

#### Eigenschaften des IEEE-P 754

Die Basis  $b$  ist gleich 2.

Das erste Bit der Mantisse wird implizit zu 1 angenommen, wenn die Charakteristik nicht nur Nullen enthält.

Normalisierung: das erste Bit der Mantisse (die implizite 1) steht vor dem Komma.

Ist die Charakteristik gleich 0, entspricht dies dem gleichen Exponenten wie die Charakteristik 1.

Das erste Bit der Mantisse wird aber dann explizit dargestellt

Auch die Null ist darstellbar

#### Eigenschaften des IEEE-P 754

Sind alle Bits der Charakteristik gleich 1, signalisiert dies eine Ausnahmesituation.

Wenn zusätzlich die Mantisse gleich Null ist, wird die Situation „overflow“ (bzw. die „Zahl“  $\pm \infty$ ) kodiert.

Dies erlaubt es dem Prozessor, eine Fehlerbehandlung einzuleiten.

Intern arbeiten Rechner nach dem IEEE-Standard mit 80 Bit, um Rundungsfehler unwahrscheinlicher zu machen.

Charakter.	Zahlenwert	Bemerkung
0	$(-1)^{1/2} 0 \cdot \text{Mantisse} \cdot 2^{1-1023}$	$-0$ ✓ oder $+0$
1	$(-1)^{1/2} 1 \cdot \text{Mantisse} \cdot 2^{1-1023}$	
...	$(-1)^{1/2} 1 \cdot \text{Mantisse} \cdot 2^{1-1023}$	
2046	$(-1)^{1/2} 1 \cdot \text{Mantisse} \cdot 2^{1023}$	
2047	Mantisse = 0: $(-1)^{1/2} \infty$	Overflow
2047	Mantisse $\neq 0$ : NaN	Not a number

## Aufgabe 1: Fließkommazahlen

### Fragen

Beantworten Sie folgende Fragen:

- Nennen Sie Beispiele für Festlegungen, die der IEEE 754 Standard mitbringt. Warum ist eine solche Standardisierung sinnvoll?
- Was ist die betragsmäßig größte bzw. kleinste darstellbare Zahl im IEEE-754 32bit Standard? Geben Sie auch die Bits an.  
8bit 1-254 0 und 255 reserviert da overflow von 127  $\Rightarrow$  von -126 bis 127

### Addition

Rechnen Sie die folgenden Zahlen in IEEE-754 32bit Darstellung um und addieren Sie sie in dieser Darstellung miteinander. Runden Sie korrekt! Stellen Sie das Ergebnis sowohl im IEEE-754 Format als auch als Dezimalzahl dar.

- 592,183940 1. IEEE umwandeln
- 0,91213 2. addieren

auf den bei 9 5 3 sonst abbrechen  
bei  $n^2$  in IEEE 1 1 1  
0 1 1  
1 1 0

### Multiplikation

Rechnen Sie die folgenden Zahlen in IEEE-754 32bit Darstellung um und multiplizieren Sie sie in dieser Darstellung miteinander. Runden Sie korrekt! Stellen Sie das Ergebnis sowohl im IEEE-754 Format als auch als Dezimalzahl dar.

- 3981.1729 char 1 - char 2
- 2.91762 Maschine 1  $\cdot$  Maschine 2

Comma = Maschine + Maschine  $\approx$  nach 64 Kommastellen

## Aufgabe 2: Floating Point Rechner

register

Implementieren Sie einen Floating Point Rechner in Software. Zwei gegebene 32-bit IEEE-754 Zahlen (*operand1* und *operand2*) sollen addiert werden können. Das Ergebnis soll ebenfalls im 32-bit IEEE-754 Floating Point Format (\*) an die Adresse *result* geschrieben werden. Verwenden Sie keine Floating Point Register/Befehle. Denken Sie auch an Normalisierung, NaN und "Unendlichkeit". Sie können von folgender Signatur ausgehen (C-Datei wird gestellt):

`void calc_add(float operand1, float operand2, float* result);`  
 xmm0 + xmm1  
 adresse = 7th  
 pointer steht in rdi da andere calling convention

\*) Sie müssen sich nicht zwingend um korrekte Rundung kümmern. Auch die korrekte Behandlung von Charakteristik=0 als Sonderfall ist ein Bonus.

