

Punto Flotante

Forma de codificar números reales, basada en la notación científica $\pm 7,13 \times 10^{-6}$ pero usando base 2. $\pm 1, \dots \times 2^e$

Representamos los números con una cantidad fija de bits. 32 bits precisión simple
64 bits precisión doble
usa octavo por defecto

\Rightarrow Solo podemos representar finitos números. n bits $\rightarrow 2^n$ combinaciones.

② Usamos precisión simple. Mantisa

$$X = \begin{matrix} + \\ - \end{matrix} (1 + F) 2^e$$

1 bit signo
23 bits mantisa
8 bits exponente
1 bit precisión doble
52 bits

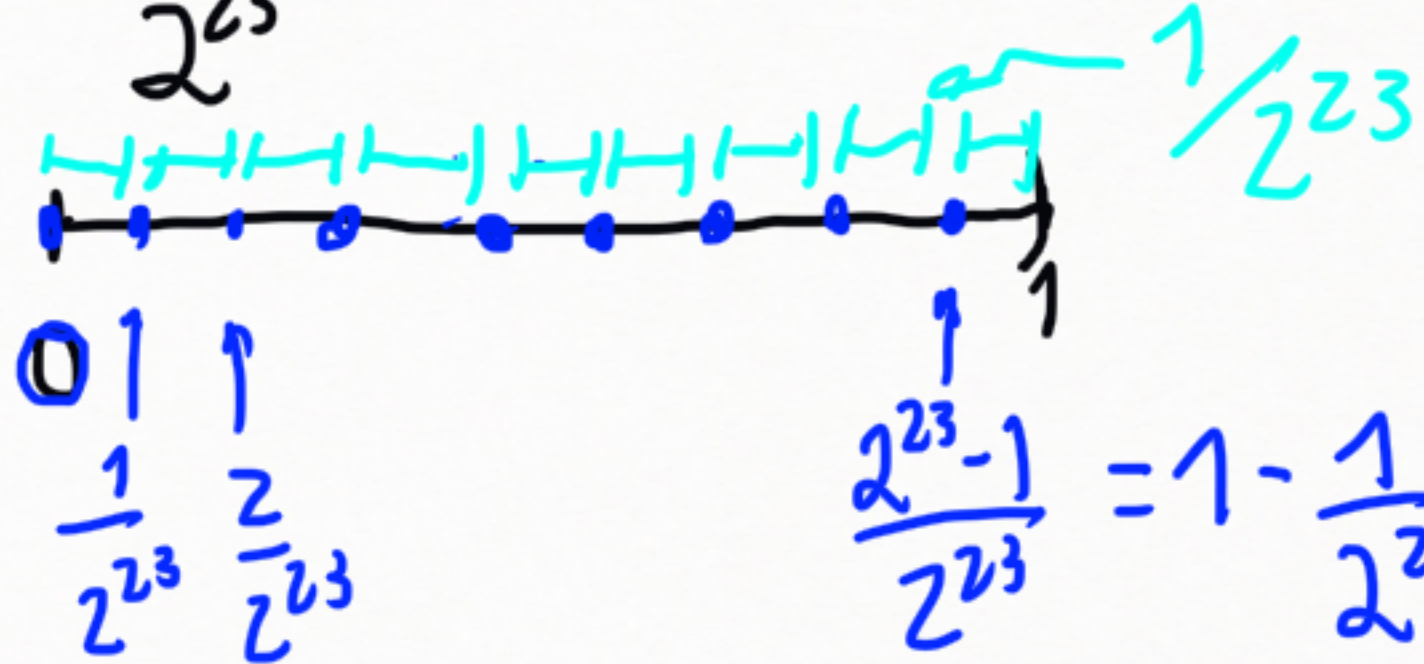
$$0 \leq F < 1$$

$$F = \frac{j}{2^{23}} \text{ con } j \in \{0, 1, \dots, 2^{23}-1\}$$

son 2^{23}

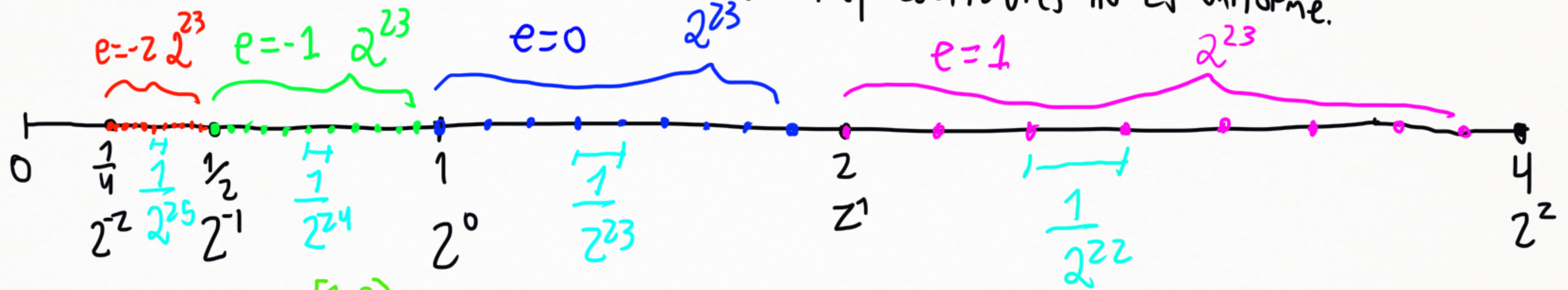
2^{23} posibles valores
equivalentes

$$0, \frac{1}{2^{23}}, \frac{2}{2^{23}}, \dots, \frac{j}{2^{23}}, \dots, \frac{2^{23}-1}{2^{23}}$$



Exponente
 $-126 \leq e \leq 127$
+ los valores 128, -127 se reservan para ciertos casos
