

7 GP
 3.1. radix 16: IBM's f.p. format. $X = (-1)^S \times 16^{X_E - 64} \times (0.X_{17}^*)$
 ++, -- $X_E \Rightarrow$ moving the binary point over 4 bits.
 normalisation: at most 3 leading 0.
 0.0001
 0.00001
 $X_E - 64$
 IBM 32
 64
 $X_E - 64 = 1$ L 1: Sign.
 7: Exp.
 24: fractional part of mantissa
 56: total part of mantissa
 X_{17}

Shifter 1: accelerator maintains alignment.

if $k \geq 0 \Rightarrow$ Right Shift M2 by $4 \cdot k$ bits.

if $k < 0 \Rightarrow$ Right Shift M1 by $-4 \cdot k$ bits.

Adder 2: 56 bits. M1 - CLA.

Zero digit checker: accelerator normalisation of result
 - determines value of \underline{l} = number of groups of 4 leading 0s

$0000 \ 0000 \ 1$
 1 group of 4 leading 0s 1 group of 4 leading 0s.
 $\underline{l} = 2$

\rightarrow Shifter 2: Left Shift result by $4 \cdot \underline{l}$ bits.

\rightarrow Adder 3: subtract \underline{l} from result's exponent.

Adder 3: chooses the result's exponent by inspecting the sign of k

3.2. Rounding: convert higher precision representation into a lower precision representation.

- for storage / data transmission.

IEEE 754 4 rounding modes

$\rightarrow 0$
 \rightarrow towards $-\infty$
 \rightarrow towards $+\infty$
 \rightarrow to nearest even

Let $X = \overbrace{x_{m-1} x_{m-2} \dots x_1 x_0}^{m \text{ bits of integer part}} . \overbrace{x_{-1} x_{-2} \dots x_{-m}}^{m \text{ bits of fractional part}}$

for brevity, consider $X^* = \text{rounded value of } X$.

$- X^* - \text{integer}$

\Rightarrow eliminate the fractional bits

! For IEEE 754 f.p. nos, the rounding does not eliminate the fractional bits.

$$X = \begin{array}{|c|c|c|} \hline S & E & \text{fractional part of } X \\ \hline 1 & 8 & 23 \\ \hline \end{array}$$

$$X = (-1)^S \times 2^{E-127} \times (1.X_{23})$$

what IEEE 754 f.p. rounding does:

if X_{23} has more than 23 bits

\Rightarrow in order to be able to store the value \Rightarrow only preserve 23 bit

$$X \times 4$$

$$X_{23} \times X_4$$

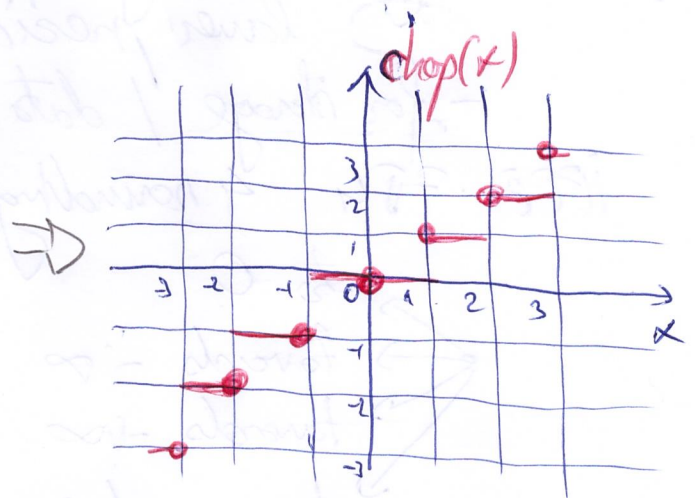
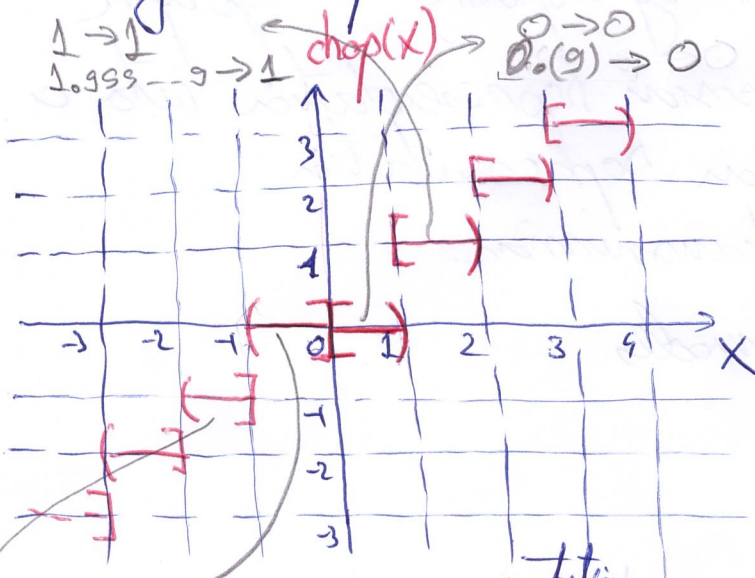
\Rightarrow 48 bit over

