

Im folgenden soll ein Überblick über die in Computersystemen bzw. Programmiersprachen verwendeten Zahlen inklusive ausgewählter Algorithmen (in neutraler Notation) gegeben werden. Des öfteren wird auch Bezug auf die verbreitete Programmiersprache Java (Version 7) genommen .

Dastellung von Zahlen:

Computer (bzw. Compiler) stellen Zahlen durch Bitfolgen im Dualsystem (Basis $b = 2$) dar. Bitfolgen sind Folgen, die ausschließlich aus den Ziffern 0 und 1 bestehen, z.B. 10011010110001 . Zum Vergleich: Das Dezimalsystem verwendet die Ziffern 0 ... 9 .

Der Wert einer Zahl im Dualsystem wird dabei genauso ermittelt wie im Dezimalsystem.

Beispiel $b=10$: $276_{10} = 2 \cdot 10^2 + 7 \cdot 10^1 + 6 \cdot 10^0 = 200 + 70 + 6$

Beispiel $b = 2$: $110101_2 = 1 \cdot 2^5 + 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 = 32 + 16 + 0 + 4 + 0 + 1 = 53_{10}$

Anmerkungen: Die tiefgestellte Zahl verdeutlicht jeweils die Basis .
Das ganz links stehende Bit ist dasjenige mit dem höchsten Wert .

Ganzzahlen (engl.: integers):

Mit einer Bitfolgenlänge von n kann man nach den Gesetzen der Kombinatorik 2^n Zahlen darstellen.
Bei natürlichen Zahlen $0\ 1\ \dots\ 2^n - 1$ und bei **ganzen Zahlen** $-2^{n-1}\ -2^{n-1}+1\ \dots\ -1\ 0\ 1\ 2\ \dots\ 2^{n-1}-1$.

Beispiel $n=4$ Bit : $0\ 1\ 2\ \dots\ 14\ 15$ (natürliche Z.)
 $-8\ -7\ -6\ -5\ -4\ -3\ -2\ -1\ 0\ 1\ 2\ 3\ 4\ 5\ 6\ 7$ (ganze Z.)
Beispiel $n=64$ Bit : $0\ 1\ \dots\ 18446744073709551615$ (natürliche Z.)
 $-9223372036854775808\ \dots\ 9223372036854775807$ (ganze Z.)

Für die ganzen Zahlen verwendet man folgende Dualdarstellung (hier für $n=4$ dargestellt):

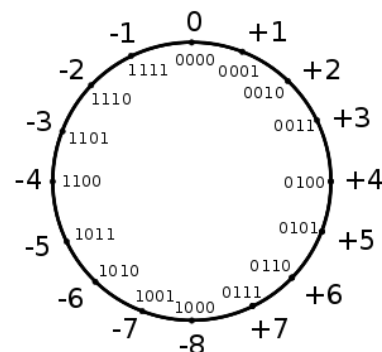
1000	1001	1010	1011	1100	1101	1110	1111	0000	0001	0010	0011	0100	0101	0110	0111
-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4	5	6	7

Man erkennt, dass jede negative Zahl im höchstwertigen Bit eine 1 hat, während bei positiven Zahlen dort eine 0 steht. Das erste Bit stellt also das Vorzeichen dar (1 negativ ; 0 positiv).

Diese Darstellung nennt man auch **Zweierkomplementdarstellung**. Die Bezeichnung rührt daher, dass jede Addition einer Zahl zu ihrem bitweisen Komplement (aus 0 wird 1 und aus 1 wird 0) stets als Ergebnis die Zahl $111\dots 1$ hat. Dies entspricht dezimal der Zahl -1 . Beispiel: $1011 + 0100 = 1111$.
Daher erzeugt man aus einer gegebenen positiven Zahl die entsprechende negative durch Komplementbildung und anschließender Addition der Zahl 1. Man beachte: Jede positive Zahl beginnt mit einer 0 !!
Beispiel: $12_{10} = 01100_2$ Komplement+1: $10011 + 00001 = 10100 = 20_{10} - 32_{10} = -12_{10}$.

Die Zweierkomplementdarstellung kann man sich mithilfe eines Kreises veranschaulichen (siehe Abbildung).

Man sieht z.B., dass einander gegenüberliegende Zahlen sich lediglich in dem höchstwertigen Bit unterscheiden.



