

Digitale Rechenanlagen

MARIÁN VAJTERŠIČ

Fachbereich Computerwissenschaften
Universität Salzburg
marian@cosy.sbg.ac.at
Tel.: 8044-6344

28. September 2017

Kapitel 2: Zahlendarstellung

Erste Digitalrechenanlagen wurden entwickelt, um einfache numerische Berechnungen durchzuführen.

[Heutzutage löst diese Aufgaben jeder Taschenrechner (sogar noch viel mehr) in wesentlich kürzerer Zeit.]

Das zentrale Problem beim Rechnen mit Computern:

Wie sollen Zahlen dargestellt werden?

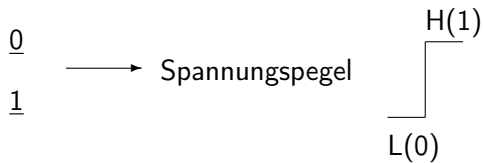
Wir wissen schon:

Die Repräsentation von Information ist gleichbedeutend mit der Codierung von Zeichen.

Es gibt aber folgende (hardwaremäßige) Einschränkung für die Repräsentation von Zahlen in Digitalrechenanlagen:

- ▶ Der Code für die Darstellung von Zahlen muss eine feste Länge haben!

Im Alltag benutzt man beim Rechnen fast immer Dezimalzahlen [d. h. diese Zahlen sind durch die Zeichen 0, 1, ..., 9 codiert], aber in Digitalrechenanlagen gibt es nur zwei Zustände:



Also Zahlen werden in Digitalrechenanlagen als Folge fester Länge von **0,1**-Zeichen dargestellt.

Die Aufgabe ist daher:

Eine Codierung $C : D \rightarrow B^n$ finden, die eine Untermenge D der Dezimalzahlen (Wörter) aus dem Alphabet $A = \{0, 1, \dots, 9\}$ auf B^n abbildet [d. h. man versucht, Dezimalzahlen binär mit n Stellen darzustellen].

Beispiel:

0 \rightarrow 1011

1 \rightarrow 1100

...

13 \rightarrow 1000

Jeder der Dezimalzahlen [aus der angegebenen Untermenge] wird ein eindeutiger Wert zugeordnet.

Mit Binärwörtern der Länge n können wir maximal 2^n verschiedene Dezimalzahlen codieren.

Prinzipiell könnten wir jeder Zahl aus der Untermenge von 2^n Dezimalzahlen eindeutig eine binäre n -Bit Zahl zuordnen, also eine tabellarische Zuordnung erzeugen. Diese wäre ineffizient, da wir dann zur Decodierung eine Tabelle mit 2^n Einträgen bräuchten!

Anstelle einer expliziten (tabellarischen) Codierung wird eine analytische Methode zur Konstruktion eines Binärcodes für Dezimalzahlen benutzt.

Einfache Binärdarstellung

[Darstellung einer Dezimalzahl im polyadischen Zahlensystem]

Die analytische Methode zur (De-)Codierung von Dezimalzahlen basiert auf folgender Definition:

Definition

[Darstellung einer Dezimalzahl im polyadischen Zahlensystem] In einem polyadischen Zahlensystem [d. h. die Basis kann mehrere Werte annehmen] mit der Basis $B \in \mathbb{N}$, $B > 1$ wird eine Dezimalzahl $m \in \mathbb{N}$ als Summe von gewichteten Potenzen von B folgendermaßen dargestellt:

$$\begin{aligned}(1) \quad m &= \sum_{i=0}^{n-1} b_i B^i = b_0 B^0 + b_1 B^1 + b_2 B^2 + \dots + b_{n-1} B^{n-1} = \\ &= b_0 + b_1 B + b_2 B^2 + \dots + b_{n-1} B^{n-1}.\end{aligned}$$

[Darstellung einer Dezimalzahl im polyadischen Zahlensystem]

Die Gewichte b_i sind die Ziffern des Zahlensystems zur Basis B [$b_i \in \{0, 1, \dots, B - 1\}$], wobei die Ziffer b_0 ganz rechts steht und die Zahl dann lautet: $b_{n-1}b_{n-2} \dots b_1b_0$.

Für uns ist das **dyadische** Zahlensystem (**Binärsystem**), aufgrund der Anwendung in Digitalrechenanlagen, von spezieller Bedeutung.

Beispiel: Einige wichtige Zahlensysteme

Zahlensystem	Basis	Ziffern
Binär	2	0,1
Oktal	8	0,1,2,3,4,5,6,7
Dezimal	10	0,1,2,3,4,5,6,7,8,9
Hexadezimal	16	0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F

[Darstellung einer Dezimalzahl im polyadischen Zahlensystem]

Um Verwechslungen zu vermeiden [z. B. kann 9 in einer Dezimalzahl, sowie auch in einer Hexadezimalzahl als Ziffer auftreten] wird das verwendete Zahlensystem [wenn nötig] als Index angegeben:

$1000_{(10)}$: Tausend im Dezimalsystem

$1000_{(2)}$: Eine „andere“ Zahl im Binärsystem

Beispiel: Zahlenumwandlung ins Dezimalsystem

$$127_{(10)} = 7 \cdot \underbrace{10^0}_{B^0} + 2 \cdot \underbrace{10^1}_{B^1} + 1 \cdot \underbrace{10^2}_{B^2} = 127_{(10)}$$

$$10101_{(2)} = 1 \cdot 2^0 + 0 \cdot 2^1 + 1 \cdot 2^2 + 0 \cdot 2^3 + 1 \cdot 2^4 = 21_{(10)}$$

$$3147_{(8)} = 7 \cdot 8^0 + 4 \cdot 8^1 + 1 \cdot 8^2 + 3 \cdot 8^3 = 1639_{(10)}$$

$$1A2C_{(16)} = 12 \cdot 16^0 + 2 \cdot 16^1 + 10 \cdot 16^2 + 1 \cdot 16^3 = 6700_{(10)}$$

Zahlenumwandlung

[Übergang zwischen Zahlensystemen]

Wie kann man Dezimalzahlen in ein anderes Zahlensystem umwandeln?

→ Sukzessive Division

[Die Gleichung (1) wird durch die Basis B dividiert.]

$$\frac{m - b_0}{B} = b_1 + b_2B + \dots + b_{n-1}B^{n-2} = m'$$

→ $m = m'B + b_0$ → Die Ziffer b_0 ist der Rest bei der Division von m durch B .

Mit m' verfährt man in gleicher Weise. Der nächste Rest ist die Ziffer b_1 , usw., bis $m' = 0$.

Zahlenumwandlung

[Übergang zwischen Zahlensystemen]

Beispiel: (Sukzessive Division) $47_{(10)} \stackrel{?}{=} x_{(2)} \rightarrow B = \underline{2}$

m		m'		b_i
$47:\underline{2}$	=	23	(+Rest $\underline{1}$)	$\underline{b_0}$
23		11		$\underline{b_1}$
11		5		$\underline{b_2}$
5		2		$\underline{b_3}$
2		1		$\underline{b_4}$
1		0		$\underline{b_5}$

$$\begin{aligned}
 47_{(10)} &= 1 \quad 0 \quad 1 \quad 1 \quad 1 \quad 1_{(2)} \\
 &\quad b_5 \quad b_4 \quad b_3 \quad b_2 \quad b_1 \quad b_0 \\
 &= 1 \cdot 2^0 + 1 \cdot 2^1 + 1 \cdot 2^2 + 1 \cdot 2^3 + 1 \cdot 2^5 = \\
 &= 1 + 2 + 4 + 8 + 32 = 47_{(10)}
 \end{aligned}$$

Zahlenumwandlung

[Wie groß ist die größte darstellbare Zahl zur Basis B mit n Stellen?]

$$m_{\max} = b_{\max} \sum_{i=0}^{n-1} B^i = b_{\max} (1 + B + B^2 + \dots + B^{n-1}).$$

$(1 + B + \dots + B^{n-1})$: Geometrische Reihe mit Quotient q :

$$\sum_{i=0}^{n-1} q^i = \frac{1 - q^n}{1 - q}.$$

Zahlenumwandlung

[Wie groß ist die größte darstellbare Zahl zur Basis B mit n Stellen?]

In unserem Fall: $q = B$, $b_{\max} = \underline{B - 1}$

[Da ein Zahlensystem zur Basis B aus den Ziffern von $0, 1, 2, \dots, B - 1$ besteht \rightarrow die größte Ziffer ist $B - 1$.]

Dann:

$$\underline{m_{\max}} = b_{\max} \sum_{i=0}^{n-1} B^i = (B - 1) \frac{1 - B^n}{1 - B} = \underline{B^n - 1}.$$

[Anders: Die größte darstellbare Zahl zur Basis B mit n Stellen ist $B^n - 1$, da B^n nicht mehr dargestellt werden kann. $B = 2 \rightarrow 2^n - 1$]

Darstellung negativer Zahlen

Negative Zahlen

[Darstellung negativer Zahlen]

Bisher haben wir nur positive [ganze] Zahlen betrachtet.

Nun: Betrachten die [binäre] Darstellung negativer Zahlen.

Codierung negativer Zahlen

Wie codiert man negative Zahlen?

Durch Voranstellen eines Minuszeichens:

$$-17 = (-1) \cdot 17 = \ominus 17.$$



Minus-Zeichen verwendet

Wie kann man das bei Binärzahlen machen?

$$\ominus 101$$



Das Minuszeichen muss integriert werden.

[d. h. auch mit 0/1 darstellbar sein.]

Naive Codierung des Minuszeichens

- ▶ MSB: Most Significant Bit b_{n-1} [Die wichtigste Stelle: ganz links]
- ▶ LSB: Least Significant Bit b_0 [ganz rechts]

Wenn MSB=0: positive Zahl
 MSB=1: negative Zahl.

