

Floating Point Representation

$$\left(0. \underbrace{d_1 d_2 d_3 \dots d_m}_{\text{mantissa bit/fraction}} \right)_{\beta} \times \beta^e$$

exponent ← e

base ↘ β

(1) Convention 1 (Standard Form)

$$\pm \left(0. \underbrace{d_1 d_2 d_3 \dots d_m}_{d_1 \neq 0} \right)_{\beta} \times \beta^e$$

(2) Convention 2 (Normalized Form)

$$\pm \left(1. \underbrace{d_1 d_2 d_3 \dots d_m}_{d_1 \in [1, \beta-1]} \right)_{\beta} \times \beta^e$$

(3) Convention 3 (Denormalized Form)

$$\pm \left(0.1 \underbrace{d_1 d_2 d_3 \dots d_m}_{d_1 \in [1, \beta-1]} \right)_{\beta} \times \beta^e$$

Ex 1. Find the Largest Positive Numer (Highest Number) using m = 2

$$\beta = 2, \quad e_{\min} = -1 \quad \& \quad e_{\max} = 2$$

(1) Convention 1

$$(0.11)_2 \times 2^2$$

(2) Normalized

$$(1.11)_2 \times 2^2$$

(3) Denormalized

$$(0.111)_2 \times 2^2$$

Ex 2. Find the Lowest Positive Numer using m = 2

$$\beta = 2, \quad e_{\min} = -1 \quad \& \quad e_{\max} = 2$$

(1) Convention 1

$$(0.1\underline{0})_2 \times 2^{-1}$$

(2) Normalized

$$(1.\underline{00})_2 \times 2^{-1}$$

(3) Denormalized

$$(0.1\underline{00})_2 \times 2^{-1}$$

Ex 3. Find the Lowest Negative Numer using m = 2

$$\beta = 2, \quad e_{\min} = -1 \quad \& \quad e_{\max} = 2$$

(1) Convention 1

$$-(0.10)_2 \times 2^{-1}$$

(2) Normalized

$$-(1.00)_2 \times 2^{-1}$$

(3) Denormalized

$$-(0.100)_2 \times 2^{-1}$$

Ex 4. Find the Largest Negative Numer (Lowest Number) using m = 2

$$\beta = 2, \quad e_{\min} = -1 \quad \& \quad e_{\max} = 2$$

(1) Convention 1

$$-(0.11)_2 \times 2^2$$

(2) Normalized

$$-(1.11)_2 \times 2^2$$

(3) Denormalized

$$-(0.111)_2 \times 2^2$$

Ex 5. Find the Number of Non-Negative Combination using m = 2

$$\beta = 2, \quad e_{\min} = -1 \quad \& \quad e_{\max} = 2$$

$$e \rightarrow \begin{bmatrix} -1 \\ 0 \\ 1 \\ 2 \end{bmatrix} \rightarrow e^{\text{num}} = 4$$

(1) Convention 1

$$0.1 _ _ \rightarrow 2^1 \times 4 \quad 2^{m-1} \times 2^{\text{num}}$$

(2) Normalized

$$1. _ _ \rightarrow 2^2 \times 4 \quad 2^m \times 2^{\text{num}}$$

(3) Denormalized

$$0.1 _ _ \rightarrow 2^2 \times 4 \quad 2^m \times 2^{\text{num}}$$

For negative and positive combination, we have to multiply 2 with each value.

Decimal Shift

Shifting Decimal Left by m → power/exponent INCREASE by m

Shifting Decimal Right by m → power/exponent DECREASE by m

IEEE Format (64 bit)

1	11	52
---	----	----

sign exponent fraction | mantissa

↓
2¹¹

only positive : 0 → 2047

positive + negative : -1022 → 0 → 1025

Highest exponent to represent infinite → 2¹⁰²⁵

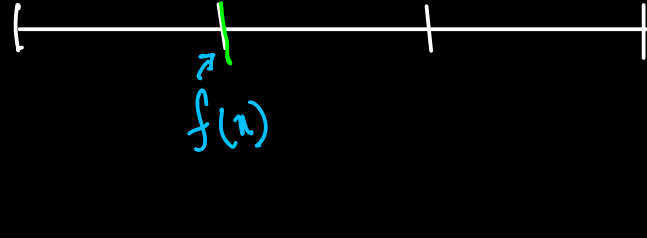
Lowest exponent to represent zero → 2⁻¹⁰²²