Umrechnungsmethoden:

① Eine dezimale Ganzzahl a lässt sich folgendermaßen in ihre Dualdarstellung $b_{n-1}b_{n-2}\ \dots\ b_1b_0$ umwandeln :

```
i = 0
wiederhole
  b[i] = a mod 2
  a = a div 2
  i = i+1
bis a = 0
```

Die Stellen der Dualdarstellung (zunächst ohne Zweierkomplement) werden also von rechts nach links ermittelt !

Beispiel: $a = 37$. Es wird fortwährend durch 2 dividiert und der Rest notiert:

```
37 / 2 = 18 Rest 1.
18 / 2 = 9 Rest 0.
9 / 2 = 4 Rest 1.
4 / 2 = 2 Rest 0.
2 / 2 = 1 Rest 0.
1 / 2 = 0 Rest 1.
```

Beim Ergebnis 0 endet der Algorithmus. Die Reste (oberster Wert ganz rechts notiert) sind 100101 .

Dies ist die Dualdarstellung von **37**.

Wählt man eine Bitlänge von $n = 8$, so ist die Dualdarstellung von 37 gleich 00100101 .

Bei negativen Zahlen muss noch das Zweierkomplement gebildet werden. Beispiel -37 :

Die Dualdarstellung von $+37$ ist 00100101 (siehe oben).

Das Zweierkomplement wird nun gebildet, indem alle Bits umgedreht werden und anschließend 1 addiert wird:
 00100101 → 11011010 11011010 + 1 = **11011011** (8bit-Zweierkomplement von -37)

⊗ Für die Umrechnung einer im Zweierkomplement vorliegenden n-bit-Dualzahl (Bitfolge $b_{n-1}b_{n-2} \dots b_1b_0$) in das Dezimalsystem kann man folgende Formel verwenden:

$$b_{n-1}b_{n-2} \dots b_1b_0 = \left(\sum_{i=0}^{n-2} b_i \cdot 2^i \right) - b_{n-1} \cdot 2^{n-1} \quad \text{Zweierkomplementzahl}$$

Ein Beispiel für eine **8bit**-Ganzzahl (n=8) :

1100 | 0101

$11000101 = 1 \cdot 2^0 + 0 \cdot 2^1 + 1 \cdot 2^2 + 0 \cdot 2^3 + 0 \cdot 2^4 + 0 \cdot 2^5 + 1 \cdot 2^6 - 1 \cdot 2^7 = 1 + 0 + 4 + 0 + 0 + 0 + 64 - 128 = 69 - 128 = -59$
 dezimal

Ein Beispiel für eine **32bit**-Ganzzahl (n=32) :

1000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000

Die Zahl ist -2147483648 dezimal

Ein Beispiel für eine **64bit**-Ganzzahl (n=64) :

0101 | 0001 | 1010 | 1111 | 1111 | 1111 | 1111 | 1111 | 1111 | 1111 | 1111 | 1111 | 1111 | 1111 | 1111 | 1111

Die Zahl ist 5886204712973238271 dezimal

Die Programmiersprache **Java** kennt verschiedene Typen von Ganzzahlen (Tabelle):

Datentyp	Format	Bereich
byte	8 Bit	$-2^7 \dots 2^7-1 = -128 \dots 127$
short	16 Bit	$-2^{15} \dots 2^{15}-1 = -32768 \dots 32767$
int	32 Bit	$-2^{31} \dots 2^{31}-1 = -2147483648 \dots 2147483647$
long	64 Bit	$-2^{63} \dots 2^{63}-1 = -9223372036854775808 \dots 9223372036854775807$

Zu beachten ist, dass bei Ganzzahloperationen ein **Überlauf** (englisch: overflow) auftreten kann, der bei Programmiersprachen (z.B. Java) nicht zu Fehlermeldungen führt !

Zum Beispiel liefert die Rechnung $2147483647+3$ beim Datentyp int (Java7) das Ergebnis -2147483646 und nicht wie erwartet $+2147483650$.

Die Erklärung dafür ist in der Zweierkomplementdarstellung zu suchen.

Demgegenüber sind Ganzzahloperationen empfindlich, wenn z.B. irgendeine Ganzzahl durch 0 geteilt wird. Java7 „wirft“ dann eine sog. „**ArithmeticException**“ .

Später wird zu sehen sein, dass Fließkommaoperationen diesbezüglich etwas anders reagieren.

