

MORE ON MACHINE NUMBERS

float

$$9.4 = (1001.\overline{0110})_2$$

$$= (1.001 \ 0110 \ 0110 \ 0110 \ 0110$$

$$0110 \ 0110 \ 0110 \ 0110$$

$$0110 \ 0110 \ 0110 \ 0110 \ 0110 \dots) \times 2^3$$

↑ 52nd bit

$$fl(9.4) = (1.001 \ 0110 \ 0110 \ 0110 \ 0110$$

$$0110 \ 0110 \ 0110 \ 0110$$

$$0110 \ 0110 \ 0110 \ 0110 \ 1) \times 2^3$$

We ① chopped $(0.110 \overline{0110}) 2^{-52} \cdot 2^3$ & ② added $2^{-51} \cdot 2^3$

$$\text{So } fl(9.4) = 9.4 - \underbrace{(0.110\overline{0110})}_{\textcircled{a}} \times 2^{-52} \cdot 2^3 + \underbrace{2^{-52} \cdot 2^3}_{\textcircled{b}}$$

$$\textcircled{a} = (0.\overline{0110}) \times 2^{-52} \cdot 2^4 = (0.4) 2^{-48}$$

$$\textcircled{b} = 2^{-49}$$

$$\begin{aligned} fl(9.4) &= 9.4 - (0.4) 2^{-48} + 2^{-49} \\ &= 9.4 + (0.2) 2^{-49} \end{aligned}$$

Created with Doceri



We saw $fl(9.4) = 9.4 + (0.2) 2^{-49}$

Define machine precision

$$\epsilon_{mach} = 2^{-52}$$

Notice $\left| \frac{fl(9.4) - 9.4}{9.4} \right| = \frac{0.2}{9.4} 2^{-49} = \frac{1}{47} 2^{-49}$

$$= \frac{8}{47} \epsilon_{mach}$$

Rule (normalized #s) $\left| \frac{fl(x) - x}{x} \right| \leq \frac{1}{2} \epsilon_{mach}$



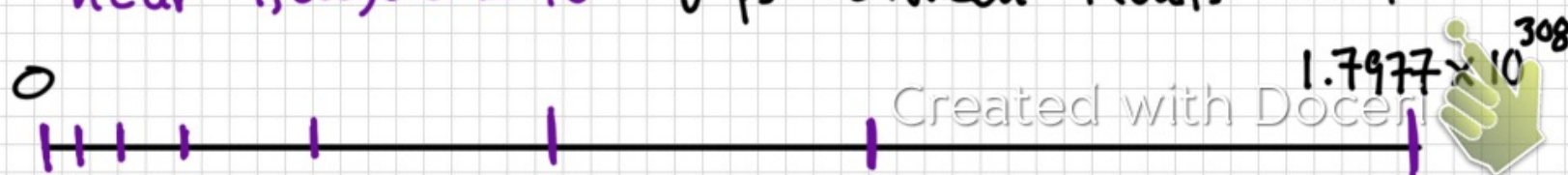
$$\epsilon_{\text{mach}} = 2^{-52} \simeq 2.22 \times 10^{-16}, \text{ so } \frac{1}{2} \epsilon_{\text{mach}} \simeq 1.11 \times 10^{-16}$$

$$\left| \frac{\text{fl}(x) - x}{x} \right| \leq \frac{1}{2} \epsilon_{\text{mach}} \simeq 1.11 \times 10^{-16}$$