Führt zu einem Problem:

$$(\underline{1})_{(10)} + (\underline{1})_{(10)} = \underline{1}01_{(2)} \oplus \underline{1}01_{(2)} = \underline{0}10_{(2)} = \underline{\pm}2_{(10)}$$

$$\begin{array}{r} 101 \quad [-1] \\ 101 \quad [-1] \\ \hline \overset{(1)}{(1)}\overset{(1)}{0}10 \quad [+2] \end{array}$$

↓

Problem: Definierte Anzahl der Stellen.

Verbesserung: Geeignete Wahl des MSB.

Die Idee

Ein Vorzeichenbit nicht nur zur Indikation des Vorzeichens zu verwenden sondern auch zur Darstellung der Zahlen.

Wir wissen: Die größte darstellbare [natürliche] Binärzahl mit n Stellen ist $2^n - 1$.

Wir verwenden diese n Stellen für die Darstellung von Zahlen aus dem Intervall $[-2^{n-1} + 1, 2^{n-1} - 1]$, d. h. insgesamt $2^n - 1$ Zahlen mit 2 Null-Darstellungen
und

interpretieren das MSB [hier das Bit links außen] als Vorzeichenbit.

Für einfache Addition [Subtraktion] von Zahlen in dieser Darstellung muss man das Gewicht des MSB richtig wählen. Eine Möglichkeit ist die **1-Komplement-Darstellung**.

1-Komplement-Darstellung

Die **1-Komplement-Darstellung** einer Zahl $z \in \mathbb{Z}$ ist folgendermaßen definiert:

$$z = -b_{n-1}(2^{n-1} - 1) + \sum_{k=0}^{n-2} b_k 2^k .$$

Wenn $b_{n-1} = 1$ [d. h. MSB ist gesetzt], dann ist z negativ $[-0$ inklusive].

[Da die größte darstellbare Zahl mit $n - 1$ Zeichen 2^{n-1} ist, gilt $\sum_{k=0}^{n-2} b_k 2^k \leq 2^{n-1} - 1$.]

Wenn $b_{n-1} = 0$ [d. h. MSB ist nicht gesetzt], dann fällt der negative Teil weg und z ist positiv $[+0$ inklusive].

1-Komplement-Darstellung

Also, wenn

1. $b_{n-1} = 1$: $-2^{n-1} + 1 \leq z \leq 0$
2. $b_{n-1} = 0$: $0 \leq z \leq 2^{n-1} - 1$.

Darstellung von z :

1. Für z sucht man das positive z' sodass $z' = z + (2^{n-1} - 1)$ und stellt dieses in einfacher Binärdarstellung mit $b_{n-1} = 1$ [negatives Vorzeichen] dar.
2. Für z ergibt sich die einfache binäre Darstellung mit $b_{n-1} = 0$ [positives Vorzeichen].

Bemerkung:

Für $z = 0$ ergeben sich zwei Darstellungen:

1. $z = -0$
2. $z = +0$

1-Komplement-Darstellung

Beispiel: [8-Bit-Zahlen im 1-Komplement]

- $z_{(10)} = 1 \rightarrow$ Fall 2. $z'_{(10)}$ wird nicht berechnet.

$$k_1(1)_{(2)} = \underset{\text{Vorzeichenbit}}{\underline{0}} 0000001$$

- $z_{(10)} = -1 \rightarrow$ Fall 1. $z'_{(10)}$ wird berechnet:

$$z'_{10} = -1 + \overbrace{(2^{(8-1)} - 1)}^{z_{(10)} + (2^{n-1} - 1)} = -1 + (2^{(8-1)} - 1) = -1 + (2^7 - 1) = -1 + 127 = 126$$

dann einfach binär codiert $[126_{(10)} = 1111110_{(2)}]$ und das negative Vorzeichen hinzugefügt.

$$k_1(-1)_{(2)} = \underset{\text{Vorzeichenbit}}{\underline{1}} 1111110$$

1. Darstellung der Zahl inklusive Vorzeichenbit = 1.
2. Einfache Codierung inklusive Vorzeichenbit = 0.

1-Komplement-Darstellung

Beispiel: [8-Bit-Zahlen im 1-Komplement]

- $z_{(10)} = 127 \rightarrow \text{Fall 2.}$ $z'_{(10)}$ wird nicht berechnet.

$$k_1(127)_{(2)} = \underline{0}1111111$$

- $z_{(10)} = -127 \rightarrow \text{Fall 1.}$ $z'_{(10)} = -127 + 127 = 0.$

$$k_1(-127)_{(2)} = \underline{1}0000000$$

1-Komplement-Darstellung

Beispiel: [Zwei 8-Bit-Zahlen im 1-Komplement]

- $z_{(10)} = 0 \rightarrow \text{Fall 2.}$

$$k_1(0)_{(2)} = \underline{00000000}$$

- $z_{(10)} = -0 \rightarrow \text{Fall 1.} \quad z'_{(10)} = -0 + 127 = 127 \xrightarrow{1.}$

$$k_1(-0)_{(2)} = \underline{11111111}$$

Bemerkung:

Es ist ersichtlich, dass für $k_1(z)_{(2)} = b_{n-1}b_{n-2}\dots b_0$ gilt:

$$k_1(-z)_{(2)} = \neg b_{n-1} \neg b_{n-2} \dots \neg b_0,$$

wobei $\neg 0 = 1$ und $\neg 1 = 0$, also das binäre Komplement darstellt.

1-Komplement-Darstellung

Bemerkung: Für eine Zahl $z_{(1)}$ im 1-er Komplement gilt: $\bar{z}_{(1)} = -z_{(1)}$.

Beweis.

Zu zeigen ist, dass $\bar{z}_{(1)} + z_{(1)} = 0$.

$$\bar{z}_{(1)} = -\overline{b_{n-1}}(2^{n-1} - 1) + \sum_{k=0}^{n-2} \overline{b_k} 2^k$$

$$z_{(1)} = -b_{n-1}(2^{n-1} - 1) + \sum_{k=0}^{n-2} b_k 2^k$$

$$\bar{z}_{(1)} + z_{(1)} = -\underbrace{(\overline{b_{n-1}} + b_{n-1})}_1 (2^{n-1} - 1) + \sum_{k=0}^{n-2} \underbrace{(\overline{b_k} + b_k)}_1 2^k$$

$$= -1(2^{n-1} - 1) + \sum_{k=0}^{n-2} 2^k = -(2^{n-1} - 1) + (2^{n-1} - 1) = 0.$$



1-Komplement-Darstellung

Beispiel:

$$z_{(1)} = 0001 = \underline{1}$$

$$z_{(1)} = 1011 = -\underline{4}$$

$$z_{(1)} = 1111 = -\underline{0}$$

$$\bar{z}_{(1)} = \overline{0001} = 1110 = -7 + 4 + 2 = -\underline{1}$$

$$\bar{z}_{(1)} = \overline{1011} = 0100 = \underline{4}$$

$$\bar{z}_{(1)} = \overline{1111} = 0000 = +\underline{0}$$

Addition im 1-Komplement

Wenn man addieren kann, kann man auch subtrahieren.

[Da die Subtraktion $a - b$ die Rückführung auf die Addition $a + (-b)$ erlaubt.]

Zu beachten: Das Bit b_n [links vom MSB] kann bei einer Addition verändert werden!

Bei der Addition im 1-Komplement müssen wir dieses Bit [im Falle von $b_n = 1$] an der Stelle b_0 hinzuaddieren [Einserrücklauf].

Es entsteht also ein Übertrag [eine $n + 1$ -stellige Zahl]. Dieser wird zu den ursprünglichen Zahlen mit jeweils verdoppeltem MSB [b_{n-1}] addiert.

Wenn dabei $b_n \neq b_{n-1}$ ist, dann kommt es zu einem Überlauf [Overflow], d. h. die Summe wird die maximal darstellbare Zahl überschreiten \rightarrow Falsches Ergebnis!

Beispiel: $[-0 \oplus -0]$

$$\begin{array}{rcccccccccccc}
 -0 & & & & & & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
 -0 & \oplus & & & & & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
 & & & & & & \hline
 & & & & & & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 \text{ Übertrag} \\
 & & & & & & \underbrace{}_{b_n} & & & & & & &
 \end{array}$$

$$\begin{array}{rcccccccccccc}
 -0 & & & & & & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
 -0 & \oplus & & & & & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
 & \oplus & & & & & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\
 & & & & & & \hline
 & & & & & & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
 & & & & & & \underbrace{}_{b_{n+1}} & \underbrace{}_{b_n} & \underbrace{}_{b_{n-1}} & & & & &
 \end{array}$$

Die Summanden
mit doppeltem MSB
[b_n an der Stelle b_0]

b_{n+1} ist uninteressant.

$b_n = b_{n-1} \rightarrow$ es ist kein Überlauf aufgetreten.

\rightarrow Das Ergebnis ist ok.

Für das Ergebnis sind nur die Stellen $b_{n-1} \dots b_0$ relevant!

$\rightarrow -0 \oplus -0 = 11111111 = -0$

Beispiel: $[100 \oplus 28]$: 8-Bit-Darstellung

$$100_{(10)} + 28_{(10)} = 128_{(10)}$$

Diese Zahl ist außerhalb des Intervalls $[-127, 127]$ [die Zahlen, die durch 8 Bit-Stellen im 1-Komplement noch darstellbar sind].

$$\begin{array}{rcccccccc}
 100 & & & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 \\
 28 & \oplus & & 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
 \hline
 & & \underbrace{0}_{b_n} & \underbrace{1}_{b_{n-1}} & 0 & 0 & 0 & 0 & 0 & 0 & 0
 \end{array}$$

$\rightarrow b_n \neq b_{n-1}$: Überlauf! $[010000000 = -127 + 0 = -127_{(1K)} \neq 128!]$

Hier ist $b_n = 0$, also müssen wir den Übertrag nicht zu den Summanden [mit verdoppeltem MSB] addieren – es würde keine Veränderung des Resultats bringen.