Dann ist die Zahl $z_{dez} = (-1)^0 \cdot (4503599627370495 + 2^{52}) \cdot 2^{1048-1075} = 9007199254740991 \cdot 2^{-27} = 6,7108863999999992 \cdot 10^7$.

Weitere Beispiele:

Die Zahl 0 wird dargestellt (E=0 ; denormalisiert ; V=0 oder V=1 für „-0“ möglich)

0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Die größte darstellbare 64bit-Gleitkommazahl (E=2046) ist

0111	1111	1110	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Wert = $(11 \dots 1) \cdot 2^{2046-1075} = (2^{52}-1+2^{52}) \cdot 2^{971} = (2^{53}-1) \cdot 2^{971} = 1,79769313486231570 \cdot 10^{308}$

Die kleinste darstellbare (denormalisierte) 64bit-Gleitkommazahl (E=0) ist

0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0001
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Wert = $(00 \dots 01) \cdot 2^{-1074} = 1 \cdot 2^{-1074} = 4,94065645841246544 \cdot 10^{-324}$

Die Darstellung von $-\infty$ bzw. -Infinity (E=2047 ; M=0) ist

1111	1111	1111	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Die Darstellung von $+\infty$ bzw. +Infinity ist ist entsprechend, mit V=0 .

Die Darstellung von NaN (E=2047 ; M≠0) ist z.B.

1111	1111	1111	1000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Es handelt sich hier um ein sog. **stilles NaN** , bei dem das höchstwertige Bit von M den Wert 1 hat . Bei stillen NaNs ist das Ergebnis nicht (mathematisch) definiert , z.B $0.0/0.0$, $\log(-1.0)$, $\sqrt{(-1.0)}$ etc .

Es gibt auch NaNs, bei denen das höchstwertige Bit von M den Wert 0 hat (**anzeigende NaNs**). Mindestens 1 anderes Bit der Mantisse M muss aber den Wert 1 besitzen. Anzeigende NaNs unterstützen fortgeschrittene Eigenschaften wie gemischte numerische und symbolische Berechnungen oder andere grundlegende Gleitkommaarithmetik .

Interessant ist, dass Programmiersprachen (z.B. Java7) bei mathematisch nicht definierten Fließkommaoperationen anders reagieren als bei entsprechenden Ganzzahloperationen. Bei Fließkommaoperationen wird in diesen Fällen keine „**ArithmeticException**“ geworfen, sondern es erscheint als Ergebnis z.B. Infinity oder NaN (vgl. Beispiele oben) .

Es ist noch zu klären, wie man zu einer gegebenen Dezimalzahl x die zugehörige Codierung VEM findet:

Beispiel: $x = 37,1$

Als erstes wird eine Dualdarstellung (Vorkomma und Nachkomma) von 37,1 erzeugt . Der Vorkommateil 37 wird verarbeitet wie bereits bei Ganzzahlen gezeigt. Es gilt $37_{10} = 100101_2$. Beim Nachkommateil $nk = 1$ wird nach folgendem Algorithmus gerechnet:

```

mult = 10^m // m ist die Länge von nk
wiederhole
    nk = nk * 2
    Schreib nk div mult
    nk = nk mod mult
bis nk = 0 oder „Periode erkannt“

```

nk = 1	mult = 10		
$1 \cdot 2 \text{ div } 10 = 0$		$nk = 1 \cdot 2 \text{ mod } 10 = 2$	
$2 \cdot 2 \text{ div } 10 = 0$		$nk = 2 \cdot 2 \text{ mod } 10 = 4$	
$4 \cdot 2 \text{ div } 10 = 0$		$nk = 4 \cdot 2 \text{ mod } 10 = 8$	
$8 \cdot 2 \text{ div } 10 = 1$		$nk = 8 \cdot 2 \text{ mod } 10 = 6$	
$6 \cdot 2 \text{ div } 10 = 1$		$nk = 2 \cdot 2 \text{ mod } 10 = 2$	
$2 \cdot 2 \text{ div } 10 = 0$		$nk = 2 \cdot 2 \text{ mod } 10 = 4$	

Man erkennt eine Periode (mit Vorperiode) ! Der duale Nachkommateil ist $000\overline{11}$. Die komplette Dualdarstellung für $37,1_{10}$ ist also $100101,000\overline{11}_2$.

Diese Dualdarstellung muss jetzt normiert werden, d.h. gesucht ist eine Form $1, M \cdot 2^E$.