of 42 bits of fractional part
one can only store 23
- discard all fr. bits

A) to 0 (inwards rounding)

X^* : largest integer, for which $|X^*| \leq |X|$

if X is represented in S.N. \Rightarrow round to 0 \equiv truncate to integer



notation
[] \equiv .

$0 \rightarrow 0$
 $-0.999 \rightarrow 0$
 $-1 \rightarrow -1$
 $-1.999 \rightarrow -1$

- ③ towards $-\infty$ (downwards rounding)
- X^* is the largest integer for which $X^* \leq X$
 - ! - for positive values rounding to 0 \equiv rounding to $-\infty$
 - if X is represented in $C2 \Rightarrow$ downward rounding \equiv eliminating all fractional bits.

$$+4_{C2} = 0100_{C2} \quad \begin{matrix} .5 & .25 \\ \underbrace{2^{-1}} & \underbrace{2^{-2}} \end{matrix}$$

$$X = +4.25_{C2} = 0100.01_{C2} \quad 4.25 \rightarrow 4$$

$$X^* = 0100_{C2} = 4$$

(downwards rounding)

$$-4_{C2} = 1100_{C2}$$

$$-4.75_{C2} = 1011.11_{C2}$$

$$+4.75_{C2} = 0100.11$$

$$-4.75_{C2} = 1011.11$$

$$X^* = 1011_{C2} = -5$$

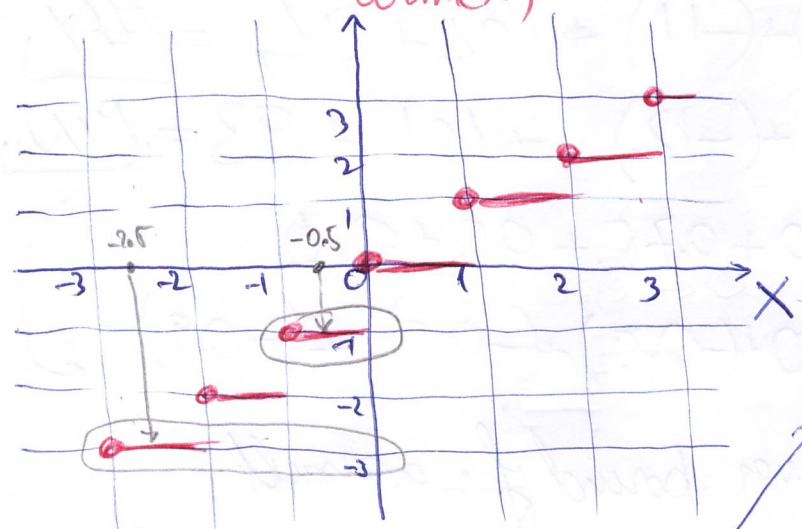
(downwards rounding)

down(X)

$$1011_{C2}$$

$$= -2^3 + 0.011_{C2}$$

$$= -8 + 3 = -5$$



if $X \geq 0$: discard (eliminate) all fractional bits.

if X is represented in SM :