Rounding & Error

Actual Value \rightarrow

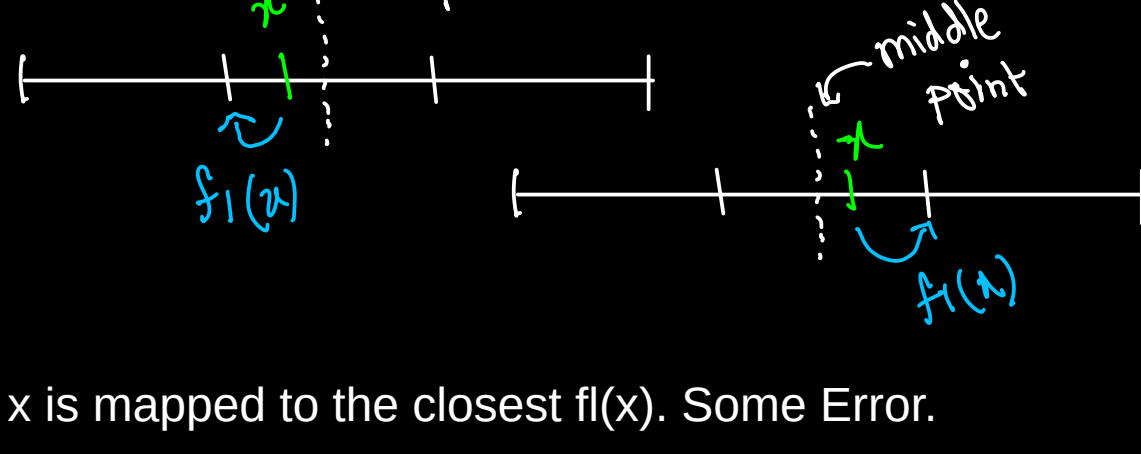
Floating Point Representation (F.P.R.) \rightarrow $fl(x)$

(1) $x = fl(x)$



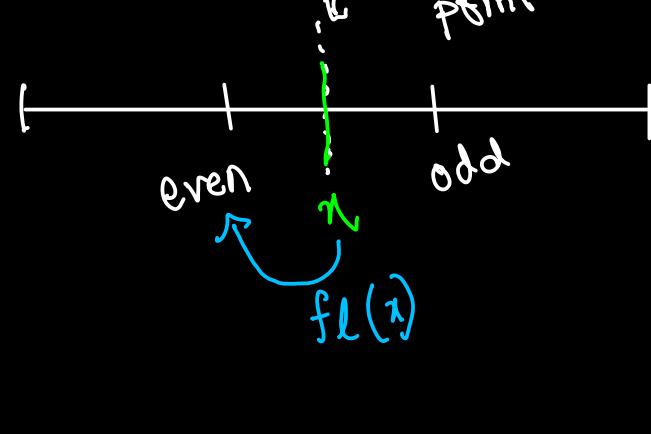
Since $x = fl(x)$, so x is mapped $fl(x)$. No Error.

(2) x not equals to $fl(x)$



x is mapped to the closest $fl(x)$. Some Error.

