

NC

A

Point I.

The answer by a computer is approximate.

Point II

$$x = \pm \left(d_1 + \frac{d_2}{\beta^3} + \dots + \frac{d_N}{\beta^{N-1}} \right) \cdot \beta^E$$
$$= \pm \sum_{k=1}^N \frac{d_k}{\beta^{k-1}} \cdot \beta^E;$$

Toy
 $\beta = 2$
 $N = 3$
 $L = -1$
 $G = 1$

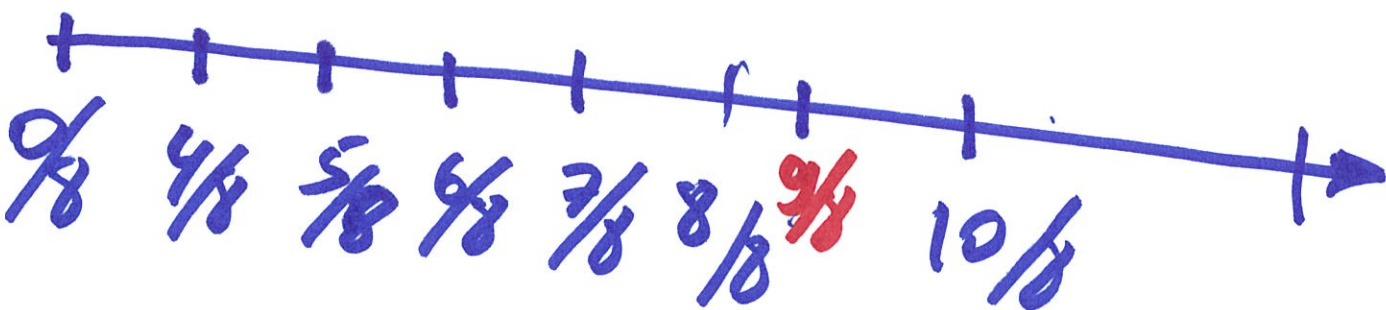
Double precision
IEEE

$\beta = 2$
 $N = 53$
 $L = -1022$
 $G = 1023$

Toy System



$4/8$	$5/8$	$6/8$	$7/8$
$8/8$	$10/8$	$12/8$	$14/8$
$16/8$	$20/8$	$24/8$	$28/8$



Finite # of #

- Almost every # can not be represented
- ∞ - number

FL(∞) - floating point presentation

$$E = \emptyset$$

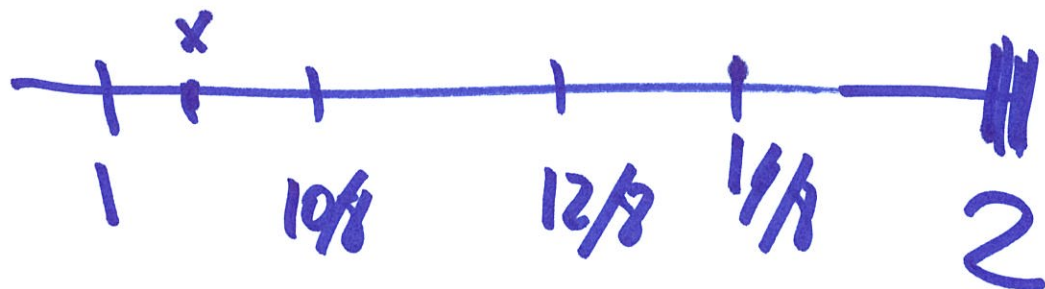
D

$$1.00 \quad 8/8$$

$$1.01 \quad 10/8$$

$$1.10 \quad 12/8$$

$$1.11 \quad 14/8$$



$$1 \leq x \leq 2$$

~~$$|x - F(x)| \leq 3/8$$~~

Round x to ~~nearest~~ $F(x)$

Chopping

E

$d_1. d_2 d_3 d_4 \dots d_N \rightarrow$

~~$d_1. d_2 d_3 \dots d_N d_{N+1} \dots d_{N+2} \dots$~~

Rounding to nearest

if there is a tie,
last digit is
even (zero $\equiv \beta=2$)

