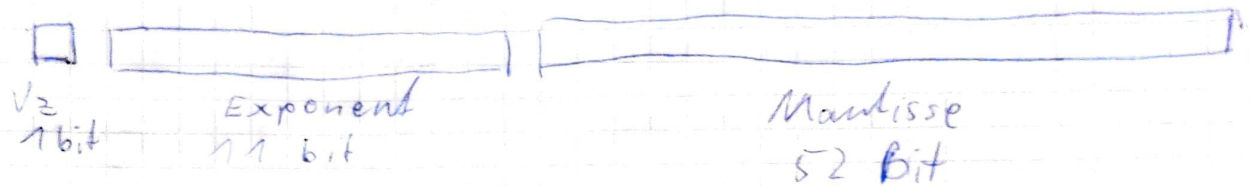
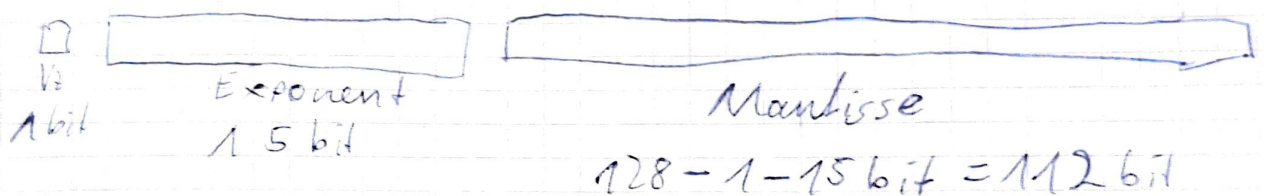


# Assignment 1

1.) double precision (64 bit)



quadruple precision



Machine Epsilon:

$$\epsilon_{\text{bit}} = b^{-(p-1)}$$

$b \dots$  Base  $\hat{=} 2$

$p \dots$  precision  $\hat{=} 113$

$$\epsilon_{128} = 2^{-(113-1)} = 2^{-112} = \underline{\underline{1,9259299e-34}}$$

~~Max significant digits (decimal) =  $2^{113} = 5.1922969$~~

~~Max significant digits (decimal) =  $\pm(2 - 2^{-112}) \cdot 2^{16383}$~~

~~range = Max sig. digit =  $\pm(2 - 2^{-p}) \cdot 2^{\text{offset}}$~~

~~range<sub>quad</sub> =  $\pm(2 - 2^{-112}) \cdot 2^{16383}$~~

significant digits  $\hat{=} \log_{10}(2^{113}) = 34,01638951 \approx \underline{\underline{34}}$

Why is quadruple precision more than twice as accurate?

~~Def~~ Precision describes the amount of bits reserved for the ~~mantissa~~ whole floating point word. ~~Doubling the bits means to double the~~  
Doubling the precision means to double the AMOUNT of bits. The decimal value of a  $n$  bits represented value gets doubled, when adding one bit (e.g.):

$$8_{10} = 1000$$

$$16_{10} = 10000$$

Therefore doubling the digits results in a much higher accuracy.

$\epsilon = b^{-(p-1)} \Rightarrow$  The precision is an exponent for calculating the machine epsilon, which describes the error.

2.)

$$a) fl(a \text{ op } b) = fl(b \text{ op } a)$$

true.



The rounding error of  $a$  and  $b$  does not change.

$$b) fl(a + a) = fl(2 * a)$$

true



The rounding error that is existing with  $a$  gets doubled, no matter what.  
"  $a$  before the operation is infinitely precise. ~~It's not~~  
It only is rounded after one operation.

$$c) fl((a+b) + c) = fl(a + (b+c))$$

false. X

$a+b$  might produce a different rounding error than  $b+c$ .  $\Rightarrow$  leads to a different result.

for  $a$  &  $b$ : we only have one operation

for  $c$ : we have two operations, dependant ~~for~~ on each other.



3.)

$$S_n = \sum_{i=1}^n X_i$$

$$\hat{S}_1 = fl(x_1) = x_1 (1 + \delta_1)$$

$$\begin{aligned} \hat{S}_2 &= fl(\hat{S}_1 + x_2) = (\hat{S}_1 + x_2)(1 + \delta_2) = \\ &= (x_1(1 + \delta_1) + x_2)(1 + \delta_2) = \end{aligned}$$

$$= x_1(1 + \delta_1)(1 + \delta_2) + x_2(1 + \delta_2)$$

$$\Downarrow \boxed{1 + \delta_i = 1 \pm \delta_i}$$

$$\hat{S}_2 = x_1(1 \pm \delta)^2 + x_2(1 \pm \delta)$$

$$\begin{aligned} \hat{S}_3 &= fl(\hat{S}_2 + x_3) \cancel{1 \pm \delta} = (\hat{S}_2 + x_3)(1 \pm \delta) = \\ &= \underline{x_1(1 \pm \delta)^3 + x_2(1 \pm \delta)^2 + x_3(1 \pm \delta)} \end{aligned}$$

$$\left( \prod_{i=1}^n (1 + \delta_i)^{p_i} = 1 + \theta_n \right) \quad \text{Backward-Error}$$

Backward-  
Error Bound:

$$\hat{S}_n = \sum_{i=1}^n X_i (1 \pm \delta)^{n-i+1}$$

↑ rounding errors  
are independent!

$$\hat{S}_n = \sum_{i=1}^n X_i (1 + O_{n-i+1})$$

$$\left| \sum_{i=1}^n x_i - f\left(\sum_{i=1}^n x_i\right) \right| \leq y_n \sum_{i=1}^n |x_i|$$

$$f(x_1, x_2) = (x_1 + \Delta x_1) + (x_2 + \Delta x_2)$$

$$\sum_{i=1}^n x_i + \sum_{i=1}^n \Delta x_i = \sum_{i=1}^n (x_i + \Delta x_i) \quad |\Delta x_i| \leq y_n |x_i|$$

forward Error

$$\sum_{i=1}^n x_i + y_n |x_i|$$

$$\begin{pmatrix} 1 & 2 \\ 0 & 3 \\ 0 & 0 \end{pmatrix}$$