Das Komma muss um 5 Stellen nach links verschoben werden. Das entspricht einer Multiplikation mit 2^5 .

Damit gilt: $100101,00011 = 1,0010100011 \cdot 2^5$. Also ist der Exponent $E = 5_{10}$.

Es muss aber noch der Biaswert von 1023 addiert werden, daher $E = 1028_{10} = 10000000100_2$.

Lässt man bei der Mantisse die führende 1 weg und ergänzt die Periode bis auf 52 Stellen, so ergibt sich als Mantisse $M = 0010100011001100110011001100110011001100110011001100110011001101$.

Anmerkung: Das letzte Bit entsteht durch **Aufrundung**, denn $1100(1) = 1101(0) = 1101$.

Für das Vorzeichen wird $V = 0$ gespeichert, denn 37,1 ist positiv.

Die dezimale 37,1 wird also im IEEE-Gleitkommaformat gespeichert als VEM =

0100	0000	0100	0010	1000	1100	1100	1100	1100	1100	1100	1100	1100	1100	1100	1101
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Weitere Codierungsbeispiele:

Die Zahl $\pi \approx 3.1415926535897932384$ wird so codiert :

0100	0000	0000	1001	0010	0001	1111	1011	0101	0100	0100	0100	0010	1101	0001	1000
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Wert = $(2570638124657944 + 2^{52}) \cdot 2^{1024-1075} = 7074237752028440 \cdot 2^{-51} = 3,1415926535897931\dots$

Nur die ersten 15 Nachkommastellen sind richtig. Das entspricht der Präzision von 64bit-Gleitkommazahlen.

Die Zahl 1.0 wird so codiert :

0011	1111	1111	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Die Zahl 2.0 wird so codiert :

0100	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Abstand benachbarter Zahlen :

Ändert man das letzte Bit auf 1, so hat man die nächste hinter 2,0 darstellbare Gleitkommazahl erzeugt. Diese hat den Wert $2,00000000000000004440892098500626\dots$, unterscheidet sich also von 2,0 um den Wert $4,440892098500626E-16$. In Java7 entspricht dies übrigens dem Aufruf von `Math.ulp(2.0)`.

Der **Abstand benachbarter Zahlen** ist bei großen Zahlen größer als bei kleinen Zahlen.

Beispiele (Java7):

`Math.ulp(1.0) = 2.220446049250313E-16`

`Math.ulp(10.0) = 1.7763568394002505E-15`

`Math.ulp(100.0) = 1.4210854715202004E-14`

Programmierung der Ausgabe auf dem Bildschirm (Java7-Methoden):

```
public static String doubleToBinaryIEEE754(double x) {
    String sBits = Long.toBinaryString(Double.doubleToLongBits(x));
    // da führende Nullen weggelassen werden ...
    while (sBits.length() < 64)
        sBits = "0" + sBits;
    return sBits;
}
public static double binaryIEEE754ToDouble(String sBin) {
    double x = Double.LongBitsToDouble(Long.valueOf(sBin, 2));
    return x;
}
```


Maschinengenauigkeit:

Die Maschinengenauigkeit eines Computers ist die **kleinste positive Zahl ϵ** , für die $1 + \epsilon - 1 > 0$ gilt.
Für alle Zahlen x zwischen 0 und ϵ gilt $1 + x - 1 = 0$.

Für den **TI84** ermittelt man als Maschinengenauigkeit $\epsilon = 2^{-39} \approx 1,8189894035458 \cdot 10^{-12}$.

Für **Calc** (Libre Office) ermittelt man als Maschinengenauigkeit $\epsilon = 2^{-47} \approx 7,1054273576010018 \cdot 10^{-15}$.

Für **Java7** (double) ermittelt man als Maschinengenauigkeit $\epsilon = 2^{-52} \approx 2,22044604925031308 \cdot 10^{-16}$.

Für **Delphi7** (extended) ermittelt man als Maschinengenauigkeit $\epsilon = 2^{-63} \approx 1,0842021724855 \cdot 10^{-19}$.

Anmerkung: extended arbeitet mit 80bit – Speicher !

Rechnet man mit Zahlen in der Nähe der Maschinengenauigkeit, so können total falsche Ergebnisse entstehen.

Beispiel: Berechne $(1 + 2^{-39} - 1) \cdot 2^{39}$. Exaktes Ergebnis müsste 1 sein.

Der TI84 berechnet aber $0,989560465$.