Chopping

$$|x - Fl(x)| < \frac{2}{8}$$

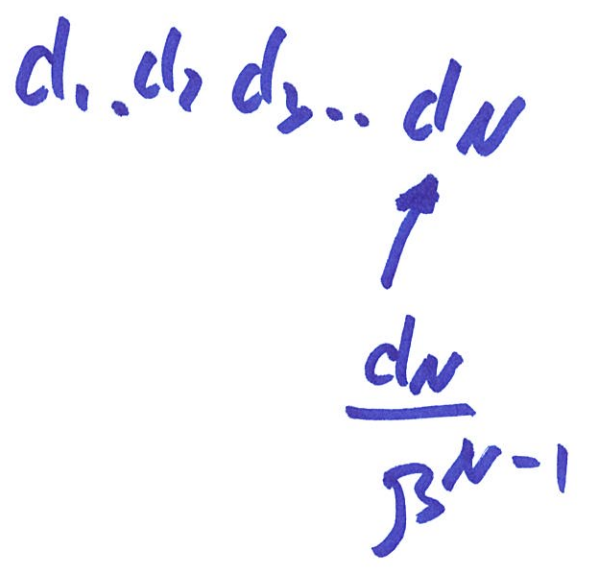
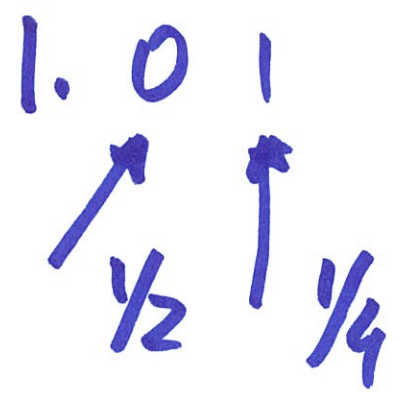
Rounding

$$|x - Fl(x)| < \frac{1}{8} ?$$

1.00
1.01

} machine epsilon
epsilon machine
unit round off

~~General~~



$\epsilon_{mach} = \beta^{1-N}$

$1 \leq x < 2$

Chopping $|x - Fl(x)| \leq \epsilon_{mach} = \beta^{1-N}$

Rounding to nearest:

$|x - Fl(x)| \leq \frac{\epsilon_{mach}}{2} = \beta^{-N}$

Rounding to
nearest $1 \leq x < 2$

G

$$\left| \frac{x - Fl(x)}{x} \right| \leq |x - Fl(x)| < \frac{\epsilon_{mach}}{2}$$

$$\left| \frac{x - Fl(x)}{x} \right| \leq \frac{\epsilon_{mach}}{2}$$

Relative error in representing
 $1 \leq x < 2$ is less than

$$\frac{\epsilon_{mach}}{2}$$

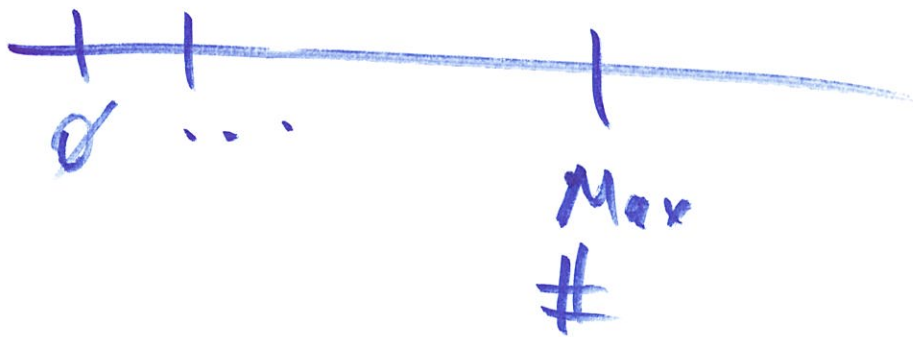
~~$x \rightarrow m$~~ $x \rightarrow m$ - mantissa