MOVQ eax, xmm0 schreibt 32 Bit von xmm in anderes 32 Bit Register  
 MOVQ ecx, xmm1  
 Ziel: MOV[rax], eax  
 Ref

VZ	char	Mankine
1 Bit	8 Bit	23 Bit

32 Bit

wird von Pc generiert (wird hinzugefügt)

$$15,25_{10} = +1111,01 \cdot 2^0 = \boxed{1,11101} \cdot 2^3 \quad \begin{matrix} 127+3 \\ 130 = \dots \end{matrix}$$

Festgelegt

Normalisiert sein  $\hat{=} 1,xxx \cdot 2^z$

char Exzees Darstellung

offset 127

10111111

addieren mit basis für char

VZ	char	Mankine
0		11101000... bis 23 Bit

Addition I 0 | 1000 0010 | 1110 1000..

II 1 | 1000 1010 | 1111 1000...

Bsp  $1,34 \cdot 10^3$   
 $+ 2,60 \cdot 10^5$

! erst gleicher Exponent

$$\approx 0,0134 \cdot 10^5$$

$$2,60 \cdot 10^5$$

Differenz

$$\begin{array}{r} 10000010 \\ - 10001010 \\ \hline = \dots = 8 \end{array}$$

Zahl mit kleineren Exponenten wird angepasst

Also I Manikine u. f. erhöhen

$\Rightarrow 0,000000011110100\dots$   
 immer wenn Zahl größer als 1  
 wird von Pc generiert

$$\begin{array}{l} II \ 1,11111000 \\ I - 0,0000000111101 \\ \hline = -1,111101100001100 \end{array}$$

VZ unterschiedlich  $\Rightarrow |x_1| \geq |x_2|$

$x_1$   
 $- x_2$   
 $=$  Umwandeln von  $x_1$  (Betrag größer - Betrag kleiner)

VZ	char	Manikine
1	10001010	111101100001100

Multiplizieren

$$\begin{array}{l} 1,34 \cdot 10^3 \\ 2,6 \cdot 10^4 \end{array} \quad \left. \begin{array}{l} 1, \\ 3 \\ 2 \end{array} \right\} \text{Exponenten egal}$$

$$\begin{array}{r} 1010 \cdot 0101 \\ \phantom{1010} 1010 \\ \phantom{1010} 0000 \\ \phantom{1010} 0000 \\ + \phantom{1010} 0000 \\ \hline 0110010 \end{array}$$

Bei 0 nur 0 schreiben  
Bei 1 linke Zahl abschreiben

→ char addieren  
Mantisse multi

1. Operanden holen

2. Operieren

rfid = VZ OP1  
rpd = char OP1  
rpid = Mantisse OP1  
esi = VZ OP2

char nach  $\leftarrow^{23}$  und den  $\rightarrow^1$

3. Exponent Vergleichen

3 Fälle: - gleich  
Op1 > Op2  
Op1 < Op2

8. explizite 1

9. Ergebnis Bauen

4. explizite 1 hinzufügen  
Mantisse +  $2^{23}$

10. zurückrechnen

5. exponenten angleichen

6. Rechnen 0 2 Fälle

- Vorzeichen =  
- Vorzeichen  $\neq$

7. Normalisieren 3 Fälle

- 1x, xx  $\leftarrow$  1 Rechts  
1, xxx  $\leftarrow$  nichts  
0, xxx einschieben

# Addition

Rechnen Sie die folgenden Zahlen in IEEE-754 32bit Darstellung um und addieren Sie sie in dieser Darstellung miteinander. Runden Sie korrekt! Stellen Sie das Ergebnis sowohl im IEEE-754 Format als auch als Dezimalzahl dar.

