

# 1.3 Zahlen-Darstellungen-Arithmetik

Stellenwertsysteme/Zahlensysteme/Konvertierung

Komplement-Zahldarstellungen

ganze pos. & neg. Zahlen (INTEGER, ...)

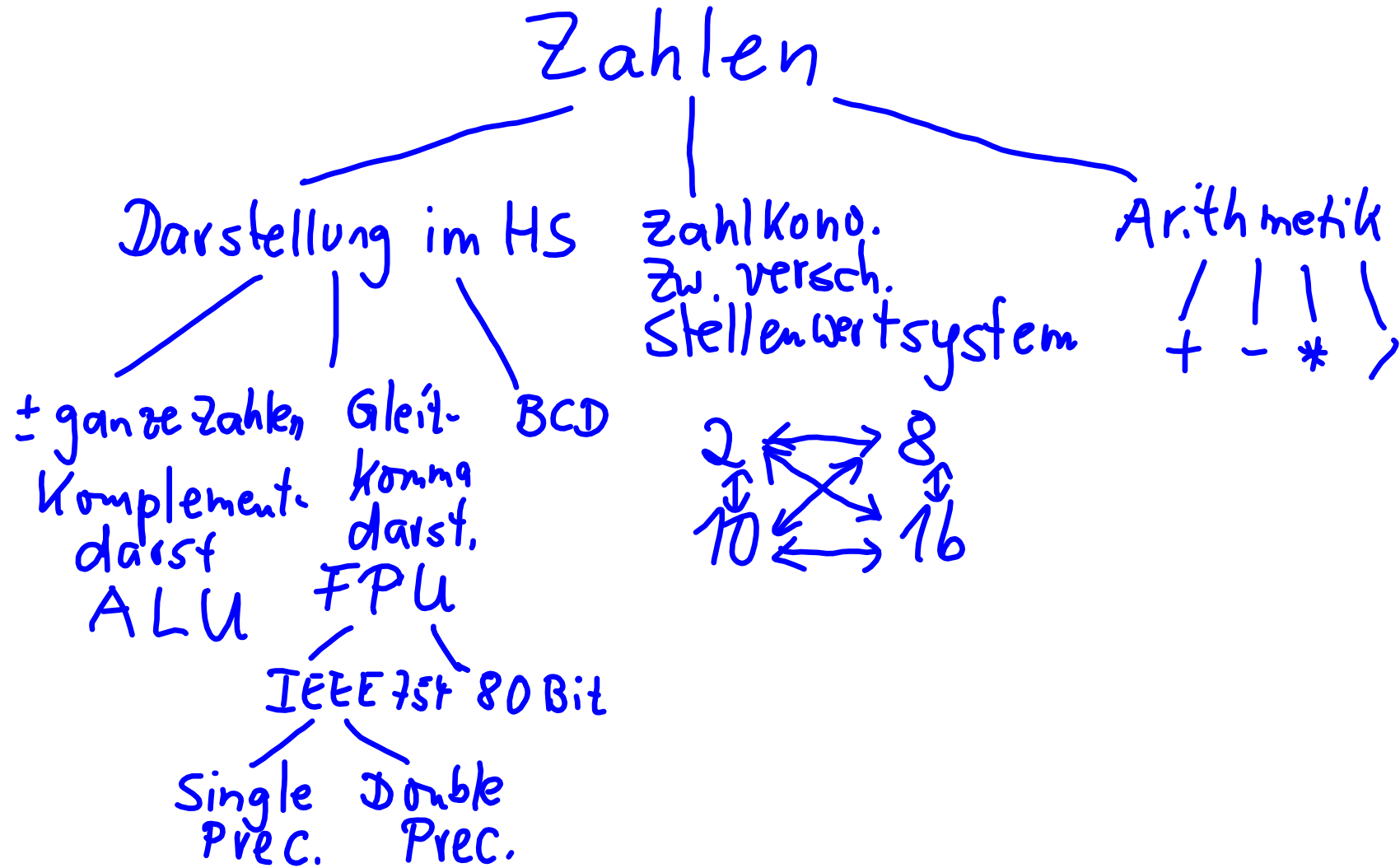
Gleitkomma-Zahldarstellung (IEEE 754)

$$Z = v \cdot m \cdot B^e \quad m = \text{Mantisse}, e = \text{Exponent}, v = +/-$$

Alternative Zahlendarstellungen

BCD-Zahldarstellung (Binary Coded Decimals)

Arithmetik (Add./Sub./Mult. im 2/8/16er-Zahlensystem)



# Stellenwertsysteme/Zahlensysteme

Die allgemeine Zahldarstellung einer Zahl  $z$  in einem Stellenwertsystem zur Basis  $B$  mit den Ziffern  $z_i \in \{0,1,2,3,\dots,B-1\}$  lautet:

$$z \Big|_B = \overbrace{B \cdot (z)}_8 = (z_n z_{n-1} \dots z_0 z_{-1} \dots z_{-m})_B = \sum_{i=-m}^n z_i B^i = z_n B^n + z_{n-1} B^{n-1} + \dots$$

für ( $B \geq 2$ ) gebr. ganz \* B

$$= \underbrace{((z_n B + z_{n-1})B + z_{n-2})B + \dots + z_1 B + z_0}_{\text{ganzzahliger Anteil}} + \underbrace{((z_{-m} B^{-1} + z_{-m+1})B^{-1} + z_{-m+2})B^{-1} + \dots + z_{-1} B^{-1}}_{\text{gebrochener Anteil}} \cdot B \quad (*)$$