Beispiel: $[-77 \oplus -52]$

$$\begin{array}{rcccccccc}
 -77 & & & & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\
 -52 & \oplus & & & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 1 \\
 \hline
 & & 1 & & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1
 \end{array}
 \quad \text{Übertrag}$$

$\underbrace{\hspace{1.5cm}}_{b_n}$

$$\begin{array}{rcccccccc}
 -77 & & & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\
 -52 & \oplus & & 1 & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 1 \\
 & \oplus & & & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\
 \hline
 & & 1 & 1 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0
 \end{array}$$

$\underbrace{\hspace{1.5cm}}_{b_{n+1}} \quad \underbrace{\hspace{1.5cm}}_{b_n} \quad \underbrace{\hspace{1.5cm}}_{b_{n-1}}$

Die Summanden
 mit doppeltem MSB
 $[b_n \text{ an der Stelle } b_0]$

b_{n+1} ist uninteressant.

$b_n \neq b_{n-1} \rightarrow \text{Überlauf}$

\rightarrow Das Ergebnis ist ungültig (und falsch).

$$01111110 = 64 + 32 + 16 + 8 + 4 + 2 = 126 \neq -129!$$

Überlauf

Überlauf: Feststellung mittels Vergleich von b_n mit b_{n-1} . Das Ergebnis dieses Vergleichs wird als Flag [Bit] in der CPU in einem Statusregister angeführt.

Statusregister 0/1 \leftarrow

Overflow Flag						
---------------	--	--	--	--	--	--

Ist das Flag = 1, dann ist ein Überlauf aufgetreten. Das Steuerprogramm kann so einfach erkennen, ob das Resultat gültig ist oder nicht.

Nachteile der 1-Komplement Darstellung

Die Nachteile der 1-Komplement Darstellung sind:

- ▶ der Einserrücklauf
- ▶ die doppelte Darstellung der 0

Diese beiden Nachteile werden eliminiert durch die **2-Komplement-Darstellung**.

2-Komplement-Darstellung

Für eine Zahl $z \in \mathbb{Z}$ gilt:

$$z = -b_{n-1}2^{n-1} + \sum_{k=0}^{n-2} b_k 2^k.$$

Die 0 hat nur eine Darstellung: 00000000 (für $n = 8$).

[Es muss $b_{n-1} = 0$ sein $\rightarrow b_{n-2} = \dots = b_0 = 0$, denn wenn $b_{n-1} = 1$
 $\rightarrow -2^{n-1} + (2^{n-1} - 1) = -1 < 0$,
weil $2^{n-1} - 1$ die größte darstellbare positive Zahl mit $n - 1$ Bits ist.]

Das ermöglicht die Darstellung einer zusätzlichen Zahl

$$\begin{aligned} \rightarrow \quad & z \in [-2^{n-1}, 2^{n-1} - 1] \\ n = 8 : \quad & z \in [-128, 127]. \end{aligned}$$

2-Komplement-Darstellung

Darstellung von z : Wir unterscheiden zwei Intervalle (Analogie zum 1-Komplement):

1. $b_{n-1} = 1$: $-2^{n-1} \leq z < 0$
2. $b_{n-1} = 0$: $0 \leq z \leq 2^{n-1} - 1$

1. Für z sucht man das positive $z' = z + 2^{n-1}$ und stellt dieses in einfacher Binärdarstellung **mit** $b_{n-1} = 1$ [negatives Vorzeichen] dar.
2. Hier wird z in einfacher binärer Darstellung **mit** $b_{n-1} = 0$ [positives Vorzeichen] dargestellt.

2-Komplement-Darstellung

Beispiel: [8-Bit-Zahlen im 2-Komplement]

$z_{(10)}$	$z'_{(10)}$	$k_2(z)_{(2)}$	
1	-	00000001	
-1	<u>127</u>	11111111	$\rightarrow z' = z + 2^{n-1} = -1 + 128 = \underline{127}$ (Vorzeichenbit: 1)
127	-	01111111	
-127	<u>1</u>	10000001	$\rightarrow z' = z + 2^{n-1} = -127 + 128 = \underline{1}$ (Vorzeichenbit: 1)
-128	0	10000000	

2-Komplement-Darstellung

Bildungsvorschrift für negative Zahlen im 2-Komplement:

$$\begin{aligned} \text{wenn } k_2(z)_{(2)} &= b_{n-1} \quad b_{n-2} \quad \dots \quad b_0 \rightarrow \\ k_2(-z)_{(2)} &= \neg b_{n-1} \quad \neg b_{n-2} \quad \dots \quad \neg b_0 \\ &\oplus \quad \quad \quad 0 \quad \quad 0 \quad \quad \dots \quad 1 \end{aligned}$$

Beispiel:

$$\begin{array}{rcll} 127 & = & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ -127 & = & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ & \oplus & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ & & \hline & & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{array}$$

Addition im 2-Komplement

Hier wird der Einserrücklauf schon bei der Komplementbildung berücksichtigt
→ die Addition/Subtraktion wird vereinfacht.

Vorgang:

Verdopplung von MSB [b_{n-1}] der beiden Summanden, danach wird die Addition durchgeführt.

Falls $b_n \neq b_{n-1}$: Überlauf aufgetreten.

Tritt kein Überlauf auf [also $b_n = b_{n-1}$], so ist das Ergebnis direkt an den Stellen der Summe $b_{n-1} \dots b_1 b_0$ abzulesen [b_{n+1} bleibt unberücksichtigt].

Beispiel: Addition im 2-Komplement

$$17_{(10)} + (-2)_{(10)} = 17_{(10)} - 2_{(10)} = 15_{(10)}$$

$$\begin{array}{rcll} 2_{(10)} & : & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ -2_{(10)} & : & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 \end{array} \quad [z' = -z + 128 = 126]$$

Verdopplung und Addition:

$$\begin{array}{rcll} 17_{(10)} & & \underline{0} & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ -2_{(10)} & \oplus & \underline{1} & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 \\ \hline & & \underbrace{1}_{b_{n+1}} & \underbrace{0}_{b_n} & \underbrace{0}_{b_{n-1}} & 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{array}$$

$b_n = b_{n-1} = 0 \rightarrow$ kein Überlauf, b_{n+1} unberücksichtigt

Das Vorzeichen ist positiv ($b_{n-1} = 0$) und $(b_{n-2}b_{n-3}\dots b_0)_{(2)} = \underline{0001111}_{(2)}$.

\rightarrow Das Ergebnis ist $\underline{00001111}_{(2)} = 15_{(10)}$.

Addition im 1- und 2-Komplement: Zusammenfassung

[1-Komplement]

- ▶ Die Summanden im 1-Komplement darstellen
- ▶ Addition der Summanden [beide \underline{n} Bit lang]

$$\begin{array}{rccccccc}
 \underline{n} & n-1 & & & & & 0 \\
 \downarrow & & & & & & \\
 & x & & \dots\dots & & & \cdot \\
 & x & & \dots\dots & & & \cdot \\
 \bullet & \cdot & & \dots\dots & & & \cdot
 \end{array}$$

- ▶ a) $\bullet = 1$: Übertrag
- ▶ b) $\bullet = 0$: kein Übertrag

Addition im 1- und 2-Komplement: Zusammenfassung

[1-Komplement]

- a) **Übertrag:** Verdopplung des Vorzeichenbits $[n - 1]$ [die Summanden sind dann $n + 1$ Bits lang], Addition von $0 \dots 01$ [aus dem Übertrag] zu den Summanden.

$$\begin{array}{rcccc}
 n + 1 & \underline{n} & \underline{n - 1} & & 0 \\
 & \times & \times & \dots & . \\
 & \times & \times & \dots & . \\
 & 0 & 0 & \dots & 1 \\
 \hline
 \bullet & \odot & \otimes & \dots & .
 \end{array}$$

- [Position $n + 1$] belanglos.

Wenn $\odot = \otimes$ [Position n und $n - 1$]: ok, kein Überlauf aufgetreten.

Wenn $\odot \neq \otimes$: Überlauf, Ergebnis ist ungültig.

Ergebnis: $\otimes \dots$ [also eine n Bit lange Zahl].

- b) **Kein Übertrag:** Wie a), aber ohne Einserrücklauf-Phase [also gleich die Summanden mit verdoppeltem Vorzeichenbit addieren]

Probe: $(Ergebnis)_{(10)} = b_0 \cdot 2^0 + b_1 \cdot 2^1 + \dots + b_{n-2} \cdot 2^{n-2} - b_{n-1} \cdot (2^{n-1} - 1)$

Addition im 1- und 2-Komplement: Zusammenfassung

[2-Komplement]

- ▶ Die Summanden im 2-Komplement darstellen
- ▶ Addition wie nach b) (vorige Folie)

Probe: $(Ergebnis)_{(10)} = b_0 \cdot 2^0 + b_1 \cdot 2^1 + \dots + b_{n-2} \cdot 2^{n-2} - b_{n-1} \cdot 2^{\underline{n-1}}$

Multiplikation von Binärzahlen

Wir können auf die bekannten Methoden für das Dezimalsystem zurückgreifen.

- ▶ Klassisches Verfahren:

Die Multiplikation wird mit den einzelnen Zahlen des Multiplikators ausgeführt und die Einzelergebnisse entsprechend ihrer Wertigkeit [Position] addiert.

Beispiel: Binärmultiplikation

Das Produkt $7_{(10)} \cdot 3_{(10)}$ in 4-Bit Binärdarstellung:

[illegible]

Also: Die Binärmultiplikation ist auf mehrere Additionen zurückzuführen, weil bei den einzelnen Multiplikationen nur zwei Teilergebnisse möglich sind - der Multiplikand (\otimes) oder die 0 (\odot). Diese Eigenschaft wird bei der Konstruktion von Multiplizierern [Hardware für die Multiplikation] ausgenützt.