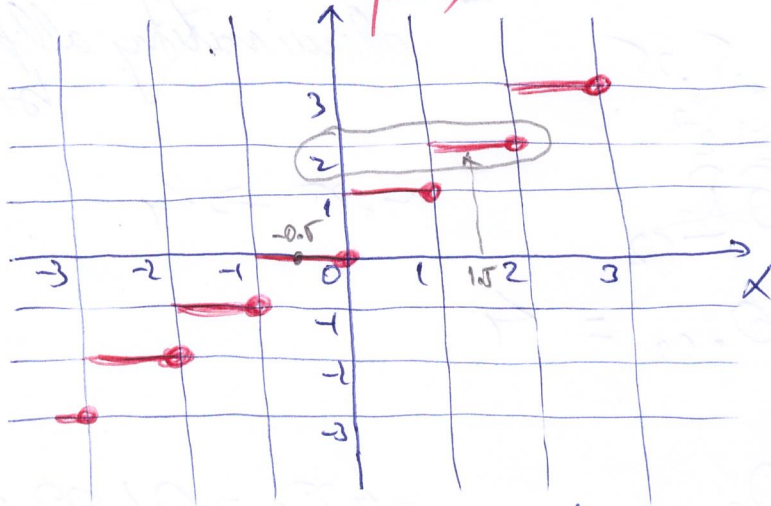
if $X < 0$:

$$X^* = \begin{cases} \overbrace{x_{m-1}x_{m-2} \dots x_1x_0}^{neg}, & \text{if } \cdot x_{-1}x_{-2} \dots x_{-m} = 0 \\ x_{m-1}x_{m-2} \dots x_1x_0 \text{ } \underline{\underline{-1}}, & \text{if } \cdot x_{-1}x_{-2} \dots x_{-m} \neq 0 \end{cases}$$

© towards $+\infty$ (upwards rounding)

X^* is the smallest integer for which $X^* \geq X$

if X is negative \Rightarrow upwards rounding \equiv inwards rounding



$$\text{error } \epsilon = X^* - X$$

upwards & downwards rounding:

- all errors are in the same direction \rightarrow either positive, or negative
 \Rightarrow errors accumulate faster

1.01×10^0

$$S = 1.1 + 2 + (-11.75) + 101.5 + 93.9$$

$$\text{upwards } 2 + 2 + (-11) + 102 + 94 = \text{MAX}$$

$$\text{downwards } 1 + 2 + (-12) + 101 + 93 = \text{MIN}$$

$$\epsilon_{\text{upwards}} = 0.9 + 0 + 0.75 + 0.5 + 0.1$$

$$\epsilon_{\text{downwards}} = (-0.1) + 0 + (-0.75) + (-0.5) + (-0.9)$$

- provide upper / lower bound for a result.

$$S \leq \text{MAX}_{\text{act}}$$

$$S \geq \text{MIN}$$

- interval arithmetic

! exponent of the values influence the error !!!

$$X_1 = 1.4 \times 2^{-32} \rightarrow 2 \times 2^{-32} \quad \epsilon = 0.6 \times 2^{-32}$$

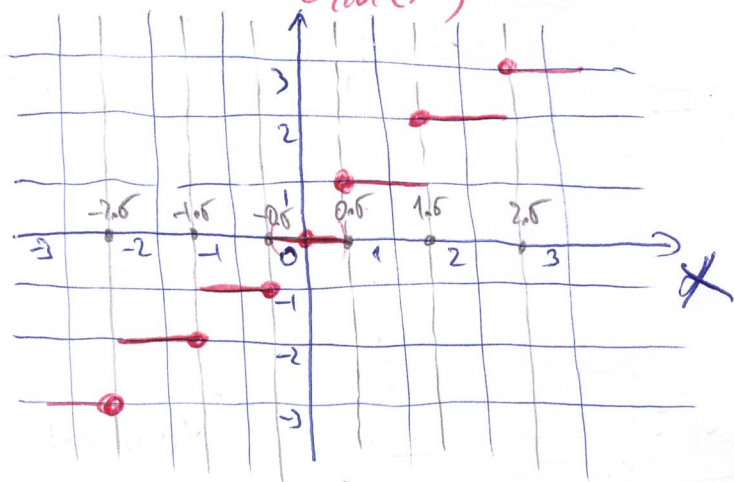
$$X_2 = 1.4 \times 2^{+32} \rightarrow 2 \times 2^{+32} \quad \epsilon = 0.6 \times 2^{+32}$$

① to nearest even
- derived from Round to nearest.