$B=8$   $\frac{4}{8}$

spezielle Zahlensysteme:

$\begin{matrix} +1 \\ \downarrow \\ VI \end{matrix}$   $\begin{matrix} -1 \\ \downarrow \\ IV \end{matrix}$

Dualsystem	$B=2$	$z_i \in \{0, 1\}$	←
Oktalsystem	$B=8$	$z_i \in \{0, \dots, 7\}$	←
Dezimalsystem	$B=10$	$z_i \in \{0, \dots, 9\}$	←
Hexadezimalsystem	$B=16$	$z_i \in \{0, \dots, 9, A, \dots, F\}$	←

Hornerschema zur Berechnung des Wertes von  $z$ :

$1989$  ,  $0,123 = \frac{123}{1000}$

	1	9	8	9
10	↓	10	190	1980
	1	19	198	<u>1989</u>

	3	2	1
$\frac{1}{10}$	↓	$\frac{3}{10}$	$\frac{23}{100}$
	3	$\frac{23}{10}$	$\frac{123}{100}$

# Zahlenkonvertierung

Konvertierung:  $\begin{matrix} \rightarrow & \text{ganzer Anteil} \\ \rightarrow & \text{gebrochener Anteil} \end{matrix}$

Umwandlung der Zahldarstellung in Stellenwertsysteme verschiedener Basen

a)  $10 \begin{matrix} \rightarrow & 2 \\ \rightarrow & 8 \\ \rightarrow & 16 \end{matrix}$  (ganzer Anteil, Divisionsmethode)

Beim Dividieren einer ganzen Zahl z durch die Zielbasis B ergeben die Divisionsreste nacheinander gelesen die von rechts gelesene ganze Zahl im Zielstellenwertesystem zur Basis B.

$$z = (1996)_{10} = (3714)_8$$

gebr. Teil	Rest
0,125	1
0,25	2 $\rightarrow \frac{2}{8} = \frac{1}{4}$
0,375	3
0,5	4
0,625	5
0,75	6
0,875	7
0	0

$$1996 : 8 = 249,5 = 249 \frac{4}{8} = 249 \cdot 8 + 4$$

$$249 : 8 = 31,125 = 31 \cdot 8 + 1$$

$$31 : 8 = 3 \text{ Rest } 7$$

$$3 : 8 = 0 \text{ Rest } 3$$

Rest  
↓



c) **gebr. Dezimalzahlen**  $\begin{matrix} \nearrow 2 \\ \rightarrow 8 \\ \searrow 16 \end{matrix}$  **Multiplikationsmethode**

Beim Multiplizieren einer echt gebrochenen Dezimalzahl mit der Zielbasis B ergeben die Überläufe über 1 nacheinander gelesen die Ziffernfolge der gebrochenen Zahl im Stellenwertsystem zur Zielbasis B.


Bsp:  $z = (0,6)_{10} = \overline{(0,4631)}_8 = 0,463146314631\dots$

$$\begin{array}{l} \rightarrow 0,6 \cdot 8 = 4,8 \\ 0,8 \cdot 8 = 6,4 \\ 0,4 \cdot 8 = 3,2 \\ 0,2 \cdot 8 = 1,6 \\ 0,6 \cdot 8 = 4,8 \end{array}$$

$$\uparrow \text{Rest} = ( \quad ) B^{-1} \cdot B$$

$$\begin{array}{l} z_{-1} = 6 \\ z_{-2} = 3 \\ z_{-3} = 3 \end{array}$$

Bsp:  $(0,\overline{3})_{10} = (0,1)_3$

d)  Ziffern  $\leftrightarrow$  Aufsplittung  
Zusammenfassung

Dualziffern  $\rightarrow$  Triple = Oktalziffer  
Dualziffern  $\rightarrow$  Tetrade = Hexadezimalziffer

Bsp:  $ABCH = (101010111100)_2$   
 $(5274)_8$

Bsp:  $(0,4631)_8 = (0,100110011001)_2$   
 $= (0,999)_{16} = (0,\bar{9})_{16}$

# Komplement-Zahldarstellungen

**B-Komplementdarstellung (echtes Komplement):**

$$z = \begin{cases} (z)_B & z \geq 0 \text{ und } z < B^{n-1} \\ \text{falls} & \\ (B^n - |z|)_B & z < 0 \text{ und } |z| \leq B^{n-1} \end{cases}$$

**(B-1)-Komplementdarstellung (unechtes Komplement):**

$$z = \begin{cases} (z)_B & z \geq 0 \text{ und } z < B^{n-1} \\ \text{falls} & \\ (B^n - |z| - 1)_B & z < 0 \text{ und } |z| < B^{n-1} \end{cases}$$