- 592,183940

1. IEEE umwandeln

591,27181

- 0,91213

2. addieren

$$592 = \frac{2^9 2^8 2^7 2^6 2^5 2^4 2^3 2^2 2^1 2^0}{10010100000} 2^0$$

$$0,183940 = \frac{2^{-1} 2^{-2} 2^{-3} 2^{-4} 2^{-5} 2^{-6} 2^{-7} 2^{-8} 2^{-9} 2^{-10} 2^{-11} 2^{-12} 2^{-13} 2^{-14} 2^{-15} 2^{-16} 2^{-17} 2^{-18} 2^{-19} 2^{-20}}{0,001010111100010101100010} 2^0$$

$$0,91213 = \frac{2^{-1} 2^{-2} 2^{-3} 2^{-4} 2^{-5} 2^{-6} 2^{-7} 2^{-8} 2^{-9} 2^{-10} 2^{-11} 2^{-12} 2^{-13} 2^{-14} 2^{-15} 2^{-16} 2^{-17} 2^{-18} 2^{-19} 2^{-20}}{0,11101001100000010011001} 2^0$$

$$0,183940 \cdot 2$$

$$0,36788 \cdot 2$$

$$0,73576$$

$$1,47152$$

$$0,91213$$

$$1,88608$$

$$-592,18394 = 1001010000,00101111000101101011 2^0$$

$$\Rightarrow 1,00101000000101111000101101011 2^9$$

$$127 + 9 = 136 = 10001000$$

$$-592,18394 = \text{Vchar}(136) \text{ Mantisse } I \text{ gerundet von 011 zu 10}$$

$$0,91213 - 592,18394 = -591,28181$$

$$\begin{array}{r} 592,18394 \\ 0,91213 \end{array} = \begin{array}{r} 1,0010100000,0010111100010110 \\ 0,0000000011101001100000 \end{array}$$

$$100100111101000101100110$$

$$\cdot 2^9$$

# Multiplikation

Rechnen Sie die folgenden Zahlen in IEEE-754 32bit Darstellung um und multiplizieren Sie sie in dieser Darstellung miteinander. Runden Sie korrekt! Stellen Sie das Ergebnis sowohl im IEEE-754 Format als auch als Dezimalzahl dar.

- 3981.1729 • = -11,615,5497
- -2.91762

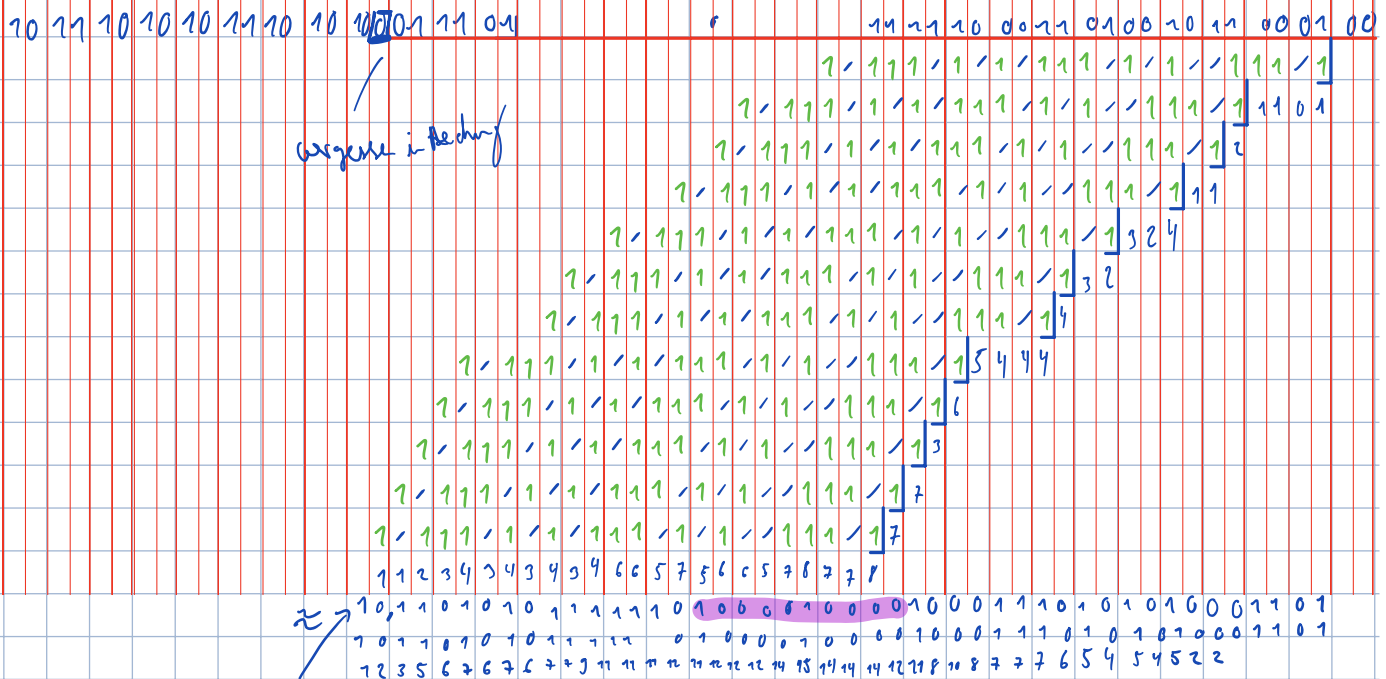
$$3981 = 1 \cdot 2^8 + 1 \cdot 2^7 + 1 \cdot 2^6 + 1 \cdot 2^5 + 0 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0$$

