

Floating Point Zahlen

→ für was stehen die einzelnen Stellen?

↳ immer für Potenzen der jeweiligen Basis

e.g. dezimal (Basis 10) 157.9

$1 \cdot 10^2 + 5 \cdot 10^1 + 7 \cdot 10^0 + 9 \cdot 10^{-1}$
 exp. der Basis pro Stelle:
 3 2 1 0 . -1 -2 -3
 $10^3 \ 10^2 \ 10^1 \ 10^0 \ 10^{-1} \ 10^{-2} \ 10^{-3}$

außerdem: $1.579 \cdot 10^2 = 157.9$ also $\cdot 10^{\text{exp}}$ heißt "punkt"
 $157.9 \cdot 10^{-2} = 1.579$ verschieben

→ genau so mit Basis 2

exponent	3	2	1	0	-1	-2	-3
Zweierpotenz	2^3	2^2	2^1	2^0	2^{-1}	2^{-2}	2^{-3}
Wert (B=10)	8	4	2	1	0.5	0.25	0.125

↳ mit $\cdot 2^{\text{exp}}$ multiplizieren heißt auch "Punkt" verschieben wenn wir in binär reden

$1.101 \cdot 2^2 = 110.1$

→ umrechnen

→ aufteilen: $1 + 0.9$ → ganz-zahliger Teil normal als Ob Zahl
 → Nachkomma-Stellen umrechnen

immer die Zahl links vom Punkt (0 oder 1)

immer $\cdot 2$ rechnen

immer in die nächste Zeile schreiben

repeating part

$0.9 - 0 \rightarrow 0.9 \rightarrow 1.8$
 $1.8 - 1 \rightarrow 0.8 \rightarrow 1.6$
 $1.6 - 1 \rightarrow 0.6 \rightarrow 1.2$
 $1.2 - 1 \rightarrow 0.2 \rightarrow 0.4$
 $0.4 - 0 \rightarrow 0.4 \rightarrow 0.8$
 $0.8 - 0 \rightarrow 0.8 \rightarrow 1.6$
 $1.6 - 1 \rightarrow 0.6 \rightarrow 1.2$
 $1.2 - 1 \rightarrow 0.2 \rightarrow 0.4$
 $0.4 - 0 \rightarrow 0.4 \rightarrow 0.8$
 $0.8 - 0 \rightarrow 0.8 \rightarrow 1.6$
 $1.6 - 1 \rightarrow 0.6$

1
0
0
0

⇒ resultat

welcher Teil ist periodisch, also wiederholt sich? (wenn überhaupt)

1.0
 $0.1110011001100...$

 1.11100

Normalized floating-point system

Basis: $b_i \in \{0, \dots, b-1\}$ → mit Basis $b=2$
 $b_0 \neq 0$ $b_i \in \{0, 1\}$

$$F(b, p, e_{\min}, e_{\max})$$

$e \in \{e_{\min}, \dots, e_{\max}\}$

Präzision:

(inkl. leading 1)

d.h. ihr könnt p Stellen benutzen:

p Stellen

$$b_0 . b_1 b_2 b_3 \dots b_{p-1} \cdot b^e$$

"normalized" heißt es ist immer $1. \dots$
und nie $0. \dots$

e.g. $0.25_{10} = 0.01_2$

$$0.01_2 \rightarrow 1 \cdot 2^{-2}$$

→ zwei floating-point Zahlen addieren.

- 1) Beide Zahlen auf den gleichen Exponenten bringen
- 2) Zahlen addieren (Punkt und 2^e beim Resultat nicht vergessen)
- 3) Re-normalisieren
- 4) Runden

↳ abrunden falls nächste Stelle 0 ist
auf 4 Stellen: $[1.001]010$

↳ aufrunden falls nächste Stelle 1 ist

$$[1.001]101$$

d.h. jetzt diese 1 "aufrunden" d.h. diese 1
wird zu 10 → 1 als carry

$$1.010$$

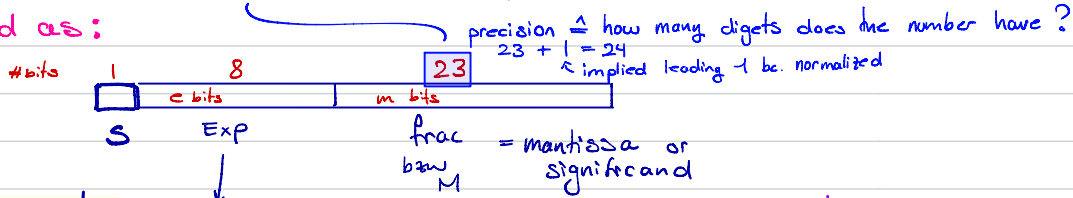
⇒ Woher kommt das? (Ihr müsst floats umrechnen etc. können, diese Seite ist einfach für die, die es interessiert)