$$x = m \cdot \beta^e$$

11

Maximum relative error
in representing "any" number
is less than $\epsilon_{\text{mach}}/2 \equiv \beta^{-N}$

$$\left| \frac{x - Fl(x)}{x} \right| < \frac{\epsilon_{\text{mach}}}{2} = \cancel{\beta^{-N}} \\ = \frac{\beta^{1-N}}{2}$$



Largest # representable
in a system:

I

$1.11 \cdot 10^4$ - Toy

Double precision

6111

$\underbrace{\hspace{10em}}_{53 \text{ bits}} \cdot 10^{21}$

$$\text{Largest number} = \beta^{e+1} (1 - \beta^{-n})$$

Special numbers

\emptyset - zero

Inf \equiv Infinity

J

If Computer obtains
number $>$ "largest #"
The number is replaced
by Inf.

overflow

~~0~~, Inf \cdot 0, Inf / Inf

NaN - not a
number

$$S_1 = \sum_{k=-1000}^{1000} k \cdot e^{k^2} \quad - \text{over flow}$$

$$= \left(\sum_{k=-1000}^{-1} + (k=0) + \sum_{k=1}^{1000} \right) k e^{k^2}$$

$$\sum_{k=-1}^{-1000}$$

~~11~~

$$S_2 = \sum_{k=1}^{1000} \frac{e^k + 2}{e^k + 3}$$

$$= \sum_{k=1}^{1000} \frac{(e^k + 2)e^{-k}}{(e^k + 3)e^{-k}} =$$

$$= \sum_{k=1}^{1000} \frac{1 + 2e^{-k}}{1 + 3e^{-k}}$$

Largest # is
called OFL

real max \equiv matlab
underflow

Smallest # = β^L

smallest normalized #

$\beta = 2; L = -1022$

any number smaller than
smallest # is replaced
by zero.

LINDERFLOW

Smallest normalized
number

N

$$\text{real min} = \cancel{2^{-53}} \\ = 2^{-1023}$$

+

0

$$x = m_1 \cdot \beta^{\epsilon_1}$$

$$y = m_2 \cdot \beta^{\epsilon_2}$$

$$x + y = ?$$

To see m_1 shift mantissas
to match the exponents:
 $\beta = 10; n = 2$

$$1 + 10 = 1 + 1.0 \cdot 10^1 =$$

$$= 1.0 + 10 = 11$$

$$1.2 + 1.0 \cdot 10^1 =$$

$$= 1.2 + 10.0 = 11.2 \cong 11.0$$

$$1.0 \cdot 10^1 = 10.0 \cdot 10^0$$

$$\equiv$$

$$1.2 \times 10^{-1} + 1.3 \cdot 10^1 =$$

$$= 0.012 \cdot 10^1 + 1.3 \cdot 10^1 =$$

$$= 1.312 \cdot 10^1 \approx 1.3 \cdot 10^1$$

$$x - \Delta x < x < x + \Delta x$$

$$y - \Delta y < y < y + \Delta y$$

~~$$x + y - \Delta x < (x + y) < (x + y) + (\Delta x + \Delta y)$$~~

$$x + y - (\Delta x + \Delta y) <$$

Multiplications

X with d digits

Y with d digits

$X \cdot Y$ may $2d$ digits

Divide

$$\frac{X}{Y}$$

may have ∞ many
digits

Q

$$X = 1.3$$

$$Y = 1.2$$



2 significant digits

$$X - Y = 0.1$$

← one digit

$$X = \left(1 + \frac{\epsilon_{mach}}{2} \right)$$

$$Y = \left(1 - \frac{\epsilon_{mach}}{2} \right)$$

$$\begin{aligned} X - Y &= \left(1 + \frac{\epsilon_{mach}}{2} \right) - \\ &\quad - \left(1 - \frac{\epsilon_{mach}}{2} \right) = \end{aligned}$$

= ~~ϵ_{mach}~~ \Rightarrow computer will give \emptyset

Subnormals:

M

IF $E = L$ then

Leading digit may
be \emptyset .