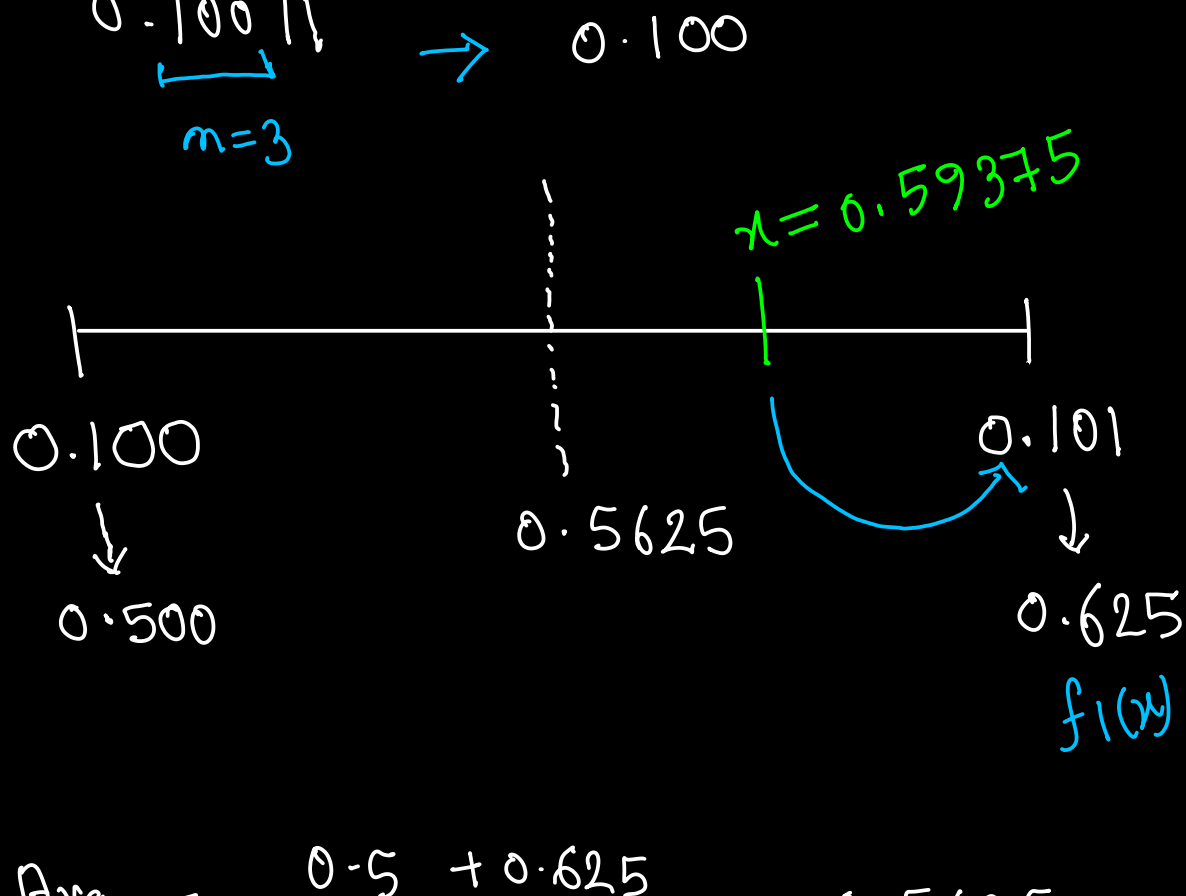
(3) x is on the middle point of two adjacent floating point



if $fl(x)$ is exactly on the middle of two floating point numbers, then x is mapped to the EVEN number.

NB: Even number in binary have LSB (Rightmost bit)