Multiplikation von Binärzahlen

Problem für die Hardware:

Die Anzahl der Stellen des Produkts ist $p = d + k$ [$5 = 3 + 2$].

[d : Stellenanzahl des Multiplikanden
 k : des Multiplikators
(0-Stellen am Anfang nicht gezählt)]

Multiplikation von Binärzahlen

[Ägyptische Multiplikation]

▶ Ägyptische Multiplikation:

Vorgang:

- ▶ a) man schreibt die beiden zu multiplizierenden Zahlen nebeneinander
- ▶ b) verdoppelt in jeder neuen Zeile die linke Zahl
- ▶ c) halbiert die rechte Zahl [ganzzahlig]
- ▶ d) streicht jene Zeilen, in denen die rechte Zahl gerade ist
- ▶ e) addiert die verbleibenden linken Zahlen

Beispiel:

a)	22		21
b)	44	c)	10 d)
b)	88	c)	5
b)	176	c)	2 d)
b)	352	c)	1
e)	462		

Multiplikation von Binärzahlen

[Ägyptische Multiplikation, Eigenschaften der Methode]

Neben der Addition wird nur die Verdopplung und Halbierung der Zahlen benötigt. Diese beiden Operationen lassen sich aber mit Binärzahlen sehr einfach realisieren.

- ▶ Verdopplung: eine Verschiebung um eine Stelle nach links [shift left, SHL]
- ▶ Halbierung: eine Verschiebung um eine Stelle nach rechts [shift right, SHR]

Multiplikation von Binärzahlen

[Halbierung/Verdopplung]

Beispiel:

$$\begin{aligned} 27 \cdot 2 &= \text{SHL}(\underline{11011}) &= \underline{110110} &= 54 \\ 33 / 2 &= \text{SHR}(\underline{100001}) &= \underline{10000} &= 16 \end{aligned}$$

[Die freiwerdende Stelle wird mit dem Bit 0 belegt.]

Bemerkung: Diese beiden Operationen gehören zu Befehlen einer gängigen Maschinensprache [Assemblerprogrammierung].

Multiplikation zur Basis $B \neq 2, 10$

► $23_{(8)} \cdot 16_{(8)}$

$$\begin{array}{r} 22_{(8)} \\ 14_{(8)} \\ \underline{23_{(8)}} \\ \textcolor{red}{4}\textcolor{blue}{1}2_{(8)} = 4 \cdot 8^2 + 1 \cdot 8^1 + 2 \cdot 8^0 = 266_{(10)} \end{array}$$

Probe:

$$23_{(8)} = 19_{(10)}$$

$$16_{(8)} = 14_{(10)}$$

$$19 \cdot 14 = 266_{(10)}$$

$$6 \cdot 3 = 18 = 2 \cdot 8 + 2 = 22_{(8)}$$

$$6 \cdot 2 = 12 = 1 \cdot 8 + 4 = 14_{(8)}$$

$$1 \cdot 23_{(8)} = 23_{(8)}$$

$$\textcolor{blue}{1}: 3 + 4 + 2 = 9 = \textcolor{green}{1}\textcolor{blue}{1}_{(8)}$$

$$\textcolor{red}{4}: \textcolor{green}{1} + 2 + 1 = 4 = 4_{(8)}$$

Multiplikation zur Basis $B \neq 2, 10$

► $41_9 \cdot 37_9$

$$\begin{array}{r}
 \underline{317_9} \\
 133_9 \\
 \hline
 1647_9
 \end{array}
 \quad
 \begin{array}{l}
 7 \cdot 1 = 7 = \underline{7_9} \\
 7 \cdot 4 = 28 = \underline{31_9}
 \end{array}$$

$$1647_9 = 1 \cdot 9^3 + 6 \cdot 9^2 + 4 \cdot 9 + 7 = 729 + 486 + 36 + 7 = 1258_{(10)}$$

Probe: $41_9 = 37_{10}$
 $37_9 = 34_{10}$

$$37 \cdot 34 = 1258_{10}$$

Fazit: Multiplikation wie üblich
 aber

- Die Darstellung von Teilprodukten zur Basis B
- Die Addition der Teilprodukte auch zur Basis B.

Division von Binärzahlen

► $111111_2 : 1001_2 = ?$

$$\begin{array}{r}
 111111_2 : 1001_2 = 111_2 \\
 \underline{-1001} \\
 01101 \\
 \underline{-1001} \\
 01001 \\
 \underline{-1001} \\
 0
 \end{array}$$

Probe: $\begin{array}{rclclcl} 11111 & = & 2^6 - 1 & = & 63_{10} \\ 1001 & = & & = & 9_{10} \end{array}$

$$63 : 9 = 7_{(10)} = 111_{(2)}$$

Division von Binärzahlen

► $1001,11_2 : 1,101_2 = ?$

$$1001,11_2 : 1,101_2 =$$

$$1001\underline{1}10_2 : 1101_2 = 110_2$$

$$\begin{array}{r} 1001110 \\ -1101 \\ \hline 001101 \\ -1101 \\ \hline 00000 \end{array}$$

Probe:
$$\begin{array}{rclcl} 1001,11_2 & = & 9 + 0,5 + 0,25 & = & 9,75_{(10)} \\ 1,101_2 & = & 1 + 0,5 + 0,125 & = & 1,625_{(10)} \end{array}$$

$$\begin{array}{l} 9,75 : 1,625 \\ 9750 : 1625 = 6 = (110)_{(2)} \\ -\underline{9750} \\ 0 \end{array}$$

Fazit: Division wie üblich mit Operationen binäre Multiplikation und binäre Subtraktion.

Division im Zahlensystem zur Basis $B \neq 2, 10$

► $11011_{(3)} : 11_{(3)} = ?$

$$\begin{array}{r} \underline{11011}_{(3)} : 11_{(3)} = \underline{1001} \\ -11_3 \\ \hline 00011 \\ -11 \\ \hline 00 \end{array}$$

Probe: $11011_{(3)} = 1 + 3 + 27 + 81 = 112_{(10)}$
 $11_{(3)} = 1 + 3 = 4_{(10)}$
 $112_{(10)} : 4_{(10)} = 28_{(10)} = 1001_{(3)}$

Division im Zahlensystem zur Basis $B \neq 2, 10$

► $1017_8 : 37_8 = ?$

$$\begin{array}{r} \underline{101}7_8 : 37_8 = \underline{2}1_8 \\ -76_8^* \\ \hline 0\underline{3}7_8 \square \\ -37_8 \\ \hline 00 \end{array}$$

$$\begin{array}{rcl} \text{Probe: } 1017_8 & = & 7 + 8 + 512 = 512_{(10)} \\ 37_8 & = & 31_{(10)} \end{array}$$

– Abschätzung für 2 : $101_8 : 37_8 = 10,1 : 3,7 \approx 10_8 : 4_8 = 8 : 4 = 21_8$

– * Erklärung für 76_8 : durch $2_8 \cdot 37_8$

$$\begin{array}{rcl} 2 \cdot 7 & = & 14 & = & 16_8 \\ 2 \cdot 3 & = & 6 & = & 6_8 \\ & & & & \hline & & & & 76_8 \end{array}$$

– □ Erklärung für 37_8 : Durch Subtraktion von $101_8 - 76_8$.
(!Also im Oktalsystem: 6 und wie viel ist $11_8 = 9_{10}$? Antwort: 3)

$$\blacktriangleright 23144,3_{(5)} : 44,321_{(5)} = ?$$

$$23144,3_{(5)} : 44,321_{(5)} = 232,4 \dots_{(5)}$$

$$\begin{array}{r}
 23144300_{(5)} : \underline{44}321_{(5)} = \underline{232},4_{(5)} \\
 \underline{-144142}_{(5)} \\
 0323010_{(5)} \\
 \underline{-244013}_{(5)} \\
 234420_{(5)} \\
 \underline{-144142}_{(5)} \\
 402230 \\
 \underline{-343334} \\
 3341
 \end{array}$$

$$\begin{array}{lll}
 \text{Probe: } 23144,3_{(5)} & = 0,6 + 4 + 4 \cdot 5 + 1 \cdot 25 + 3 \cdot 125 + 2 \cdot 625 = \\
 & = 1674,6_{(10)} \\
 44,321_{(5)} & = 2 \cdot 0,04 + 3 \cdot 0,2 + 4 + 4 \cdot 5 = \\
 & = 24,68_{(10)}
 \end{array}$$

$$\begin{aligned}
 1674,6 : 24,68 &\approx 67,85 &= 2 + 3 \cdot 5 + 2 \cdot 25 + 4 \cdot 0,2 = \\
 & &= 67,8 = \underline{\underline{232,4_{(5)}}}
 \end{aligned}$$

-Abschätzung für 2 : $23_5 : 4,4_5 \approx 13_{10} : 4,8_{10} = 2$

-Erklärung * für $32301_{(5)}$

$$\begin{array}{rclcl}
 2 & + & \underline{1} & = & 3 \\
 4 & + & \underline{0} & = & 4 \\
 1 & + & \underline{3} & = & 4 \\
 4 & + & \underline{2} & = & 11_5 = 6 \\
 (4 + 1) & + & \underline{3} & = & 13_5 = 8 \\
 (1 + 1) & + & \underline{0} & = & 2
 \end{array}$$

-Abschätzung für 3 : $32_5 : 4,4_5 \approx 17 : 4,8 = 3$

-Abschätzung für 2 : $23_5 : 4,4_5 \approx 13 : 4,8 = 2$

-Abschätzung für 4 : $40_5 : 4,4_5 \approx 20 : 4,8 = 4$

Fazit:

- Abschätzung aufgrund weniger (2 ?) Ziffern
- Multiplikation zur Basis B
- !Subtraktion zur Basis B

Darstellung rationaler Zahlen

Darstellung rationaler Zahlen

In numerischen Berechnungen ist die Verwendung rationaler Zahlen $q \in \mathbb{Q}$

[Kommazahlen] unumgänglich.