0, INF, NaN

La distribución de los números representables no es uniforme.



Exponente = e : $x = [1, 2) 2^e$

$0 \leq F < 1$ \Rightarrow $2^e \leq x < 2^{e+1}$

$F=0 \rightarrow x = 2^e$

$F=1 \rightarrow x = (1+1)2^e = 2 \cdot 2^e = 2^{e+1}$
(nota sup)

$e \rightarrow x \in [2^e, 2^{e+1})$

por los distintos valores de F , tenemos 2^{23} números equiespaciados.

separación entre uno y el siguiente: 2^{e-23}

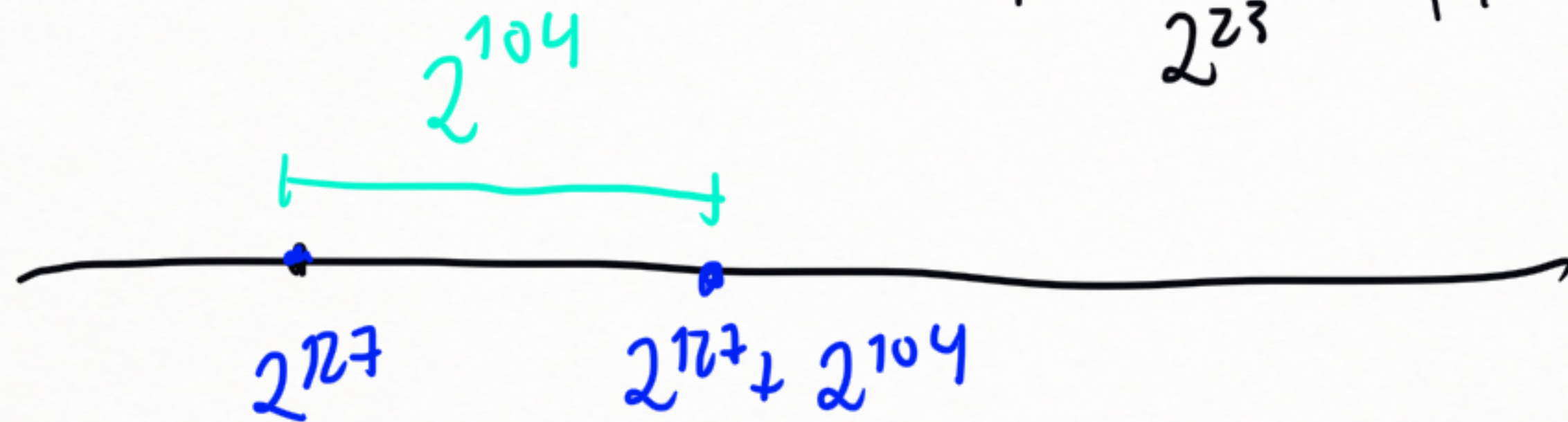
$F=0 \rightarrow 2^e$

$F = \frac{1}{2^{23}} \rightarrow (1 + \frac{1}{2^{23}})2^e = 2^e + 2^{e-23}$

$F = \frac{2}{2^{23}} \rightarrow (1 + \frac{2}{2^{23}})2^e = 2^e + 2 \times 2^{e-23}$

La máxima separación es con el exponente más grande.

$e = 127$: $x = (1+F)2^{127}$ $F=0$: 2^{127}
 $F = \frac{1}{2^{23}}$: $(1 + \frac{1}{2^{23}})2^{127} = 2^{127} + \frac{2^{127}}{2^{23}} = 2^{127} + 2^{104}$



Normalmente nos importa más el error relativo que el absoluto.