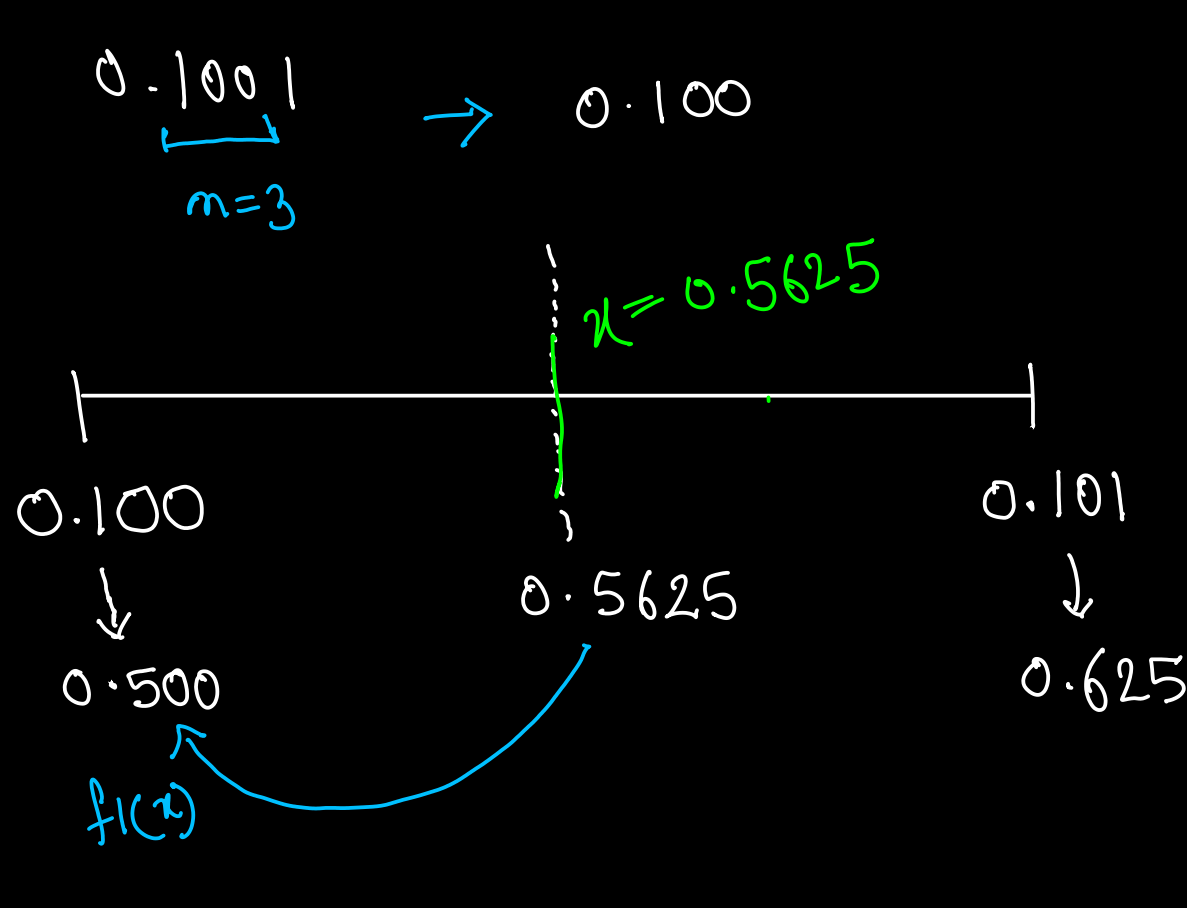
Ex.1 Using Convention 1, $m=3$ & $\beta=2$.
Round $(0.10011)_2$ to the appropriate f.p.r.



$$(0.10011)_2 \rightarrow (0.59375)_{10}$$

Floating point representation = $(0.101)_2$

Ex 2. Using Convention 1, $m = 3$ and $\beta=2$.
Round $(0.1001)_2$ to the appropriate f.p.r.