"Subnormals" ? 1986

"Toy" system

$0.01 \cdot 10^{-1}$

$0.10 \cdot 10^{-1}$

$0.11 \cdot 10^{-1}$

gradual
underflow

$$x = 3\left(\frac{4}{3} - 1\right) - 1$$

S

$$\beta = 10, N = 2$$

\emptyset, NaN, Inf

$$\frac{4}{3} = 1.3$$

$$\frac{4}{3} - 1 = 0.3$$

$$3\left(\frac{4}{3} - 1\right) = 0.9$$

$$3\left(\frac{4}{3} - 1\right) - 1 = -0.1$$

$$\left| \frac{x - Fl(x)}{x} \right| < \frac{E_{mach}}{2}$$

$$f(x) = \frac{1 - (1-x)}{x}$$

$$x > 0, \quad 1-x < 1,$$

$$\frac{\epsilon_{mach}}{2} \text{ apart}$$

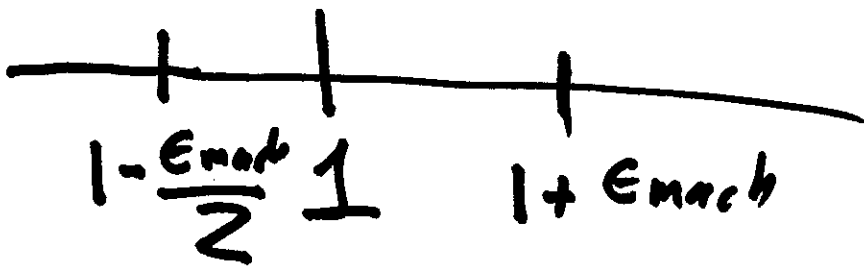
$$x < 0, \quad 1-x > 1$$

$$\epsilon_{mach} \text{ apart}$$

$$f(x) = \frac{1 - (1-x)}{x}$$

0

$$x = \frac{\epsilon_{mach}}{4}$$



$$\left(1 - \frac{\epsilon_{mach}}{4}\right) \approx 1 - \frac{\epsilon_{mach}}{2}$$

$$\begin{aligned} 1 - (1-x) &= 1 - \left(1 - \frac{\epsilon_{mach}}{2}\right) \\ &= \frac{\epsilon_{mach}}{2} \end{aligned}$$