[Eine rationale Zahl ist durch einen Bruch einer ganzen Zahl $[\mathbb{Z}]$ und einer natürlichen Zahl $[\mathbb{N}]$ darstellbar.]

- ▶ natürliche Zahlen $\mathbb{N} = \{1, 2, \dots\}$
- ▶ ganze Zahlen $\mathbb{Z} = \{\dots - 2, -1\} \cup \{0\} \cup \mathbb{N} = \{\dots - 2, -1, 0, 1, 2, \dots\}$

Im Rechner können nur endliche Zahlen dargestellt werden. Eine Erweiterung des polyadischen Zahlensystem ermöglicht die Darstellung [gewisser] rationaler Zahlen:

$$q = \sum_{k=-m}^{n-1} b_k B^k = b_{-m} B^{-m} + b_{-m+1} B^{-m+1} + \dots + b_0 + \dots + b_{n-1} B^{n-1}$$

Festpunktdarstellung

Jede rationale Zahl q lässt sich als Summe einer ganzen Zahl $u \in \mathbb{Z}$ und einer rationalen Zahl v ($1 > v \geq 0$) mit $q = u + v$ angeben:

$$\begin{aligned} \blacktriangleright \quad -\frac{4}{1} &= (-4.0)_{(10)} = -4 + \underline{0.0} = u + \underline{v} \\ \blacktriangleright \quad -\frac{37}{10} &= (-3.7)_{(10)} = -4 + \underline{0.3} = u + \underline{v} \\ \blacktriangleright \quad \frac{2625}{1000} &= (2.625)_{(10)} = 2 + \underline{0.625} = u + \underline{v} \end{aligned}$$

Einfaches Darstellungsformat

Festpunktdarstellung

$$\blacktriangleright q_{(10)} = (\underbrace{u_n u_{n-1} \dots u_0}_u \odot \underbrace{v_{-1} \dots v_{-m}}_v)_{(2)}$$

wobei

- ▶ u_n : das Vorzeichenbit (0 für nichtnegative Zahl, 1 für negative Zahl)
- ▶ $u = u_{n-1} \dots u_0$: eine ganze n -Bit Zahl
- ▶ $v = v_{-1} \dots v_{-m}$: eine rationale m -Bit Zahl

Für diese Darstellung verwendet man die **Schreibweise**:

- ▶ $(+\underline{n}.\underline{m})$: n positives Vorzeichenbit, n Vorkommaziffern, m Nachkommaziffern
- ▶ $(-\underline{n}.\underline{m})$: n negatives Vorzeichenbit, n Vorkommaziffern, m Nachkommaziffern.

Beispiel:

$$(0 \underbrace{1010}_u . \underbrace{11}_v)_{(2)} \quad \text{hat das Fixpunktdarstellungsformat } +4.2.$$

Beispiele: Festpunktdarstellung

1. Umwandlung der Zahl $(0.125)_{(10)}$ in eine Binärzahl in $+\underline{2}.\underline{4}$ -Festpunktdarstellung [d. h. positives Vorzeichenbit, 2 Vorkomma- ($n = 2$) und 4 Nachkomma-Ziffern ($m = 4$)].

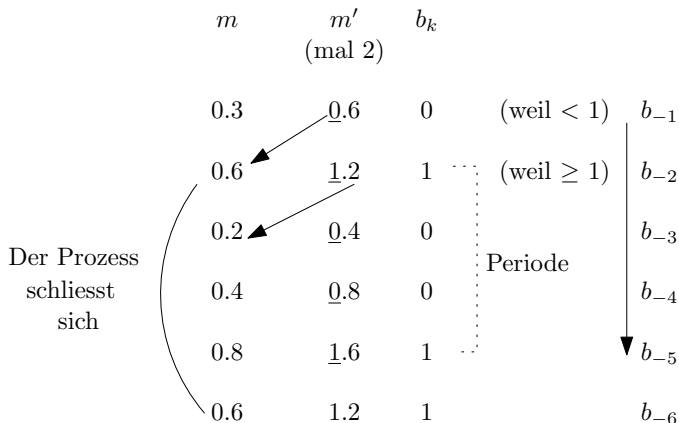
$$(0.125)_{(10)} = 0 + \frac{1}{8} = 0 + 2^{-3} = \underbrace{0}_{\text{Vorzeichen}} \underbrace{00}_{2 \text{ Stellen}} . \underbrace{0010}_{4 \text{ Stellen}}$$

Hier sehr einfach, da die Zahl durch eine einzige [negative] Zweierpotenz darstellbar ist. Die Umwandlung ist auch mittels sukzessiver Division für den Vorkommateil und sukzessiver Multiplikation für den Nachkommateil möglich (siehe folgende Beispiele).

Beispiele: Festpunktdarstellung

2. Umwandlung der Zahl $-1.3_{(10)}$ in eine Binärzahl in (-2.8) -Festpunktdarstellung.

Zuerst die Umwandlung des Nachkommateils 0.3 (sukzessive Multiplikation):



Beispiele: Festpunktdarstellung (Fortsetzung Beispiel 2.)

- ▶ Die erste abgespaltete Ziffer b_k entspricht b_{-1} .
- ▶ Der Prozess terminiert nur dann, wenn die rationale Zahl als Summe negativer Zweierpotenzen darstellbar ist.
- ▶ Bei 0.3 ist dies nicht der Fall. Hier erhalten wir eine periodische Binärzahl.

[Der Prozess schließt sich bei 0.6.]

→ daher die Periode 1001.

Weil diese periodische Binärzahl auf $m = \underline{8}$ Stellen beschränkt ist:

$$\text{▶ } 0.3 = 0.0\underbrace{1001}_{b_{-2} \dots b_{-5}} 100$$

- ▶ Unter Beachtung des Vorzeichens ergibt sich

$$\text{▶ } -1.3_{(10)} = (\underbrace{1}_{\text{Vorzeichen}} \quad \underbrace{01}_{2 \text{ Stellen}} \quad . \underbrace{01001100}_{8 \text{ Stellen}})_{(2)}$$

Beispiele: Festpunktdarstellung

3. Umwandlung der Zahl $0.25_{(10)}$ in $+0.2$ -Festpunktdarstellung.

$$\begin{array}{rcccl}
 m & m' & & & \\
 \hline
 \blacktriangleright & 0.25 & \underline{0.5} & 0 & b_{-1} \quad | \\
 & 0.5 & \underline{1.0} & 1 & b_{-2} \quad \downarrow
 \end{array}
 \rightarrow 0.25_{(10)} = (\underbrace{0}_{\text{Vorzeichen}} . \underbrace{01}_2)_{(2)}$$

Vorzeichen, 0 Vorkommastellen, 2 Nachkommastellen

Normierung

Die Zahl im Beispiel 2. können wir nicht exakt darstellen, auch wenn wir beliebig viele Stellen zur Verfügung hätten. Das führt zu Fehlern in numerischen Berechnungen in digitalen Rechenanlagen.

Eine Möglichkeit für die Verbesserung der Berechnungen mit Zahlen in Festpunktdarstellung ist die Normierung:

Die Nachkommastellen werden vor jedem Rechenschritt solange nach links geschoben, bis $b_{-1} \neq 0$ wird. Dies führt meist zu genauerer Berechnung, jedoch müssen alle Normierungen aufgezeichnet werden, um das Resultat korrekt denormieren zu können.

Zahlenumwandlung $B \rightarrow B^k$

B : Basis des Zahlensystems (z. B. $B = 2, 7, \dots$)
 B^k : Basis eines „verwandten Zahlensystems“ (z. B. $2^3 = 8, 7^3 = 343$)

Umwandlung: Angenommen d ist eine rationale Zahl zur Basis B

$$\begin{array}{ccc}
 (d_{n-1} \dots d_4 d_3 d_2 d_1 d_0 & \cdot & d_{-1} d_{-2} d_{-3} d_{-4} \dots d_{-m})_{(B)} \\
 \leftarrow & & \rightarrow
 \end{array}$$

Die Basis (neu): $B^{\underline{k}}$ $\underline{k} = 2, 3, \dots$

Zahlenumwandlung $B \rightarrow B^k$

- ▶ Für den Vorkommateil: Beginnend von $\underline{d_0}$ werden Teilketten von Bits der Länge \underline{k} genommen und zur Basis B^k dargestellt \rightarrow neue Ziffern für den Vorkommateil.
- ▶ Für den Nachkommateil: Beginnend mit $\underline{d_{-1}}$ werden analog k -Bit-Ketten zur Basis B^k dargestellt \rightarrow neue Ziffern für den Nachkommateil.

Zahlenumwandlung $B \rightarrow B^k$

Begründung

- Vorkommateil transformiert [zur Basis B^k] durch sukzessive Division durch B^k [d. h. Schieben um k Stellen nach rechts].

Ermittlung der ersten Ziffer:

$$(d_{n-1} \dots d_0)_B: B^k = d_{n-1} \dots d_k \cdot d_{k-1} \dots d_0$$

→ Rest der Division [d. h. die erste Ziffer z_0] zur Basis B^k ist $(d_{k-1} \dots d_0)_B = (z_0)_{B^k}$

Also die erste k -Bit-Kette (Vorkomma, $\leftarrow \cdot$) dargestellt zur Basis B^k .

- Nachkommateil transformiert [zur Basis B^k] durch sukzessive Multiplikation mit B^k [d. h. Schieben um k Stellen nach links].