$$0,1729 = 0 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} + 0 \cdot 2^{-4} + 1 \cdot 2^{-5} + 0 \cdot 2^{-6} + 0 \cdot 2^{-7} + 0 \cdot 2^{-8} + 1 \cdot 2^{-9} + 0 \cdot 2^{-10} + 0 \cdot 2^{-11} + 0 \cdot 2^{-12} + 1 \cdot 2^{-13} + 1 \cdot 2^{-14} + 0 \cdot 2^{-15} + 0 \cdot 2^{-16} + 1 \cdot 2^{-17} + 0 \cdot 2^{-18} + 1 \cdot 2^{-19} + 0 \cdot 2^{-20}$$

$$-2,91762 = 1 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3} + 1 \cdot 2^{-4} + 0 \cdot 2^{-5} + 1 \cdot 2^{-6} + 0 \cdot 2^{-7} + 1 \cdot 2^{-8} + 1 \cdot 2^{-9} + 0 \cdot 2^{-10} + 0 \cdot 2^{-11} + 0 \cdot 2^{-12} + 1 \cdot 2^{-13} + 1 \cdot 2^{-14} + 0 \cdot 2^{-15} + 0 \cdot 2^{-16} + 0 \cdot 2^{-17} + 1 \cdot 2^{-18} + 1 \cdot 2^{-19} + 0 \cdot 2^{-20}$$

$$3981,1729 = 1 \cdot 2^8 + 1 \cdot 2^7 + 1 \cdot 2^6 + 1 \cdot 2^5 + 0 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 + 0 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} + 0 \cdot 2^{-4} + 1 \cdot 2^{-5} + 0 \cdot 2^{-6} + 0 \cdot 2^{-7} + 0 \cdot 2^{-8} + 1 \cdot 2^{-9} + 0 \cdot 2^{-10} + 0 \cdot 2^{-11} + 0 \cdot 2^{-12} + 1 \cdot 2^{-13} + 1 \cdot 2^{-14} + 0 \cdot 2^{-15} + 0 \cdot 2^{-16} + 1 \cdot 2^{-17} + 0 \cdot 2^{-18} + 1 \cdot 2^{-19} + 0 \cdot 2^{-20}$$

$$-2,91762 = 1 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3} + 1 \cdot 2^{-4} + 0 \cdot 2^{-5} + 1 \cdot 2^{-6} + 0 \cdot 2^{-7} + 1 \cdot 2^{-8} + 1 \cdot 2^{-9} + 0 \cdot 2^{-10} + 0 \cdot 2^{-11} + 0 \cdot 2^{-12} + 1 \cdot 2^{-13} + 1 \cdot 2^{-14} + 0 \cdot 2^{-15} + 0 \cdot 2^{-16} + 0 \cdot 2^{-17} + 1 \cdot 2^{-18} + 1 \cdot 2^{-19} + 0 \cdot 2^{-20}$$



addiere den exponenten für das  $1+1+11 = 13$

char 130