

BAG

1

240 hours

Bit; 0 or 1
8 bit = 1 byte

$2^{8-1} = 255$

$2^{32-1} = 4 \text{ byte}$

$= 30^3 \times 4 - 1$
 $= (02)^3 \times 4 - 1$
 $\approx 00000000000000000000000000000000$

$\begin{array}{r} 0 \\ 1 \\ 2 \\ 3 \\ 4 \end{array}$

$$0100 = 4$$

$$1011 = 8 + 0 + 2 + 1 = 11$$

$$\text{stored no.} = (2^t - 1) + n$$

$$11 = 16 - 1 + n$$

$$\Rightarrow n = -4$$

#2

$$\begin{array}{r}
 1000 \\
 0^5 = 2^0 \\
 0^2 = 2^0 \\
 0^4 = 2^1 \\
 0^6 = 2^2 \\
 \hline
 1010 \\
 1111 \\
 \hline
 1001 \\
 0110
 \end{array}$$

0110

0010

110

m

$$\begin{array}{r}
 1000 \\
 0111 \\
 + 1 \\
 \hline
 1000
 \end{array}$$

$$\left\{
 \begin{array}{l}
 \text{rep } m \\
 0001 \\
 1110 \\
 1111
 \end{array}
 \right.$$

Estimate error
problem

- ① smallest # when added to 1 that increases \sum
- ② largest # when added to 1 that decreases \sum

$$(2^m - 1) - m + 1$$

$$x = (e, f) = f b^{e-9}$$

e - 4 bits - exponent

f - 3 bits

Sig - 1 bit

8 byte -

e - 11 bits

Let $b = 10$

$$\text{Then } 0.5 = 0.5 \times 10^0$$

$$\Rightarrow f = 0.5$$

$$e = 9$$

$$y_b < f < 1$$

$$2^4 (2^{20})$$

$$1.6 \times 10^7$$

$$E = 255(2047)$$

$f \neq 0$ NaN

$f = 0$

$$E = 255(2047)$$

$\pm \infty$

↓
sign

8 by 2
 $\frac{1}{2}$
 $\frac{1}{2}$

Bagla

2.1

$$x_1 = f_1 b^{e_1 - q}$$

$$x_2 = f_2 b^{e_2 - q}$$

Deputize exponents

0.32×10^{-1}	0.32×10^{-2}	0.64×10^{-2}
0.64×10^0	0.64×10^0	0.64×10^0
0.032	0.0032	0.0064
0.64	0.64	0.64
0.672	0.6432	0.6464
0.67×10^0	0.64×10^0	0.65×10^0

(exponents differ by more than 2 have no effect.)

Defⁿ: relative round off error
of e_{max} den

overflow - $\frac{f_{\text{max}}}{b^{e_{\text{max}} - q}}$
else if $e < 0$ then $\frac{f_{\text{min}}}{b^{e - q}}$
underflow - 0

else

fast

1001	
1000	
0000	00
0000	X X
0000	X X
1001	X X X
1001000	

$$1) e = \underbrace{e_1 + e_2 - q}_{\text{in } t+2} = f b^{e - q}$$

1 for sign

1 for overflow

$$\begin{array}{r} 0 \dots \\ 0 \dots \\ \hline 5 \times 2 \\ 0101 \\ 0010 \\ \hline 0000 \\ 0101 \times \\ 0000 \times \times \\ \hline 001010 = 10 \end{array}$$

2t-1 bit registers
are needed to
store multiply.

this can be
skipped
but will
give errors

$$x = s_i f_i b^{e_i - q}$$

$$x^{-1} = s_i f_i^{-1} b^{q - e_i}$$

$$= s_i \left(f_i^{-1} \cdot \frac{1}{b} \right) b^{q - e_i + 1}$$

$$1 \leq f_i \leq b$$

$$b^{-1} \leq f_i^{-1} \leq 1$$

$$|e_i| = |e|$$

$$s_n = \overline{s_n} + E_n$$

$$1 \geq |x_i|$$

Error Analysis

Comparing two in additn

$$S = \sum_{i=1}^n x_i$$

$$= \underbrace{x_1 + x_2 + x_3 + x_4 + \dots}_{\sum_{i=1}^n x_i}$$

$$S_i = S_{i-1} + x_i$$

$$S_1 = x_1$$

$$S_2 = x_1 + x_2$$

$$= (x_1 + x_2)(1 + \epsilon_2)$$

$$S_3 = ((x_1 + x_2)(1 + \epsilon_2) + x_3)(1 + \epsilon_3)$$

$$S_j = x_1 \prod_{i=2}^j (1 + \epsilon_i) + \sum_{k=2}^j x_k \prod_{i=k}^j (1 + \epsilon_i)$$