$$\frac{1 - (1-x)}{x} \approx 2$$

$$\epsilon_{\text{mach}} = 2.2E-16$$

$$\epsilon_{\text{ps}} = 2.2E-16$$

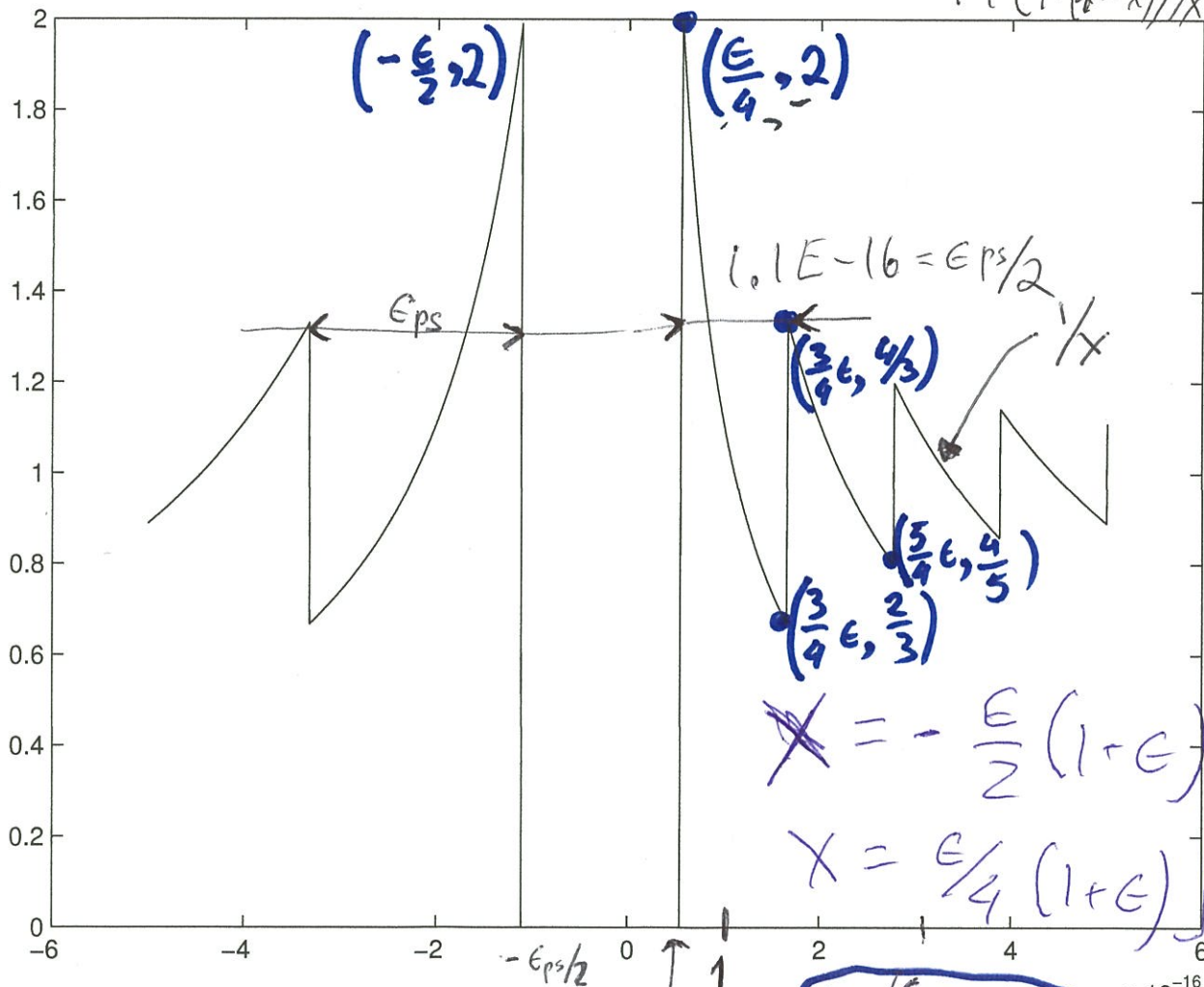
$$X = -\epsilon_{\text{ps}}/2 - \overset{\text{lost}}{\downarrow} \mu, \quad 1 - (1 - X) = 1 - (1 + \frac{\epsilon_{\text{ps}}}{2} + \mu) = -\epsilon_{\text{ps}}; \quad \frac{1 - (1 - X)}{X} = 2$$

$$\frac{1 - (1 - X)}{X}$$

$$X = -\frac{\epsilon_{\text{ps}}}{2} - 10^{-30}$$

$$(1 - (1 - X))/X = 1.99981$$

$$\frac{1 - (1 - X)}{X} = 2$$



$$X = -\frac{\epsilon}{2}(1 - \epsilon) \quad f(x) = 2$$

$$X = \frac{\epsilon}{4}(1 + \epsilon) \quad f(x) = 2$$

$$X = \left(\frac{\epsilon}{4} + \frac{\epsilon}{2}\right) \times 10^{-16}; \quad f(x) = \frac{1 - (1 - \epsilon)}{\frac{3}{4}\epsilon} = \frac{4}{3}$$