Ermittlung der ersten Ziffer:

$$\cdot d_{-1} d_{-2} \dots d_{-k} d_{-k-1} \dots d_{-m} \cdot B^k = d_{-1} d_{-2} \dots d_{-k} \cdot d_{-k-1} \dots d_{-m}$$

→ Die Zahl $(d_{-1} \dots d_{-k})_B$ ist die erste Nachkommaziffer $(z_{-1})_{(B^k)}$.

Zahlenumwandlung $B \rightarrow B^k$

Beispiel 1: $B = 2$

$$(1000000, 1011)_2 \rightarrow \begin{array}{ll} d_{n-1} = 1 & n = 8 \\ d_{-m} = 1 & m = 4 \end{array}$$

- oktale Darstellung (Umwandlung zur Basis $8 = 2^3 \rightarrow \underline{k = 3}$)

$$\begin{array}{ccccccc} 1 & | & 000 & | & 000, & 101 & | & 1^{00} & = (100, 54)_8 \\ \downarrow & & \downarrow & & \downarrow & & \downarrow & & \downarrow \\ 1_8 & 0_8 & 0_8 & & 0_8 & 5_8 & 4_8 & & \\ & & & & \downarrow & & & & \\ & & & & (000)_2 & = 0_8 & & & \end{array}$$

- hexadezimale Darstellung (neue Basis = $16 = 2^4 \rightarrow \underline{k = 4}$)

$$\begin{array}{ccccccc} 100 & | & 0000, & \underline{1011} & = (40, B)_{16} & = & 4 \cdot 16^1 + 0 + \frac{11}{16} = 64 \frac{11}{16}_{(10)} \\ \downarrow & & \downarrow & & \downarrow & & \\ 4_{16} & 0_{16} & B_{16} & & & & \\ & & & & = & & (\underbrace{1000000}_2, \underbrace{1011}_2) \\ & & & & & & \downarrow \quad \downarrow \\ & & & & 64 & \frac{1}{2} + \frac{1}{8} + \frac{1}{16} = \frac{11}{16} \end{array}$$

- Umwandlung zur Basis $64 = 2^6 \rightarrow \underline{k = 6}$

$$\begin{array}{ccccccc} 1 & | & \underline{000000}, & \underline{1011}^{00} & = 10, 44_{(64)} \\ \downarrow & & \downarrow & & \downarrow \\ 1_{64} & 0_{64} & & & 44_{64} \end{array}$$

Zahlenumwandlung $B \rightarrow B^k$

Beispiel 2: $B = 7$

214_7

- Darstellung zur Basis $B^2 = 49$ ($k = 2$)

$$\begin{array}{c} 2 \overline{) 14_7} \\ \downarrow \downarrow \\ 2_{49} \quad 11_{49} \end{array} \Rightarrow 214_7 = 2 \cdot 11_{49} = 11 + 2 \cdot 49^1 = 109_{10} = 2 \cdot 7^2 + 1 \cdot 7^1 + 4 \cdot 7^0 = 214_7$$

Beispiel 3: $B = 3$

$22102, 2110_3$

- Basis $9 = 3^2 \rightarrow k = 2$

$$\begin{array}{ccccc} 2 & | & 21 & | & 02, 21 & | & 10 \\ \downarrow & & \downarrow & & \downarrow & & \downarrow \\ 2_9 & & 7_9 & & 2_9 & & 7_9 & & 3_9 \end{array} \Rightarrow 272, 73_{(9)}$$

Zahlenumwandlung $B1 \rightarrow B2$ (Allgemein)

$B1, B2$: Basen zweier unterschiedlicher Zahlensysteme.

Umwandlung $B1 \rightarrow B2$ mittels sukzessiver Division [Vorkommateil] und sukzessiver Multiplikation [Nachkommateil].

→ Division (Multiplikation) mit $B2$ aber zur Basis $B1$ dargestellt, also mit $(B2)_{B1}$ [damit erfolgt die Operation im Zahlensystem zur Basis $B1$].

Die Ziffern der neuen transformierten Zahl zur Basis $B2$ sind:

- ▶ Rest bei Division (in $B2$ -Darstellung) [Vorkommateil]
- ▶ Ganzzahliger Teil des Produkts ($B2$ -Format) [Nachkommateil]

Beispiel 1: $B1 = 3$ $B2 = 7$ $110_3 \stackrel{?}{=} x_7$

$$110_3 : 7 = 110_3 : 21 = 1$$

$$\begin{array}{r} -21_3 \\ \hline \end{array}$$

$$12_3 = \text{Rest}$$

$$1_3 : 21_3 = 0$$

$$1_3 = \text{Rest}$$

$$12_3 = 5_7$$

$$1_3 = 1_7$$



$$\Rightarrow x_7 = 15_7 = 12_{10} = 110_3$$

Beispiel 2: $B1 = 5$ $B2 = 4$ $4013_5 \stackrel{?}{=} x_4$

$$4013_5 : 4 = 4013_5 : 4_5 = 1002_5$$

$$\begin{array}{r} -4_5 \\ \hline 0013_5 \\ -13_5 \\ \hline 0_5 = Rest \end{array}$$

$$0_5 = 0_4$$

$$1002_5 : 4_5 = 111_5$$

$$\begin{array}{r} -4_5 \\ \hline 10_5 \\ -4_5 \\ \hline 12_5 \\ -4_5 \\ \hline 3_5 = Rest \end{array}$$

$$3_5 = 3_4$$

$$111_5 : 4_5 = 12_5$$

$$\begin{array}{r} -4_5 \\ \hline 21_5 \\ -13_5 \\ \hline 3_5 = Rest \end{array}$$

$$3_5 = 3_4$$

$$12_5 : 4_5 = 1_5$$

$$\begin{array}{r} -4_5 \\ \hline 3_5 = Rest \end{array}$$

$$3_5 = 3_4$$

$$1_5 : 4_5 = 0$$

$$1_5 = Rest$$

$$1_5 = 1_4$$

⇒

$$x_4 = 13330_4 =$$

$$= 0 + 12 + 3 \cdot 16 + 3 \cdot 64 + 256 = 508_{(10)}$$

$$4013_5 = 3 + 5 + 4 \cdot 125 = 508_{(10)}$$

Beispiel 3: $B1 = 5$ $B2 = 10$ $0,4_5 \stackrel{?}{=} x_{10}$

$$0,4_5 \cdot 10 = 0,4_5 \cdot 20_5 =$$

$$\begin{array}{r} 13_5 \\ 0_5 \\ \hline \end{array}$$

$\underline{13}, 0_5 \Rightarrow \underline{8}_{(10)}$ ist die erste Nachkomma Ziffer zur Basis 10

$$\Rightarrow 0,4_5 = 0,8_{10}$$

Probe: $0,4_5 = \frac{4}{5} = \underline{0,8_{10}} = x_{10}$

Beispiel 4: $B1 = 4$ $B2 = 6$ $0,312_4 \stackrel{?}{=} x_6$

$$0,312_4 \cdot 6 = 0,312 \cdot 12_4$$

$$\begin{array}{r} 1230_4 \\ \hline 11,010_4 \end{array} \Rightarrow 11_4 = 5_6$$

$$0,010_4 \cdot 12_4 =$$

$$\begin{array}{r} 020_4 \\ \hline 0,120_4 \end{array} \Rightarrow 0_4 = 0_6$$

$$0,120_4 \cdot 12_4 =$$

$$\begin{array}{r} 300_4 \\ \hline 0,2100_4 \end{array} \Rightarrow 0_4 = 0_6$$

$$0,21_4 \cdot 12_4 =$$

$$\begin{array}{r} 102_4 \\ \hline 0,312_4 \end{array} \Rightarrow 0_4 = 0_6$$

.... \Rightarrow Periode

$$\Rightarrow x_6 = 0,\overline{5000}_6$$

$\Rightarrow 5000$ ist die Periode

Probe: $0,312_4 = \frac{3}{4} + \frac{1}{16} + \frac{2}{64} \approx \frac{13}{16}$

$$0,\overline{5000}_6 \approx \frac{5}{6}$$

$$\frac{13}{16} - \frac{5}{6} = \frac{78-80}{96} = \frac{-2}{96} = \frac{-1}{48} \rightarrow \text{ok.}$$

Gleitpunktdarstellung rationaler Zahlen

Effiziente Darstellung rationaler Zahlen durch Gleitpunktdarstellung
(Floating-Point-Representation)

Grundidee:

- ▶ $q = \pm M \cdot B^C$
- ▶ B : Basis [z. B.: $B = 2$, $B = 10$]
- ▶ M : Mantisse
- ▶ C : Charakteristik

Beispiel: (Fix-Point versus Floating-Point-Darstellung)

- weniger Zeichen, $q = 2^{-16}$
Wahl $B = 2$, $M = 1.0$

Festpunkt: $q_{(2)} = (0.00 \dots \underbrace{1}_{b_{-16}})_{(2)} \rightarrow 16 \text{ Nachkommastellen}$

Floating-Point: Für $B = 2$ und $M = 1.0$ benötigt man nur 6 Bit:

5 Bit für C [$16_{(10)} = (10000)_2$]

1 Bit für M [$1.0_{(10)} = (1)_2$]

+ 1 Bit [Vorzeichenbit für C].

Eine andere Wahl von M [z. B. $M = 2.0$] würde zu einer anderen Darstellung führen.

Floating Point Standards

IEEE-Standards [„EI-TRIPL-I“] - Institute of Electrical and Electronic Engineers [weltweit renommierter Fachverein]. Diese Standards gelten für die Repräsentation von Zahlen in modernen Rechnern.

Format:

- ▶ 32 Bit Wortlänge: single precision
- ▶ 2 · 32-Bit-Wörter [= 64 Bit]: double precision
- ▶ 4 · 32-Bit-Wörter [oder 2 · 64-Bit-Wörter = 128 Bit]: quadruple precision

Wir zeigen das 32- und das 64-Bit-Format.

