

WEEK 2

- 10/4 - 1, 2, 3, 4, ... are representations of numbers ($\mathbb{R} - \mathbb{Q}$)
- base 10, 60, 2, 8, etc.
- | | |
|---|--|
| <ul style="list-style-type: none">- $\mathbb{N} = \{1, 2, 3, \dots\}$- $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$- $\mathbb{Q} = \left\{ \frac{a}{b} : a, b \in \mathbb{Z} \right\}$ | <ul style="list-style-type: none">- $\mathbb{R} = \mathbb{Q} \cup (\text{irrational \#s})$- $\mathbb{C} = \mathbb{Q} \cup (\text{imaginary \#s})$- $\pi, e, \phi \in \mathbb{Q}$- $\mathbb{C} = i, a+bi$ tho $a, b \in \mathbb{R}$ |
|---|--|
- $\mathbb{Z} = \text{int}, \mathbb{Z}^+ = + \text{int}, \mathbb{Z}^- = - \text{int}$
- 0 \neq + or -
- natural \#s $\mathbb{N} = \{0\} \cup \mathbb{Z}^+$

bits	min	max
8	-128	127
16	-32768	32767
32	-2147483648	2147483647
64	-9223372036854775808	9223372036854775807

bits	min	max
8	0	255
16	0	65535
32	0	4294967295
64	0	18446744073709551615

- computers do arithmetic in finite fields (large set of #s)
- digit for computer = bit (either 0 or 1)

INTS

unsigned } short
 } long
 } long long } int

unsigned } char

char = 8 bits
short = 16 bits
int = 16 bits
long = 32 bits
long long = 64 bits

#include <stdint.h>		
int8_t	uint8_t	8
int16_t	uint16_t	16
int32_t	uint32_t	32
int64_t	uint64_t	64
<u>signed</u>	<u>unsigned</u>	<u>bit size</u>

Binary Arithmetic

$$0+0=0 \quad 1+1=10$$

$$1+0=1$$

$$0 \times 0 = 0$$

$$1 \times 0 = 0$$

$$1 \times 1 = 1$$

$$101+11=1000$$

$$101 \times 101 = 11001$$

- finite fields have additive inverses for every $\#$
 - ex) $i = \#, \bar{i} = \text{add. inverse} \Rightarrow i + \bar{i} = 0$
- k bits \Rightarrow can have every $(+)$ int from 0 to $2^k - 1$
 - \hookrightarrow half will be add. inv so range = -2^{k-1} to 2^{k-1}

- Two's Complement (2C)


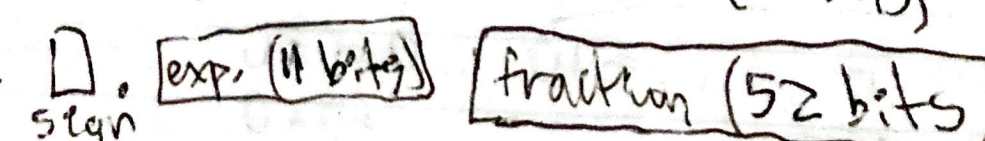
- to get add inverse of m :
 - flip bits of m (0 to 1 & 1 to 0) (\sim operator)
 - add 1 to result \rightarrow add. inverse
 - if 2C bits $>$ 06 $\#$, drop/carry highest-order bit
- Real $\#$ s = continuous or uncountably infinite

- Floating Point $\#$ s

- proper subset of real $\#$ s ($F \subsetneq \mathbb{R}$)
- proper subset of rational $\#$ s ($F \subsetneq \mathbb{Q}$)
- proper subset of ints ($F \subsetneq \mathbb{Z}$)
- not reals or rationals \Rightarrow FPNs are approximations

- float = single precision = 32 bits (IEEE 754 single)
- double = double precision = 64 bits (IEEE 754 double)
- long double = extended precision = 128 bits (IEEE 754 quad)

- float $\neq \mathbb{Q}$

- Single Precision = 
 - sign
 - exponent (8 bits)
 - fraction (23 bits)
- double precision = 
 - sign
 - exp. (11 bits)
 - fraction (52 bits)

- intel extended precision $\boxed{10}$ sign $\boxed{\text{exp. (15 bit)}}$ $\boxed{(1 \text{ bit})}$ int part $\boxed{\text{fraction (63 bit)}}$

- quad precision = 1 bit sign \cdot 15 bit exponent \cdot 112 bit fraction

- floating point limits

- float = 8 bit exponent, 23 bit mantissa

- double = 11 exp, 52 mant

- long double = 15 exp, 112 mant

- Little Endian = lowest address byte } has least sig bit
- Big Endian = highest address byte }

→ LE = $\boxed{78} \boxed{56} \boxed{34} \boxed{12}$ vs. BE = $\boxed{12} \boxed{34} \boxed{56} \boxed{78}$

- modulo = % = ints only

$\begin{array}{l} - \times / \% \\ + - \\ \ll \gg \\ < < = > > = \\ == != \\ \& \\ \wedge \\ \\ \&\& \\ \end{array}$	$\left. \begin{array}{l} \\ \\ \\ \\ \\ \\ \\ \\ \end{array} \right\} \text{left associativity}$	$\begin{array}{l} ? : \\ = \\ += \\ -= \\ \times = \\ /= \\ \% = \end{array}$	$\left. \begin{array}{l} \\ \\ \\ \\ \\ \\ \end{array} \right\} \text{right associativity}$
---	--	---	---

- 8 bits = 1 byte

- REMEMBERS FOR ACNZ
- $\log(x) = -\log(\frac{1}{x})$ for $x > 0$ } invert exponents
 - $\log(e^x) = e^{\log(x)}$
 - $x^{1/k} = \sqrt[k]{x}$
 - binary search solutions = monotonic
 - $x_{k+1} = \frac{1}{2}(x_k + \frac{n}{x_k}) \Rightarrow \text{sqr_newton}$
 - computer arithmetic approximates $\#s$
 - \hookrightarrow round to finite representations
 - round \uparrow = ceiling, \downarrow = floor, to 0 = truncate

- $(\&)$ = shift & ~~addition~~ addition
- $(/)$ = shift & subtraction
- $(+)$ = more than exclusive - XOR

- $\text{abs}(x) = \text{if } x < 0, \text{ then } -x; \text{ else } x$

- libraries should be numerical precise, fast, & portable

- Taylor series $\Rightarrow f(x) = \sum_{n=0}^{\infty} \frac{f^n(a)}{n!} (x-a)^n$

- ex) e^x centered at $a=0$ is $e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} \dots$ $-\infty < x < \infty$

- \sum should be for loops, \sum should be while loop

- computing e^x

- $k! > e^k$ for all x

- when k is large enough

- $\frac{x^k}{k!}$ gets smaller w/ each step, can stop computing when small enough

- Computing \sin (periodic $[-2\pi, 2\pi]$)

$$-\sin(x) = x - \frac{x^3}{6} + \frac{x^5}{120} - \frac{x^7}{5040} + \frac{x^9}{362880} \dots + O(x^{12})$$

- series converges vvv slowly unless $x \approx 0$

- don't directly compare floating point #s
↳ check to epsilon ϵ (absolute error)

$$-\frac{|x^* - x|}{|x|} < \epsilon \quad (\text{relative error} \rightarrow \text{more accurate})$$