$x \in \mathbb{R}$ $x \neq 0$

$e_x := \hat{x} - x$ error absoluto

\hat{x} aproximación.

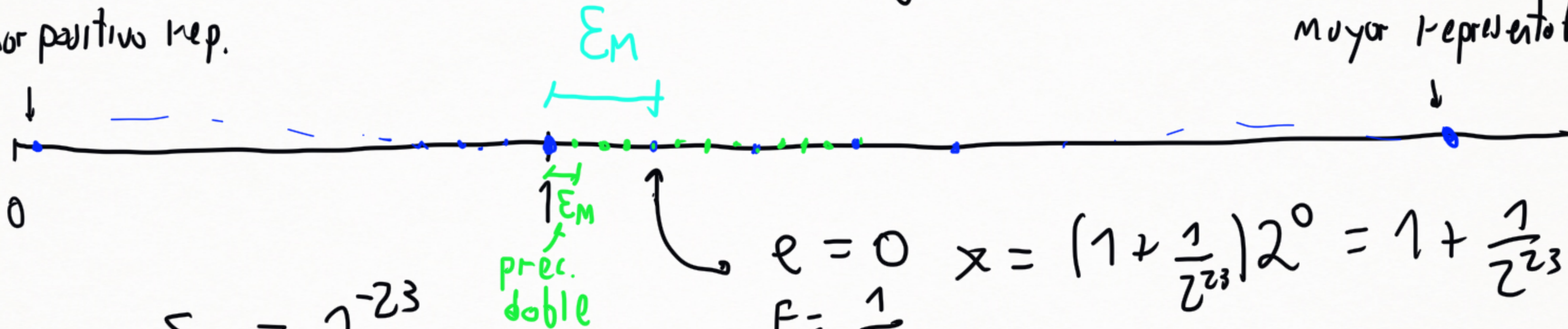
$\epsilon_x := \frac{\hat{x} - x}{x}$ error relativo.

$x = 30$ $\hat{x} = 32$ $\epsilon_x = \frac{32-30}{30} = \frac{2}{30} = \frac{1}{15}$ $x = 1000000000$ $\epsilon_x = \frac{2}{1000000000}$ chico
 $\hat{x} = 1000000002$

ϵ_M : la separación entre 1 y el siguiente número representable.

menor positivo rep.

maya representable



en este caso $\epsilon_M = 2^{-23}$

$e = 0 \quad x = (1 + \frac{1}{2^{23}})2^0 = 1 + \frac{1}{2^{23}}$
 $F = \frac{1}{2^{23}}$

en precisión doble, 52 bits para mantisa, $\epsilon_M = 2^{-52}$

Mayor representable

Menor representable

$(1+F)2^e$ F, e la más grandes posible

$e = 127$
 $F = \frac{2^{23}-1}{2^{23}} = 1 - \frac{1}{2^{23}}$

$x = (1 + 1 - \frac{1}{2^{23}})2^{127}$
 $= (2 - \frac{1}{2^{23}})2^{127} = 2^{128} - 2^{104}$

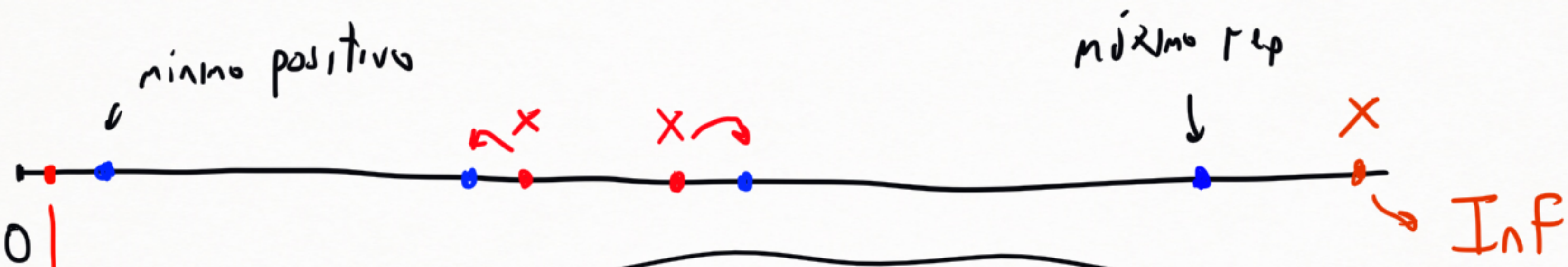
$(1+F)2^e$
 $F=0, e=-126$
 $x = 2^{-126}$

