Ë | 203 | Euml | E umlaut | Û | 219 | Ucirc | U circumflex | ë | 235 | euml | e umlaut | û | 251 | ucirc | u circumflex |
| Ì | 204 | Igrave | I grave | Ü | 220 | Uuml | U umlaut | ì | 236 | igrave | i grave | ü | 252 | uuml | u umlaut |
| Í | 205 | Iacute | I acute | Ý | 221 | Yacute | Y acute | í | 237 | iacute | i acute | ý | 253 | yacute | y acute |
| Î | 206 | Icirc | I circumflex | Þ | 222 | THORN | THORN | î | 238 | icirc | i circumflex | þ | 254 | thorn | thorn |
| Ï | 207 | Iuml | I umlaut | ß | 223 | szlig | sharp s | ï | 239 | iuml | i umlaut | ÿ | 255 | yuml | y umlaut |

| Character set | Description | Covers |
|---|---|---|
| ISO-8859-1 | Latin alphabet part 1 | North America, Western Europe, Latin America, the Caribbean, Canada, Africa |
| ISO-8859-2 | Latin alphabet part 2 | Eastern Europe |
| ISO-8859-3 | Latin alphabet part 3 | SE Europe, Esperanto, miscellaneous others |
| ISO-8859-4 | Latin alphabet part 4 | Scandinavia/Baltics (and others not in ISO-8859-1) |
| ISO-8859-5 | Latin/Cyrillic part 5 | The languages that are using a Cyrillic alphabet such as Bulgarian, Belarusian, Russian and Macedonian |
| ISO-8859-6 | Latin/Arabic part 6 | The languages that are using the Arabic alphabet |
| ISO-8859-7 | Latin/Greek part 7 | The modern Greek language as well as mathematical symbols derived from the Greek |
| ISO-8859-8 | Latin/Hebrew part 8 | The languages that are using the Hebrew alphabet |
| ISO-8859-9 | Latin 5 part 9 | The Turkish language. Same as ISO-8859-1 except Turkish characters replace Icelandic ones |
| ISO-8859-10 | Latin 6 Lappish, Nordic, Eskimo | The Nordic languages |
| ISO-8859-15 | Latin 9 (aka Latin 0) | Similar to ISO 8859-1 but replaces some less common symbols with the euro sign and some other missing characters |
| ISO-2022-JP | Latin/Japanese part 1 | The Japanese language |
| ISO-2022-JP-2 | Latin/Japanese part 2 | The Japanese language |
| ISO-2022-KR | Latin/Korean part 1 | The Korean language |

From w3schools.com

# But

- 256 characters is not enough...

- if you need to cover more than one language

- if you're a mathematician

- if you're Chinese

# And so we have unicode

- Character set used in Java

- First 128 characters are ASCII

- First 256 are ISO-8859-1

- Total number of characters available is: 1,114,112 = 21 (just over 20) bits

- Lots of segments correspond to particular languages

- See http://www.unicode.org/

# Unicode

- Distinguishes between the number of the character (the code point) and the way it is represented.

- So we have UTF-32 (represents character in 32 bits = 4 bytes)

- But also UTF-16 and UTF-8 (16 and 8 bits).

- UTF-8 is now standard on web applications.

- UTF-32 and UTF-16 have "big endian" and "little endian" variants, referring to the order of the bytes.

# Unicode encodings

**Figure 2-12.** Unicode Encoding Schemes

| | | | | |
|---|---|---|---|---|
| A | Ω | 語 | Ⅲ | UTF-32BE |
| 00 00 00 41 | 00 00 03 A9 | 00 00 8A 9E | 00 01 03 84 | |

| | | | | |
|---|---|---|---|---|
| A | Ω | 語 | Ⅲ | UTF-32LE |
| 41 00 00 00 | A9 03 00 00 | 9E 8A 00 00 | 84 03 01 00 | |

| | | | | |
|---|---|---|---|---|
| A | Ω | 語 | Ⅲ | UTF-16BE |
| 00 41 | 03 A9 | 8A 9E | D8 00 DF 84 | |

| | | | | |
|---|---|---|---|---|
| A | Ω | 語 | Ⅲ | UTF-16LE |
| 41 00 | A9 03 | 9E 8A | 00 D8 84 DF | |

| | | | | |
|---|---|---|---|---|
| A | Ω | 語 | Ⅲ | UTF-8 |
| 41 | CE A9 | E8 AA 9E | F0 90 8E 84 | |

# UTF-8

**Preferred Usage.** UTF-8 is typically the preferred encoding form for HTML and similar protocols, particularly for the Internet. The ASCII transparency helps migration. UTF-8 also has the advantage that it is already inherently byte-serialized, as for most existing 8-bit character sets; strings of UTF-8 work easily with C or other programming languages, and many existing APIs that work for typical Asian multibyte character sets adapt to UTF-8 as well with little or no change required.

From Unicode Standard v 6.2

# Translation

- We want you to use UTF-8 for internet and web-based protocols

- here are some reasons

# UTF-8

- But UTF-8 uses different numbers of bytes to encode individual characters.

- It therefore has to be cleverly and carefully designed so that it is not ambiguous, and so that as much existing software continues to work with it as possible.

- The designer (Ken Thompson) was one of the lead designers of Unix.

# From the rfc for UTF-8

UTF-8 encodes UCS characters as a varying number of octets, where the number of octets, and the value of each, depend on the integer value assigned to the character in ISO/IEC 10646 (the character number, a.k.a. code position, code point or Unicode scalar value).  This encoding form has the following characteristics (all values are in hexadecimal):

o  Character numbers from U+0000 to U+007F (US-ASCII repertoire) correspond to octets 00 to 7F (7 bit US-ASCII values).  A direct consequence is that a plain ASCII string is also a valid UTF-8 string.

# From the rfc for UTF-8

o US-ASCII octet values do not appear otherwise in a UTF-8 encoded
  character stream. This provides compatibility with file systems
  or other software (e.g., the printf() function in C libraries)
  that parse based on US-ASCII values but are transparent to other
  values.

o Round-trip conversion is easy between UTF-8 and other encoding
  forms.

o The first octet of a multi-octet sequence indicates the number of
  octets in the sequence.

o The octet values C0, C1, F5 to FF never appear.

o Character boundaries are easily found from anywhere in an octet
  stream.

o The byte-value lexicographic sorting order of UTF-8 strings is the
  same as if ordered by character numbers. Of course this is of
  limited interest since a sort order based on character numbers is
  almost never culturally valid.

o The Boyer-Moore fast search algorithm can be used with UTF-8 data.

o UTF-8 strings can be fairly reliably recognized as such by a
  simple algorithm, i.e., the probability that a string of
  characters in any other encoding appears as valid UTF-8 is low,
  diminishing with increasing string length.