$$0.55E(-16) = \epsilon_{\text{ps}}/4; \quad X = \left(\frac{\epsilon}{4} + \frac{\epsilon}{2}\right)(1 - \epsilon) \quad f(x) = \frac{2}{3}$$



$$X = \epsilon_{\text{ps}}/4$$

$$1 - X = 1 - \frac{\epsilon_{\text{ps}}}{4} = \frac{1.000}{0.000} \quad 1 = 1 - \frac{\epsilon_{\text{ps}}}{2}$$

$$1.0111 \cdot 10^{-1}$$

$$\frac{1 - (1 - X)}{X} = 2$$

2 digits, $X = 0.05$

$$1 - (1 - 0.05) = 1 - (1 - 0.1) = 0.05$$

Matlab command used:

```
< M A T L A B (R) >  
    Copyright 1984-2023 The MathWorks, Inc.  
    R2023b (23.2.0.2365128) 64-bit (glnxa64)  
    August 23, 2023
```

To get started, type doc.
For product information, visit www.mathworks.com.

```
>> Inf/Inf  
Inf/Inf
```

```
ans =
```

```
NaN
```

```
>> -Inf+Inf  
-Inf+Inf
```

```
ans =
```

```
NaN
```

```
>> k = [ 1:1000];  
k = [ 1:1000];
```

```
>> y = (exp(k)+2)./(exp(k)+3);  
y = (exp(k)+2)./(exp(k)+3);  
>> S1 = sum(y)  
S1 = sum(y)
```

```
S1 =
```

```
NaN
```

```
>> Inf+Inf  
Inf+Inf
```

```
ans =
```

```
Inf
```

```
>> Inf/Inf  
Inf/Inf
```

```
ans =
```


NaN

```
>> y = (2*exp(-k)+1)./(3*exp(-k)+1);  
y = (2*exp(-k)+1)./(3*exp(-k)+1);  
>> sum(y)  
sum(y)
```

ans =

999.6577

```
>> Maximum = 2^(1023+1)*( 1-2^(-53))  
Maximu
```

```
>> m = 2^(1023+1)*( 1-2^(-53))
```

Maximum =

Inf

```
>> Maximum = 2^(1023)*( 1-2^(-53))*2  
m = 2^(1023+1)*( 1-2^(-53))
```

m =

Inf

```
>> Maximum = 2^(1023)*( 1-2^(-53))*2  
Maximum = 2^(1023)*( 1-2^(-53))*2
```

Maximum =

1.7977e+308

```
>> realmax  
realmax
```

ans =

1.7977e+308

```
>>
```

```
>> 2^(-1022)  
2^(-1022)
```

ans =

2.2251e-308

```
>> realmin  
realmin
```

ans =

2.2251e-308

```
>> 2 * realmax  
2 * realmax
```

ans =

Inf

```
>> realmin/2  
realmin/2
```

ans =

1.1125e-308

```
>> 2^(-52)  
2^(-52)
```

ans =

2.2204e-16

```
>> eps  
eps
```

ans =

2.2204e-16

```
f(3)
```

y =

1

```
>>
```

```
>>
```

```
>> f(eps*223/888)
```

y =

1.9910

>> f(-eps*449/888)

y =

1.9777

>> f(eps*667/888)

y =

1.3313

>> f(eps*664/888)

y =

0.6687

ans =

0.6687

>> f(eps*667/888)

y =

1.3313

>> f(1.25*eps)

y =

0.8000

f(1.25*eps)

y =

0.8000

f(-1.49*eps)

y =

0.6711

ans =

0.6711

>> f(-1.51*eps)

y =

1.3245

ans =

1.3245

=====

Here is how the file f.m looks like:

```
>> !more f.m
function y=f(x)
y = (1-(1-x))./x
end
```