$\epsilon_M \approx 1,2 \times 10^{-7}$

$M_{\max} \approx 3,4 \times 10^{38}$

$M_{\min} \approx 1,2 \times 10^{-38}$

e	sep.	2^{e-23}	$(1+F)2^e$
$e=23$	sep = 1		$(1 + \frac{j}{2^{23}})2^e$
$e=24$	sep = 2		$= 2^e + 2^{e-23}j$



Usar subnormales

$$g(x) = \frac{x^3}{x - \sin(x)}$$

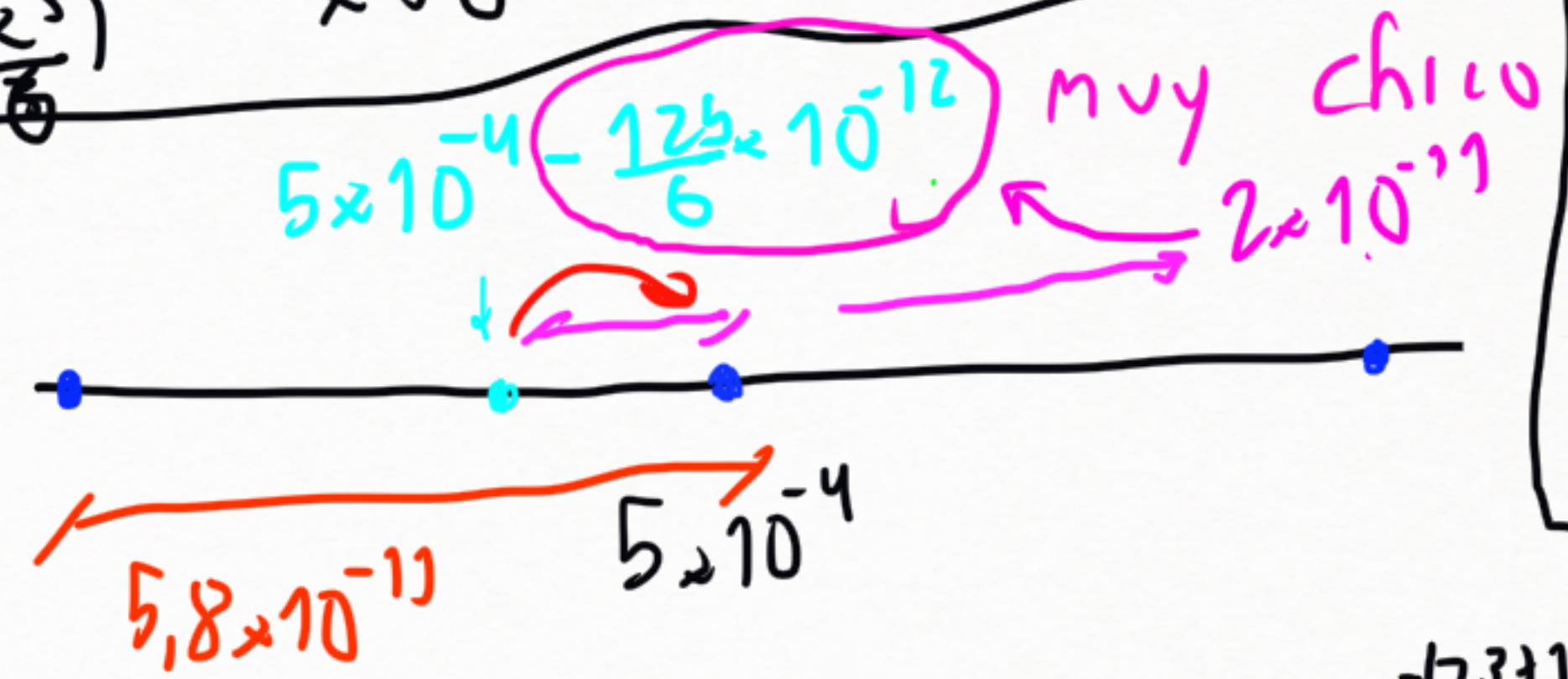
$$\sin(x) = x - \frac{x^3}{6} + r(x) \quad (\text{Taylor order 3})$$

$$\lim_{x \rightarrow 0} g(x) = \lim_{x \rightarrow 0} \frac{x^3}{x - \sin(x)} = \lim_{x \rightarrow 0} \frac{x^3}{x - (x - \frac{x^3}{6})} = \lim_{x \rightarrow 0} 6 = 6$$

$$x - \sin(x) \approx x - (x - \frac{x^3}{6})$$

$$5 \times 10^{-4} - \left(5 \times 10^{-4} - \frac{125 \times 10^{-12}}{6} \right)$$

5×10^{-4}



$$5 \times 10^{-4} \in [2^{-11}, 2^{-10})$$

separación

$$2^{-(23+11)} = 2^{-34} \approx 5.8 \times 10^{-11}$$

e sep 2^{e-23}