near 1 $= 10^0$ gaps between floats about 10^{-16}

near 1000 $= 10^3$ gaps between floats about 10^{-13}

near 1,000,000 $= 10^6$ gaps between floats about 10^{-10}



Note: ϵ_{mach} is the smallest positive # such that

$$\text{fl}(1 + \epsilon_{\text{mach}}) \neq 1.$$

$$1 + \epsilon_{\text{mach}}$$

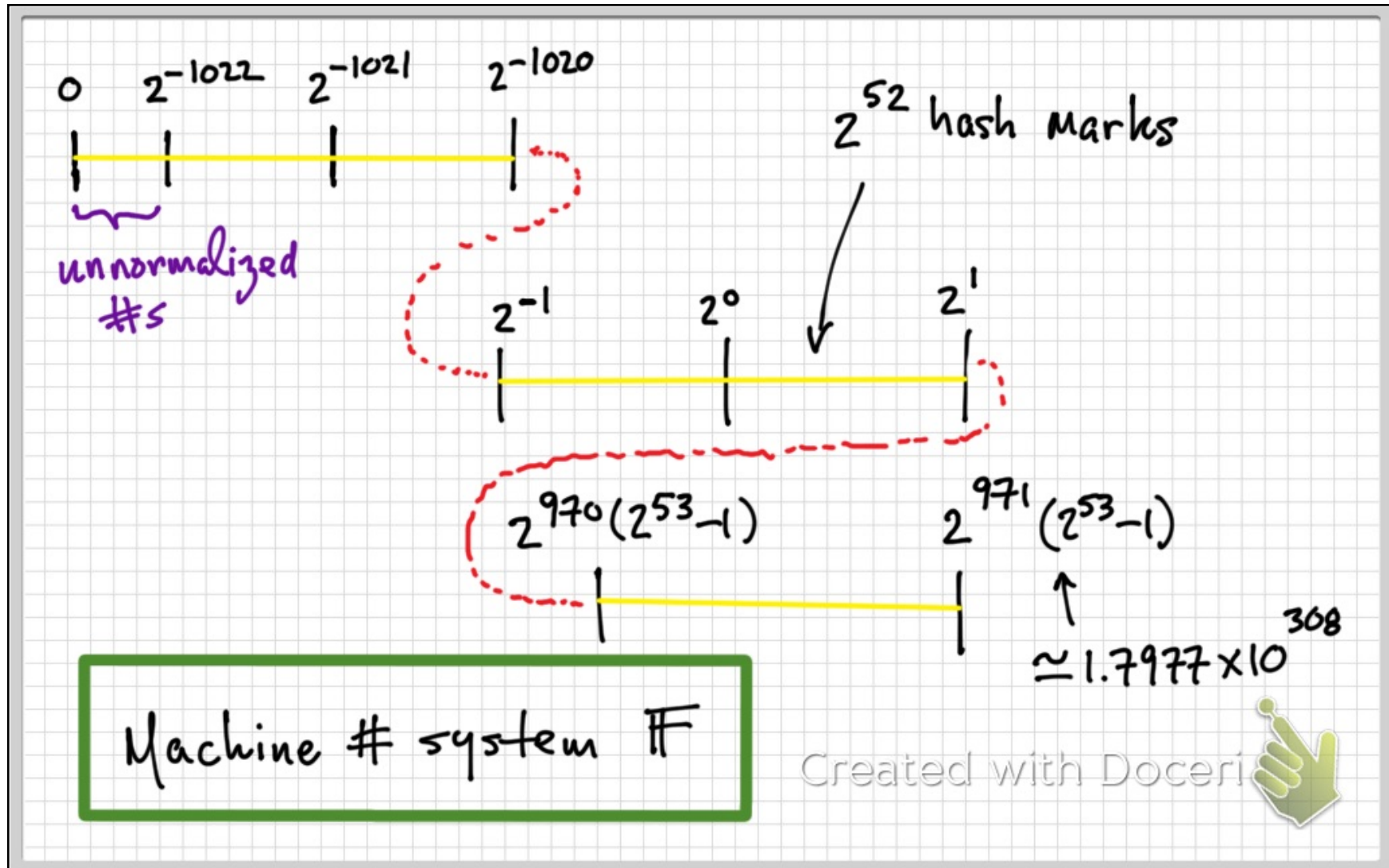
$$= (1.000 \dots \dots \dots 0) \times 2^0$$

$$+ (0.000 \dots \dots \dots 1) \times 2^0$$

$$= (1.000 \dots \dots \dots 1) \times 2^0$$

Created with Doceri





$$V_{\max, \text{norm}} = \underbrace{(1.111\dots1)_2}_{52 \text{ bits}} \times 2^{2046-1023}$$

Notice $(1.11\dots1)_2 + \underbrace{(0.00\dots1)_2}_{2^{-52}}$

$$= (10.00\dots0)_2$$

$$= 2 \Rightarrow (1.11\dots1)_2 = 2 - 2^{-52} = (2^{53} - 1) 2^{-52}$$

$$V_{\max, \text{norm}} = (2^{53} - 1) 2^{-52} 2^{1023} = \underline{\underline{(2^{53} - 1) 2^{471}}}$$



Arithmetic

The computer makes mistakes!

$+$, $-$, \times , $/$

basic ops

\oplus , \ominus , \otimes , \oslash

IEEE implementation

In general $x \oplus y \neq x + y$ even if $x = fl(x)$
 $y = fl(y)$

Model:

$$x \oplus y = (x + y)(1 + \mu)$$

where $|\mu| \leq \epsilon_{mach}$

$$\left| \frac{(x \oplus y) - (x + y)}{x + y} \right| = |\mu| \leq \epsilon_{mach}$$

Created with Doceri