single precision: float

$$\rightarrow (-1)^S \cdot M \cdot 2^E$$

$$F^*(2, 24, -126, 127)$$

stored as:



exponent

$$E = \text{Exp} - \text{Bias} \quad \text{Bias} = 2^{e-1} - 1$$

(→ a way of encoding E to be \ominus and \oplus , different from two's complement → this way is better for comparison)

→ single p. → Exp is the binary number stored in those 8 bits ⇒ Exp is in range $\{0, \dots, 255\}$
 as we have 8 bits ⇒ 0 to 255. $\text{Exp} = 0 \rightarrow$ denorms, $\text{Exp} = 255 \rightarrow$ special vals

→ so to get the E from your Exp bits, given $\text{Exp} \in \{1, \dots, 254\}$ do:

$$E = \text{Exp} - \text{Bias}, \quad \text{Bias for float} = 2^{8-1} - 1 = 2^7 - 1 = 128 - 1 = 127$$

⇒ range:

$$e_{\min} = 1 - 127 = -126$$

$$e_{\max} = 254 - 127 = 127$$

⇒ i.e. to get our number $(-1)^S \cdot M \cdot 2^E$ we use $E = \text{Exp} - \text{Bias}$ but E is encoded in memory through Exp.

mantissa / significand

⇒ encoded with an implicit leading 1 → "normalized"

$$M = 1. \underbrace{\text{xxxx} \dots \text{x}}_{23 \text{ digits}}$$

→ these are the bits that are actually stored.

min : all zeros $\hat{=} M = 1.0$

max : all ones $\hat{=} M = 1.1 \dots 1$

Denorms: Exp is all zeros

$$\rightarrow E = -\text{Bias} + 1 = -126$$

→ frac all zeros → ± 0

→ frac $\neq 0$ → small values closer to 0
 "subnorms"

Special Vals Exp all ones

↳ frac all zero → $\pm \text{inf}$

↳ frac $\neq 0$ → NaN

$$F^*(2, 4, -2, 2)$$

↑ Basis 2
p=4

$$E \in \{e_{\min}, e_{\max}\}$$

-2 2

Welche Zahlen sich darstellen?

State the following numbers in $F^*(2, 4, -2, 2)$:

1. the largest number;
2. the smallest number;
3. the smallest non-negative number.

! heisst normalisiert
 • → die Zahlen müssen also immer 1.xxx sein

Compute how many numbers are in the set $F^*(2, 4, -2, 2)$.

Wir können in unserem System mit $p=4$ genau 3 bits variieren.

(das 4te bit ist ja schon durch die leading 1 gefüllt)

⇒ müssen bei jeder Frage überlegen wie wir die 3 bits und den exponent wählen

1) grösste Zahl? → frac feld so gross wie mögl. → 111 } $1.111 \cdot 2^2$
 E so gross wie mögl. → 2

2) kleinste Zahl? einfach ⊖ die grösste Zahl d.h. sign bit s = -1

3) Zahl > 0, "am nächsten bei 0" → E am kleinsten wählen

→ die 3 bits in frac so klein wie mögl.

⇒ $\left. \begin{matrix} 000 \\ -2 \end{matrix} \right\} 1.000 \cdot 2^{-2}$

↑ nicht vergessen, es muss ja immer norm. sein.

→ wie viele Zahlen können dargestellt werden?

80 → $5(2 \cdot 2^3)$

pro E $\left\{ \begin{array}{l} 2^3: \square\square\square \rightarrow 3 \text{ bits, jeweils zwei Mögl. (entweder 0,1)} \\ 2: \text{immer } \oplus \text{ und } \ominus \end{array} \right.$

5: die mögl. E sind -2, -1, 0, 1, 2 i.e. 5