$$(0.1001)_2 \rightarrow (0.5625)_{10}$$

Floating point representation = $(0.100)_2$

Absolute Rounding Error

$$= |fl(x) - (x)|$$

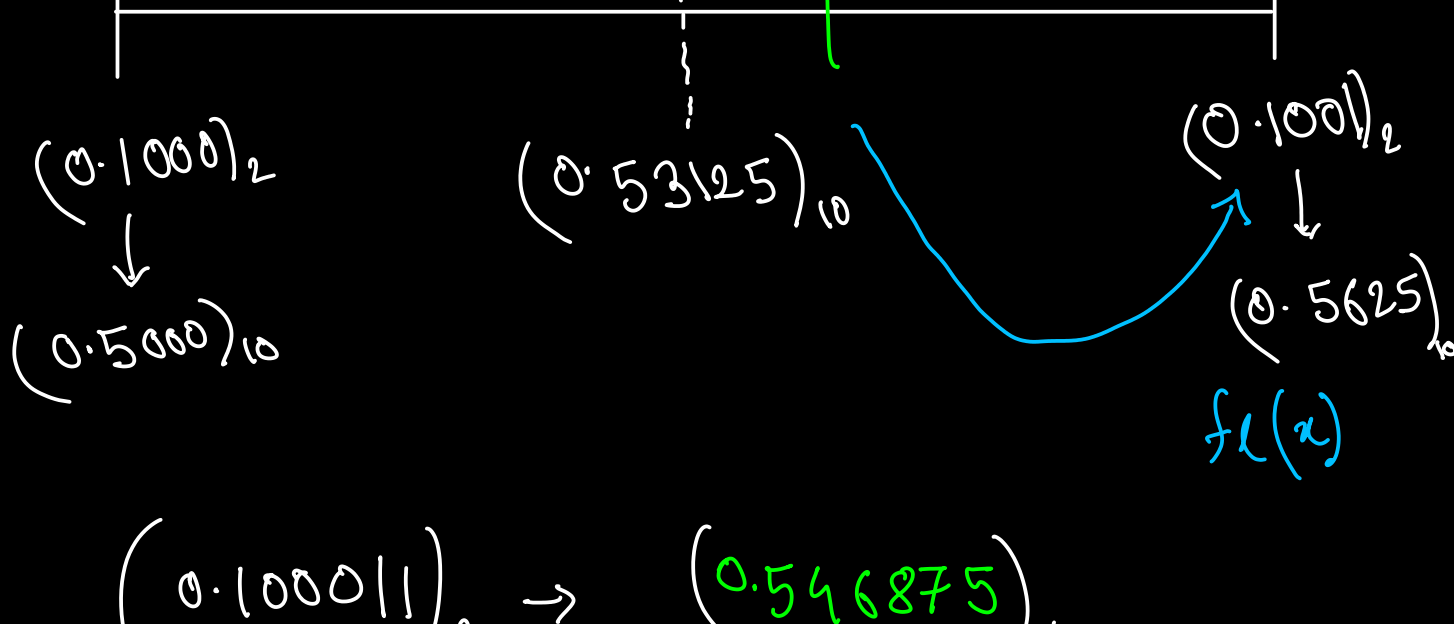
Scale Invariant Rounding Error (Relative rounding Error),

$$\delta = \frac{|fl(x) - (x)|}{|x|}$$

Ex 3. Let's say we have $x = \frac{5}{8}$ and $y = \frac{7}{8}$. Using Convention 1, $m = 4$ and $\beta=2$. Find $Fl(x.y)$. And find the relative rounding error.

$$x.y = \frac{5}{8} \times \frac{7}{8} = \frac{35}{64} = (0.546875)_{10} = (0.100011)_2$$

$$(0.100011)_2 \rightarrow (0.1000)_2$$



$$(0.100011)_2 \rightarrow (0.546875)_{10}$$

F.P.R. = $(0.1001)_2$

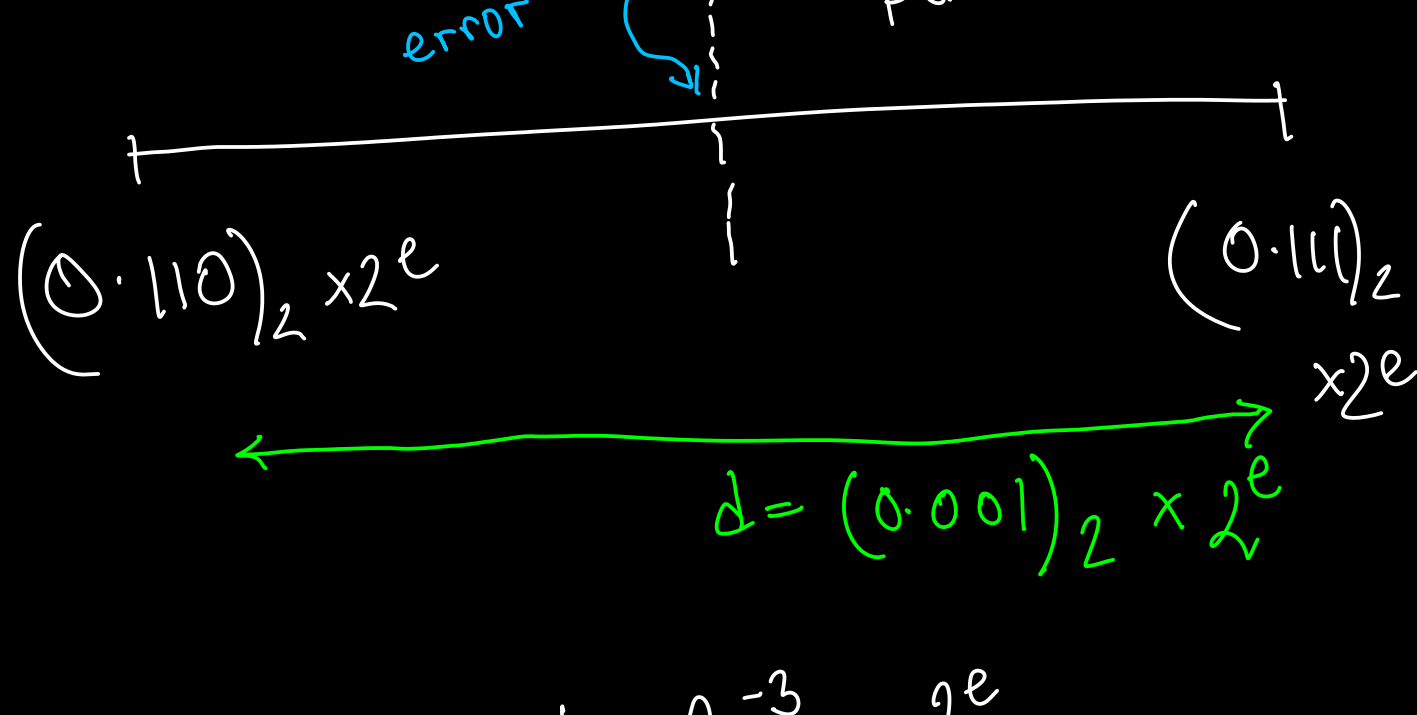
$$\delta = \frac{|0.5625 - 0.546875|}{|0.546875|} \times 100 = 2.86\%$$