Round to nearest

Let X is positive $X^* = \begin{cases} X_{m-1}X_{m-2} \dots X_1X_0, & \text{if } 0.X_1X_2 \dots X_m < \frac{1}{2} \\ X_{m-1}X_{m-2} \dots X_1X_0 + 1, & \text{if } 0.X_1X_2 \dots X_m \geq \frac{1}{2} \end{cases}$

Similarly, can be defined round to nearest for negatives



Error analysis

- for brevity $X_{-3} = X_{-4} = \dots = X_{-m} = 0$
- only X_{-1} and X_{-2} can be ± 1

if $X > 0$

inputs		Output	
X_{-1}	X_{-2}	$X^* = rtm(X)$	$\epsilon = X^* - X$
0	0	$X_{m-1}X_{m-2} \dots X_1X_0$	0
0	1	$X_{m-1}X_{m-2} \dots X_1X_0$	$-\frac{1}{4}$
1	0	$X_{m-1}X_{m-2} \dots X_1X_0 + 1$	$+\frac{1}{4}$
1	1	$X_{m-1}X_{m-2} \dots X_1X_0 + 1$	$+\frac{1}{4}$

with equal probabilities

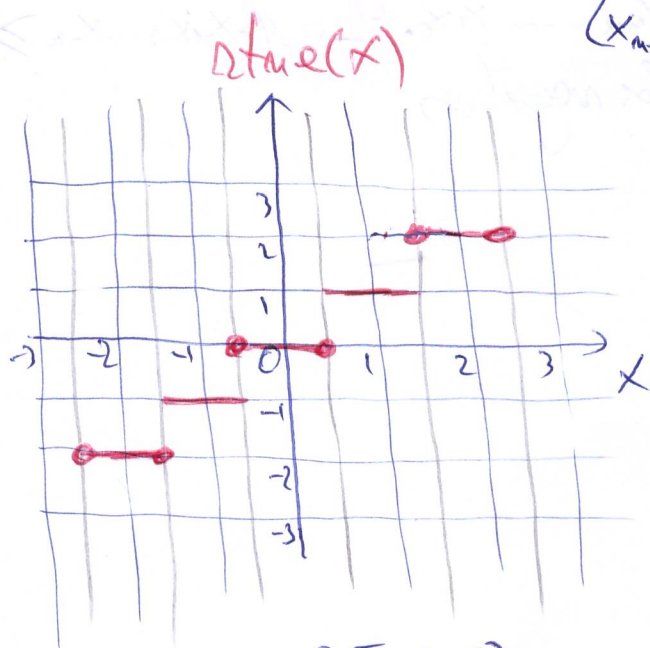
Consider all 4 cases \uparrow to be equally probable

$$\epsilon_{\text{mean}} = \frac{0 + (-\frac{1}{4}) + \frac{1}{4} + \frac{1}{4}}{4} = \frac{1}{8}$$

after 8 operations ϵ can become 1

! - if the case $0.10_{(2)} (0.5_{(10)})$ if more probable
 $\Rightarrow \epsilon_{\text{mean}} > \frac{1}{8}$

- round the case $\cdot 10_{(2)}$ with equal probability \nearrow upwards $(+\frac{1}{2})$
 \searrow downwards $(-\frac{1}{2})$
 rounding to nearest even: dropped bit X_0
 if X is positive $X^* = \begin{cases} X_{n-1}X_{n-2} \dots X_0, & \text{if } X_{n-1}X_{n-2} \dots X_m < \frac{1}{2}, \text{ or} \\ & X_{n-1}X_{n-2} \dots X_m = \frac{1}{2}, \text{ and } X_0 = 0 \\ X_{n-1}X_{n-2} \dots X_0 + 1, & \text{if } X_{n-1}X_{n-2} \dots X_m > \frac{1}{2}, \text{ or} \\ & X_{n-1}X_{n-2} \dots X_m = \frac{1}{2}, \text{ and } X_0 = 1 \end{cases}$



$0.5 \rightarrow 0$
 $-0.5 \rightarrow 0$
 $1.5 \rightarrow 2$
 $-1.5 \rightarrow -2$