$$\text{Def}^n A \oplus B = (A + B)(1 + \epsilon)$$

$$-\{x_1(n-1)\epsilon + x_2(n-1)\epsilon + \dots + x_{n-1}(n-1)\epsilon\} \leq E_n \leq \{ \dots \}(1+\delta)$$

$$x_{n-1}\epsilon\}(1+\delta)$$

$$-|x_1|(1+\delta)\{ (n-1) + (n-1) \dots + (2+1) \} \leq E_n \leq |x_1|(1+\delta)\{ (n-1) + (n-1) \dots + (2+1) \}$$

$$n-1 + \frac{n(n-1)}{2} = \underbrace{\frac{(n-1)(n+2)}{2}}_{\sim n^2}$$

$$|E_n| \leq 1 \cdot \frac{n^2}{2} \in \Theta(1+\delta)$$
$$\leq \Sigma_n \frac{n \epsilon}{2} (1+\delta) \quad \left| \begin{array}{l} \text{claim: realistic} \\ \sqrt{n} \end{array} \right.$$

If you go $x_1 + x_2 + \underbrace{x_3 + x_4 + \dots + x_{n-1} + x_n}_{\text{error}}$

down: memory

up: accuracy + parallel

$$\leq \bar{c}_n (\log_2 n) \in \Theta(1+\delta)$$

bagla #2.2

Roundoff Errors

$$\bar{x} = f(a_1, a_2, \dots, a_n)$$

$$x = f(a_1(1+\epsilon_1), a_2(1+\epsilon_2), \dots, a_n(1+\epsilon_n))$$

$$E = x - \bar{x} \leq \sum_{i=1}^n \epsilon_i \left| \frac{\partial f}{\partial a_i} \right| \leq \bar{x} \sum_{i=1}^n \left| \frac{\epsilon_i}{\bar{x}} \right| \left| \frac{\partial f}{\partial a_i} \right|$$

$$\frac{|E|}{\bar{x}} \leq \sum_{i=1}^n \left| \frac{\epsilon_i}{f} \right| \left| \frac{\partial f}{\partial a_i} \right| \approx \# \text{ condition } \# \ll 1 \text{ (for small error)}$$

Mercury:

- Precession of perihelion
- Quadrupole moment of the sun
- Many body effects
- Relativistic effects

$$\text{Area}_2 = \frac{[f(x) + f(x + \Delta x)] \Delta x}{2}$$

Modeling error: Implicit

Mathematical Model [Modeling error]

computation I Model error [Truncation error]

Numerical Implementation [round off error]

Area₃ = $\frac{[f(x) + f(x + \frac{\Delta x}{2})]}{2} \Delta x + \frac{[f(x + \frac{\Delta x}{2}) + f(x + \Delta x)]}{2} \Delta x$



Target is not to do worse computationally, than the math model you start with.

NR: Truncation error will increase if you decrease the steps.

Round off errors will decrease if you increase steps?

Corcl: So there's a trade-off:

Interval Arithmetic

* overestimates at times

* singularities

$$a_i : [a_{i,\min}, a_{i,\max}]$$

$$x : [x_{\min}, x_{\max}]$$

Example: Planetary Motion

(a) Gravity [Newtonian]

(b) Point mass

(c) Two body problem

(d) Radiation Pressure.

$$= [f(x) + f(x + \Delta x)] \frac{\Delta x}{2} + \frac{\Delta x}{2} f(x + \frac{\Delta x}{2})$$

$$\Delta A = \text{Area}_3 - \text{Area}_2 = \frac{\Delta x}{4} \left(2f(x + \frac{\Delta x}{2}) - f(x) - f(x + \Delta x) \right)$$

$$\approx - \frac{(\Delta x)^3}{4} f''(x + \frac{\Delta x}{2})$$