Machine Epsilon = Maximize scale invariant rounding error

$$\epsilon = \delta_{\max} = \frac{|fl(x) - x|}{|x|}$$

Convention 1

$$(0.d_1d_2d_3\ldots)_\beta \times \beta^e$$



$$m=3: 1 \times 2^{-3} \times 2^e$$

$$m=4 [0.1100 \rightarrow 0.1101]: 0.0001 \times 2^e = 1 \times 2^{-4} \times 2^e$$

$$m=2 [0.10 \rightarrow 0.11]: 0.01 \times 2^e = 1 \times 2^{-2} \times 2^e$$

$$\beta^{-m} \beta^e$$

$$\text{max error} = |fl(x) - x|_{\max} = \frac{1}{2} \times \beta^{-m} \times \beta^e$$

Let's take $m=3$,

$$|x|_{\min} = 0.100 \times \beta^e = 1 \times 10^{-1} \times \beta^e$$

irrespective of the value of m , $|x|_{\min}$ will be same.

$$|x|_{\min} = \beta^{-1} \beta^e$$

$$\epsilon = \frac{\frac{1}{2} \beta^{-m} \beta^e}{\beta^{-1} \beta^e}$$

$$\epsilon = \frac{1}{2} \beta^{1-m} \quad [\text{for convention 1}]$$

Convention 2 (Normalized)

