



Zahldarstellung in Digitalrechnern

Die kleinste Informationseinheit in Digitalrechnern ist **1 Bit (binary digit)** mit den Werten 0 oder 1. Üblicherweise werden 8 Bit zu einem **Byte** zusammengefasst und Information nur in Vielfachen davon gespeichert.

Die Zahldarstellung ist auf verschiedenen Rechnern unterschiedlich. Die folgenden Angaben beziehen sich auf die x86-64-Architektur und richten sich nach der C-Notation:

Integer-Zahlen

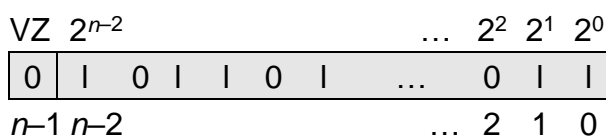
Mit INTEGER-Zahlen können ganze Zahlen des erlaubten Zahlenbereichs **exakt** dargestellt werden.

Formate

- Char (int8_t): 1 Byte = 8 Bits
- Short Integer (int16_t): 2 Bytes = 16 Bits
- Integer (int32_t): 4 Bytes = 32 Bits
- Long Integer (int64_t): 8 Bytes = 64 Bits

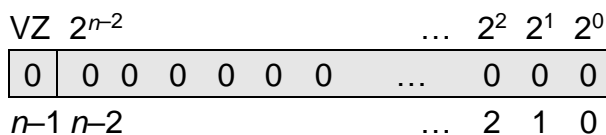
Darstellung

- positive Zahlen



$$0 < x \leq 2^{n-1} - 1$$

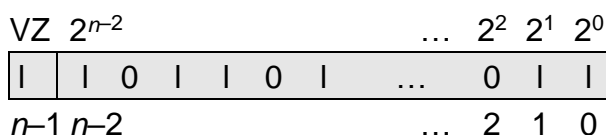
- Null



$$x = 0$$

- negative Zahlen

Darstellung im Allgemeinen im Zweierkomplement, d.h. Invertieren der Darstellung von $(-x)$ ($1 \rightarrow 0, 0 \rightarrow 1$, Einserkomplement) und Addition von 1.



$$-2^{n-1} \leq x < 0$$



Beispiel: Darstellung von $\pm 1234_{10}$ als Short Integer

$$\begin{array}{r}
 1234_{10} = 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 0\ 0\ 1\ 0\ 2 \\
 \text{Einserkomplement } 1\ 1\ 1\ 1\ 1\ 0\ 1\ 1\ 0\ 0\ 1\ 0\ 1\ 1\ 0\ 1\ 2 \\
 + 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 2 \\
 \text{Zweierkomplement } 1\ 1\ 1\ 1\ 1\ 0\ 1\ 1\ 0\ 0\ 1\ 0\ 1\ 1\ 1\ 0\ 2 = -1234_{10}
 \end{array}$$

Gleitpunktzahlen (ANSI/IEEE 754 Standard)

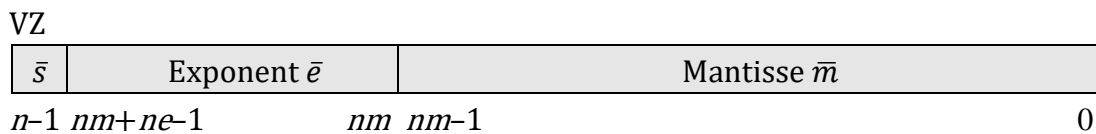
Mit Gleitkommazahlen lassen sich reelle Zahlen im Allgemeinen nur **fehlerbehaftet** darstellen, da sie zu diesem Zweck gerundet werden müssen. Die auf $(nm + 1)$ Stellen gerundete Dualzahl sei wie folgt normiert (Eindeutigkeit der Darstellung):

$$x = \pm \underbrace{1.z_{-1}z_{-2} \dots z_{-nm}}_m \cdot 2^e, \quad z_{-i} \in \{0,1\}.$$

Formate

- einfache Genauigkeit (float): 4 Bytes = 32 Bits
- doppelte Genauigkeit (double): 8 Bytes = 64 Bits

Darstellung



Vorzeichen

Das Vorzeichen wird in einem Vorzeichenbit festgehalten. Ist das Vorzeichenbit gesetzt ($\bar{s} = 1$), so ist die Zahl negativ.

Mantisse

Bei der Speicherung der Mantisse wird die Ziffer vor dem Dualpunkt als gesetzt angenommen und nur die Ziffern hinter dem Dualpunkt gespeichert (1 Bit zusätzliche Genauigkeit):

$$\bar{m} = m - 1.$$

Exponent

Der Exponent wird um einen Offset (Bias = $2^{ne-1} - 1$) verschoben, so dass nur positive Werte auftreten:

$$\begin{aligned}
 \bar{e} &= e + (2^{ne-1} - 1), \\
 e &= \bar{e} - (2^{ne-1} - 1), \quad -2^{ne-1} + 1 \leq e \leq 2^{ne-1}.
 \end{aligned}$$