IEEE-32-Bit-Format

v	C	m
1	8	23

v : das Vorzeichenbit [$v = 0$ oder $v = 1$]

C : die 8-stellige Charakteristik

$$C = (c_7 c_6 \dots c_0)_{(2)} = c_7 \cdot 2^7 + c_6 \cdot 2^6 + \dots + c_0 \cdot 2^0, \quad c_i \in \{0, 1\}$$

$$\rightarrow C = 0, 1, \dots, 2^8 - 1 = 255$$

[also eine übliche binäre Darstellung]

IEEE-32-Bit-Format

v	C	m
1	8	23

m : 23-stelliger Binärbruch $m_{22}m_{21}...m_0$ der Mantisse

$M = (1.m)_{(2)}$, $m_i \in \{0, 1\}$ [mit Komma in der Darstellung ganz links]

d. h. $M = (\underline{1}.m)_{(2)} = (\underline{1}.r)_{(10)} = (\underline{1} + r)_{(10)}$, $0 \leq \underline{r} < 1$

weil

$$\begin{aligned}
 M &= (\underline{1}.m)_{(2)} = (\underline{1} + \underbrace{m_{22} \cdot 2^{-1} + m_{21} \cdot 2^{-2} + \dots + m_0 \cdot 2^{-23}}_r)_{(10)} \\
 &= (\underline{1} + r)_{(10)} = (\underline{1}.r)_{(10)}
 \end{aligned}$$

also $r_{(10)} = (m)_{(2)}$

IEEE-32-Bit-Format

Die Darstellung [Repräsentation] einer rationalen Zahl q im 32-Bit-IEEE-Format:

$$q_{\text{IEEE-32-Bit}} = \left(\underbrace{v}_{1} \underbrace{c_7 \dots c_0}_{8} \underbrace{m_{22} m_{21} \dots m_0}_{23} \right)_{\text{IEEE-32-Bit}}$$

Die Zahl q , repräsentiert durch die obige Darstellung, wird dezimal interpretiert als

$$q_{(10)} = (-1)^v \cdot 2^{(C-127)_{(10)}} \cdot (1.m)_{(2)} = (-1)^v \cdot 2^{C'_{(10)}} \cdot (1.r)_{(10)},$$

wobei $C' = (C - 127)_{(10)}$ die Interpretation von $C \in [0 \dots 255]$ in Excess-127-Code ist.

IEEE-32-Bit-Format

Also:

Die Charakteristik C [repräsentiert in $q_{\text{IEEE-32-Bit}}$ durch 8 Bits $c_7 \dots c_0$] wird in der Dezimaldarstellung $q_{(10)}$ von q interpretiert als $C' = (C - 127)_{(10)}$.

Konkret:

C	C'
0	-127
1	-126
...	...
127	0
...	...
255	128

Ausnahmefälle

Wenn $C = \overbrace{(00000000)}^8 = 0_{(10)}$, d. h. $C' = (-127)_{(10)}$, wird M ohne 1 in der Mantisse [also $M = (\underline{0}.m)_{(2)} = (0.r)_{(10)}$] interpretiert.

1. Speziell [Darstellung von $0_{(10)}$]:

Wenn $C = 00000000$ und $m = \underbrace{00000000000000000000000}_{23}$, dann ist die

dargestellte Zahl $q_{(10)} = \underline{0}$

[0 hat zwei Darstellungen: $v = 0 : +0$, $v = 1 : -0$]

2. $C = 00000000 (= 0)$, $m \neq 0$ [„denormalized numbers to fill in gap around zero“]
 $\rightarrow q_{(10)} = (-1)^v \cdot 2^{-126} \cdot (0.m)_{(2)} = (-1)^v \cdot 2^{-126} \cdot (0.r)_{(10)}$

Ausnahmefälle

$$3. \ C = 11111111 \ (= 255), \ m = \underbrace{000000000000000000000000}_{23} \ (= 0)$$

$$\rightarrow q_{(10)} = (-1)^v \cdot \infty \ [v = 0 : +\infty, v = 1 : -\infty]$$

$$4. \ C = 11111111 \ (= 255), \ m \neq 0$$

$\rightarrow q_{(10)} = NaN$ [Not a Number], symbolisiert ungültige Ergebnisse
[z. B. Division durch 0].

Beispiele: Umwandlung von $q_{\text{IEEE-32-Bit}}$ in die Dezimalform

► Interpretation der 32-Bit-Binärzahl

$$q_{\text{IEEE-32-Bit}} = \overset{v}{1} \overset{C}{10000111} \overset{m}{101000000000000000000000} \text{ im Dezimalsystem.}$$

- $v = 1$ [\rightarrow negatives Vorzeichen]
- $C = (10000111)_{(2)} = (135)_{(10)}$
- $M = (1.m)_{(2)} = (1.\underbrace{10100\dots 0}_m)_{(2)} = 1 + \underbrace{2^{-1} + 0 + 2^{-3} + 0\dots}_r = (1.\underline{625})_{(10)}$

$$\rightarrow q_{(10)} = (-1)^1 \cdot 2^{135-127} \cdot 1.625 = -2^8 \cdot 1.625 = (-416)_{(10)}$$

- Wäre $C = (01111000)_{(2)}$ [$= 120_{(10)}$], dann ist $C' = C - 127 = (-7)_{(10)}$ und die Floating-Point-Darstellung wäre $-2^{-7} \cdot 1.625 \approx (-.0127)_{(10)}$.

- Falls $C = (00000000)_{(2)}$ wäre \rightarrow [laut Ausnahme]

$$q_{(10)} = (-1)^1 2^{-126} \cdot \underline{0.625} \approx -7.34 \cdot 10^{-39}$$

[also: 1 [Vorkommastelle] ist in der Mantisse weggelassen]

Beispiele: Umwandlung von $q_{\text{IEEE-32-Bit}}$ in die Dezimalform

► Interpretation der Zahl

$(0\ 01111000\ \underline{1}11110000000000000000000)_{\text{IEEE-32-Bit}}$ im Dezimalsystem.

- $v = (0)_{(2)}$
- $C = (01111000)_{(2)} = 2^6 + 2^5 + 2^4 + 2^3 = (120)_{(10)}$
 $\rightarrow C' = C - 127 = 120 - 127 = (-7)_{(10)}$
- $r_{(10)} = (.m)_{(2)} = \underline{1} \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3} + 1 \cdot 2^{-4} + 1 \cdot 2^{-5} = (0.96875)_{(10)}$

$$\rightarrow q_{(10)} = (-1)^0 \cdot 2^{-7} \cdot 1.96875 \approx 0.015$$

Umwandlung von Dezimalzahlen ins IEEE-Format

1. $q_{(10)}$ in die gewünschte Form bringen:

$$\square \quad q_{(10)} = (-1)^{v(2)} \cdot 2^{C'_{(10)}} \cdot (1.m)_{(2)}$$

Dazu wandeln wir $q_{(10)}$ zunächst in die Binärform um:

$$q_{(10)} = q_{(2)} = (\pm b_\ell b_{\ell-1} \dots b_0 . b_{-1} b_{-2} \dots b_k \dots)_{(2)}$$

[durch sukzessive Division (Vorkommateil) von $q_{(10)}$ und sukzessive Multiplikation (Nachkommateil)], wobei $b_\ell = 1$

$$\rightarrow q_{(10)} = (\pm 1) \cdot 2^{(\ell)_{(10)}} \cdot (1 . b_{\ell-1} \dots b_0 b_{-1} b_{-2} \dots b_k \dots)_{(2)}$$

Aus \square folgt:

- ▶ $v = 0$ (+) oder $v = 1$ (-)
- ▶ $C' = (\ell)_{(10)}$
- ▶ $(1.m)_{(2)} = (1 . b_{\ell-1} b_{\ell-2} \dots b_k \dots)_{(2)}$

Umwandlung von Dezimalzahlen ins IEEE-Format

2. aus der gewünschten Form \square das entsprechende IEEE-32- bzw. 64-Bit-Format generieren:

z. B. 32 Bit $q_{\text{IEEE-32-Bit}} = (v \ C \ m)_{\text{IEEE-32-Bit}}$

- ▶ $v = v_{(2)}$ (1 Bit)
- ▶ $C = (C' + 127)_{(10)} = (\ell + 127)_{(2)}$ (8 Bit)
- ▶ $m = m_{(2)}$ (23 Bit)

Umwandlung von Dezimalzahlen ins IEEE-Format

Beispiel: Darstellung von $(29.0)_{(10)}$ im IEEE-32-Bit-Format

1. Umwandlung von $q_{(10)} = (29.0)_{(10)}$ in die gewünschte Form

$$\square q_{(10)} = (-1)^{v_{(2)}} \cdot 2^{C'_{(10)}} \cdot (1.m)_{(2)}$$

Zunächst die Darstellung von $q_{(10)}$ in Binärform:

29	:	2	=	14		1
14	:	2	=	7		0
7	:	2	=	3		1
3	:	2	=	1		1
1	:	2	=	0		1

$$(29.0)_{(10)} = (11101.0)_{(2)} = (+b_4b_3b_2b_1b_0)_{(2)} \rightarrow \ell = (4)_{(10)}$$

$$\square (29.0)_{(10)} = (-1)^0 \cdot 2^4 \cdot (1.1101)_{(2)}$$

$$\rightarrow v = 0, \quad C' = 4_{(10)}, \quad (1.m)_{(2)} = (1.1101)_{(2)}.$$

Umwandlung von Dezimalzahlen ins IEEE-Format

Beispiel: Darstellung von $(29.0)_{(10)}$ im IEEE-32-Bit-Format (Fortsetzung)