$$(1.d_1d_2d_3\ldots)_\beta \times \beta^e$$

$$\text{max error} = |fl(x) - x|_{\max} = \frac{1}{2} \beta^{-m} \beta^e$$

$$|x|_{\min} = (1.0)_2 \times 2^e = 1.0 \times 2^0 \times 2^e$$

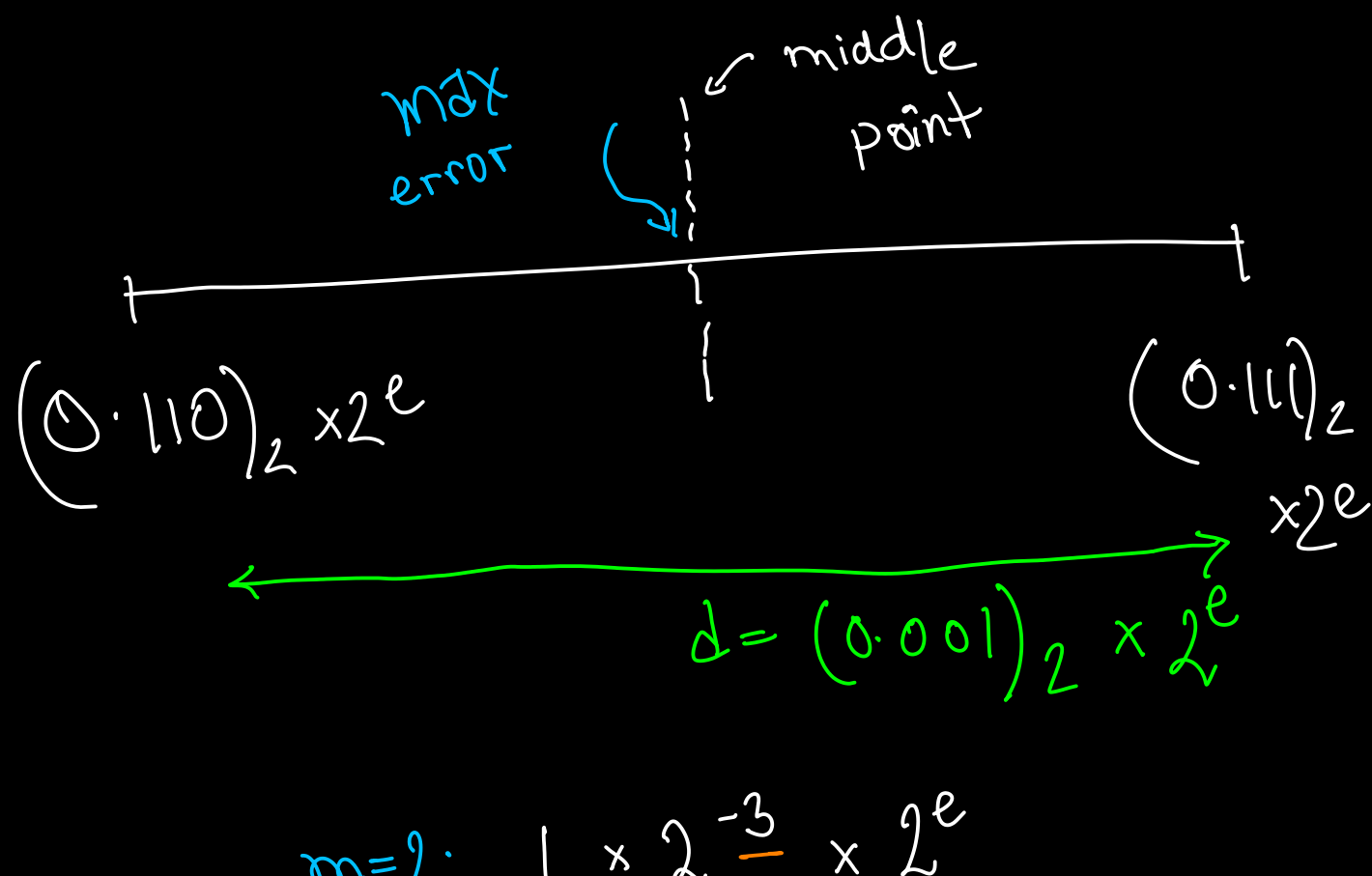
$$|x|_{\min} = \beta^0 \beta^e = \beta^e$$

$$\epsilon = \frac{\frac{1}{2} \beta^{-m} \beta^e}{\beta^e}$$

$$\epsilon = \frac{1}{2} \beta^{-m} \quad [\text{for Normalized}]$$

Convention 3 (Denormalized)

$$(0.1d_1d_2d_3\ldots)_\beta \times \beta^e$$



$$m=2: 1 \times 2^{-3} \times 2^e$$

$$m=3 [0.1100 \rightarrow 0.1101]: 0.0001 \times 2^e = 1 \times 2^{-4} \times 2^e$$

$$m=1 [0.10 \rightarrow 0.11]: 0.01 \times 2^e = 1 \times 2^{-2} \times 2^e$$

$$\text{max error} = |fl(x) - x|_{\max} = \frac{1}{2} \times \beta^{-\underline{m-1}} \times \beta^e$$

$$|x|_{\min} = (0.10)_e \times 2^e = 1 \times 2^{-1} \times 2^e$$

$$|x|_{\min} = \beta^{-1} \times \beta^e$$

$$\epsilon = \frac{\frac{1}{2} \times \beta^{-\underline{m-1}} \times \beta^e}{\beta^{-1} \beta^e} = \frac{1}{2} \beta^{-m}$$