2. Aus der gewünschten Form \square das IEEE-32-Bit-Format für $(29.0)_{(10)}$ ableiten.

- ▶ $v = 0_{(2)}$ (1 Bit)
- ▶ $C = (C' + 127)_{(10)} = (4 + 127)_{(10)} = (131)_{(10)} = (10000011)_{(2)}$ (8 Bit)
- ▶ $m = (11010 \dots 0)_{(2)}$ (23 Bit)

$$\begin{aligned} \rightarrow q_{\text{IEEE-32-Bit}} &= (v \ C \ m)_{\text{IEEE-32-Bit}} \\ &= (0 \ 10000011 \ 11010000000000000000000)_{\text{IEEE-32-Bit}} \end{aligned}$$

Umwandlung von Dezimalzahlen ins IEEE-Format

Beispiel: Darstellung von $(0.5)_{(10)}$ im IEEE-32-Bit-Format

1. Gewünschte Form

Die binäre Form:

$$(0.5)_{(10)} = (0.1)_{(2)} \quad 0.5 \cdot 2 = 1.0 \quad 1 = b_{-1}$$

$$\rightarrow \ell = (-1)_{(10)}.$$

$$\rightarrow \text{Die gewünschte Form: } (0.5)_{(10)} = (-1)^0 \cdot 2^{-1} \cdot (1.0)_{(2)}$$

$$\rightarrow v = 0, \quad C' = (-1)_{(10)}, \quad (1.m)_{(2)} = (1.0)_{(2)}$$

2. $\blacktriangleright v = 0$

$$\blacktriangleright C = (C' + 127)_{(10)} = (-1 + 127)_{(10)} = (126)_{(10)} = (01111110)_{(2)}$$

$$\blacktriangleright m = (0 \dots 0)_{(2)} \text{ (23 Bit)}$$

$$(0.5)_{(10)} = (0 \ 01111110 \ 000000000000000000000000)_{\text{IEEE-32-Bit}}$$

Umwandlung von Dezimalzahlen ins IEEE-Format

Beispiel: Darstellung von $(-1.0)_{(10)}$ im IEEE-SP-Format

1. Gewünschte Form

$$(-1.0)_{(10)} = (-1)^1 \cdot 2^0 \cdot (1.0)_{(2)} \rightarrow$$

$$v = 1, \quad \ell = C' = 0, \quad .m = 0 \dots 0$$

2.

$$v = 1 \quad C = C' + 127 = 127 = 011 \dots 1 \quad m = \underbrace{0 \dots 0}_{23}$$

$$(-1.0)_{(10)} = 1 \ 01111111 \ 000000000000000000000000 \quad \circledast$$

Bemerkung:

Darstellung von $(1.0)_{(10)}$ wie \circledast nur $v=0$

Umwandlung von Dezimalzahlen ins IEEE-Format

Beispiel: Darstellung von $(-47.7)_{(10)}$ im IEEE-32-Bit-Format

m				m'	b_i		
Periode		<u>47</u>	: 2 =	23	1	b_0	↑
		23	: 2 =	11	1	b_1	
		11	: 2 =	5	1	b_2	
		5	: 2 =	2	1	b_3	
		2	: 2 =	1	0	b_4	
		1	: 2 =	0	1	b_5	
	[<u>0.7</u>	· 2 =	1.4	1	b_{-1}	
		0.4	· 2 =	0.8	0	b_{-2}	
		0.8	· 2 =	1.6	1	b_{-3}	
		0.6	· 2 =	1.2	1	b_{-4}	
		0.2	· 2 =	0.4	0	b_{-5}	
		0.4	· 2 =	0.8	0	b_{-6}	
		0.8	· 2 =	1.6	1	b_{-7}	
		0.6	· 2 =	1.2	1	b_{-8}	↓
				\vdots	\vdots		

Umwandlung von Dezimalzahlen ins IEEE-Format

Beispiel: Darstellung von $(-47.7)_{(10)}$ im IEEE-32-Bit-Format (Fortsetzung)

1. Gewünschte Form:

$$\begin{aligned}
 (-47.7)_{(10)} &= (\underbrace{b_5 b_4 b_3 b_2 b_1 b_0}_{101111} . \underbrace{b_{-1} b_{-2} b_{-3} b_{-4} b_{-5} \dots}_{10110\dots})_{(2)} \\
 &= (-1) \cdot 2^5 \cdot (1.b_4 b_3 b_2 b_1 b_0 b_{-1} b_{-2} b_{-3} \dots)_{(2)} \\
 &= (-1) \cdot 2^5 \cdot (1.01111 \underbrace{1}_{b_{-1}} \underline{01100} \dots)_{(2)}
 \end{aligned}$$

- 2.
- ▶ $v = (1)_{(2)}$ (1 Bit)
 - ▶ $C' = \ell = (5)_{(10)} \rightarrow C = 127 + 5 = (132)_{(10)} = (10000100)_{(2)}$ (8 Bit)
 - ▶ $m = 011111011001 \dots = b_{-1} b_{-2} \dots b_k \dots = b_4 b_3 b_2 \dots$ (23 Bit)

$$(-47.7)_{(10)} = (1 \ 10000100 \ 011111 \ \underline{0110} \ \underline{0110} \ \underline{0110} \ \underline{0110} \ 0)_{\text{IEEE-32-Bit}}$$

Umwandlung von Dezimalzahlen ins IEEE-Format

[Die größte positive Single-Precision-Zahl im IEEE-Format: 32 Bit]

Positive Zahl:

- ▶ $v = 0$
- ▶ $C = (1111111\underline{0})_{(2)} = \underline{(254)}_{(10)}$
[weil $C = 11111111$ und $m \neq 0$ für den Ausnahmefall 4. reserviert ist]
- ▶ $M = (1.m)_{(2)} = (1.\underbrace{111\dots 1}_{23\times})_{(2)}$

$$\begin{aligned}
 q_{\text{IEEE-32-Bit}} &= (0 \ 11111110 \ 111111111111111111111111) \\
 [\text{largest positive}] &= (-1)^0 \cdot 2^{(C-127)_{(10)}} \cdot (1.m)_{(2)} \\
 &= 2^{(254-127)_{(10)}} \cdot (1.\underbrace{111\dots 1}_{23\times})_{(2)} \\
 &= 2^{127} \cdot (2 - 2^{-23})_{(10)} \approx (3.4028 \cdot 10^{38})_{(10)}
 \end{aligned}$$

Umwandlung von Dezimalzahlen ins IEEE-Format

[Die größte positive Single-Precision-Zahl im IEEE-Format: 32 Bit (Fortsetzung)]

Beweis.

$$\text{für } (1.\underbrace{111\dots 1}_{23\times})_{(2)} = (2 - 2^{-23})_{(10)}$$

$$\begin{aligned} (1.\underbrace{111\dots 1}_{23\times})_{(2)} &= 1 + 2^{-1} + 2^{-2} + 2^{-3} + \dots + 2^{-23} = \\ &= 2^0 + 2^{-1} + 2^{-2} + 2^{-3} + \dots + 2^{-23} = \sum_{i=0}^{23} \left(\frac{1}{2}\right)^i \end{aligned}$$

Es gilt für die Summe der geometrischen Reihe: $\sum_{i=0}^n q^i = \frac{1-q^{n+1}}{1-q} \rightarrow$

$$\sum_{i=0}^{23} \left(\frac{1}{2}\right)^i = \frac{1 - \left(\frac{1}{2}\right)^{24}}{1 - \frac{1}{2}} = \frac{1 - \left(\frac{1}{2}\right)^{24}}{\frac{1}{2}} = 2(1 - 2^{-24}) = (2 - 2^{-23})_{(10)}.$$

Die 32-Bit-Darstellung erlaubt nur eine eingeschränkte Genauigkeit \rightarrow daher definiert man auch ein 64-Bit-IEEE-Format (Double Precision). □

IEEE-64-Bit-Format

Die binäre Darstellung [Repräsentation]:

$$q_{\text{IEEE-64-Bit}} = \underbrace{v}_{1} \underbrace{c_{10} \dots c_0}_{11} \underbrace{m_{51} \dots m_0}_{52} \text{ wird als Dezimalzahl}$$

$$q_{(10)} = (-1)^v \cdot 2^{C-1023} \cdot (1.m)_{(2)} = (-1)^v \cdot 2^{C'} \cdot (1.r)_{(10)}$$

interpretiert.

Bemerkung: Im IEEE-Format ist die Bezeichnung „Floating-Point“ etwas irreführend, da in dieser Darstellung das Komma gar nicht explizit auftritt.

IEEE-64-Bit-Format

Darstellung von $(-3.5)_{(10)}$ im IEEE-64-Bit-Format

$$\begin{array}{rcl}
 (-3.5)_{(10)} & = & (11.1)_{(2)} \\
 & = & (b_1 b_0 . b_{-1})_{(2)}
 \end{array}
 \quad \left| \quad
 \begin{array}{rcl}
 3 : 2 & = & 1 \quad 1 \\
 1 : 2 & = & 0 \quad 1 \\
 0.5 : 2 & = & 1.0 \quad 1
 \end{array}$$

$$\rightarrow \ell = \underline{1}$$

$$\rightarrow (-3.5)_{(10)} = (-1)^1 \cdot 2^{\underline{1}} \cdot (1.11)_{(2)}$$

$$\rightarrow v = 1 \quad C' = \ell = 1 \quad (1.m)_{(2)} = (1.11)_{(2)}$$

$$\begin{aligned}
 \rightarrow C &= (C' + 1023)_{(10)} = (1 + 1023)_{(10)} = (1024)_{(10)} = 2^{10} \\
 &= (10000000000)_{(2)} \quad (11 \text{ Bit})
 \end{aligned}$$

$$\rightarrow m = (110 \dots 0)_{(2)} \quad (52 \text{ Bit})$$

$$(-3.5)_{(10)} = (1 \quad 1 \underbrace{000 \dots 0}_{10 \times} \quad 11 \underbrace{000 \dots 0}_{50 \times})_{\text{IEEE-64-Bit}}$$