Ganze Zahlen in **Zweier-Komplementdarstellung** zur Basis 2 mit n-Bitstellen

$$z = \begin{cases} (z)_2 & z \geq 0 \text{ und } z < 2^{n-1} \\ \text{falls} & \\ (2^n - |z|)_2 & z < 0 \text{ und } |z| \leq 2^{n-1}, \text{ d.h. } z \geq -2^{n-1} \end{cases}$$

Ganze Zahlen im **Einerkomplement** (n-Bitdarstellung)

$$z = \begin{cases} (z)_2 & z \geq 0 \text{ und } z < 2^{n-1} \\ \text{falls} & \\ (2^n - |z| - 1)_2 & z < 0 \text{ und } |z| < 2^{n-1}, \text{ d.h. } z > -2^{n-1} \end{cases}$$

# Komplement-Zahldarstellungen

für ganze Zahlen

B-Komplement-Darstellung ganzer Zahlen (echtes Komplement)

$$z'' = \begin{cases} z|_B & \text{falls } z \geq 0, z < B^{n-1} \\ (B^n - |z|)_B & \text{falls } z < 0, |z| \leq B^{n-1} \end{cases}$$

(B-1)-Komplement-Darstellung ganzer Zahlen (unechtes Komplement)

$$z' = \begin{cases} z|_B & \text{falls } z \geq 0, z < B^{n-1} \\ (B^n - |z| - 1)_B & \text{falls } z < 0, |z| < B^{n-1} \end{cases}$$

speziell für B=2:

Zweier-Komplement

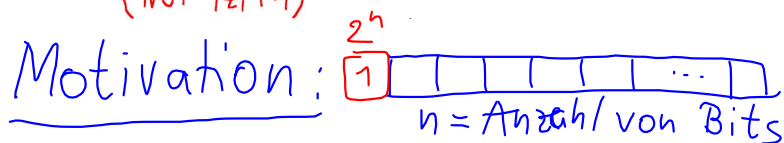
$$z'' = \begin{cases} z|_2 & \text{falls } z \geq 0, z < 2^{n-1} \\ (2^n - |z|)_2 & \text{falls } z < 0, |z| \leq 2^{n-1} \end{cases}$$

(NOT |z| + 1)

Einer-Komplement

$$z' = \begin{cases} z|_2 & \text{falls } z \geq 0, z < 2^{n-1} \\ (2^n - |z| - 1)_2 & \text{falls } z < 0, |z| < 2^{n-1} \end{cases}$$

(NOT |z|)

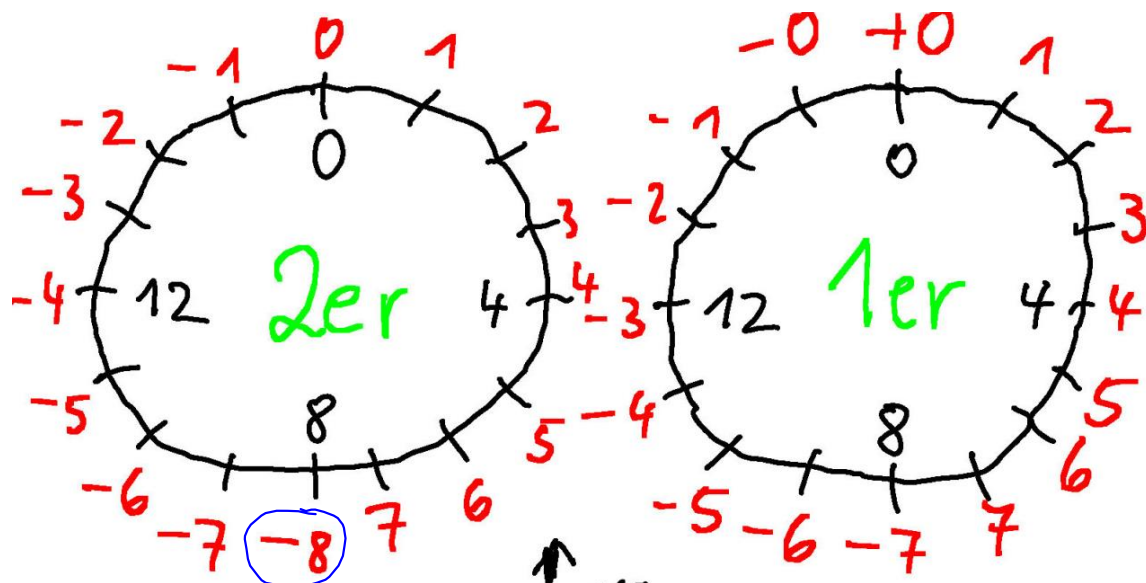


z.B. nur 8 Bit zur Speicherung einer ganzen Zahl  
→ 256 ganze Z.

Subtraktion:  $d = a - b + \underbrace{2^n - 2^n}_0$

$$= a + (2^n - b) = a + (-b)$$

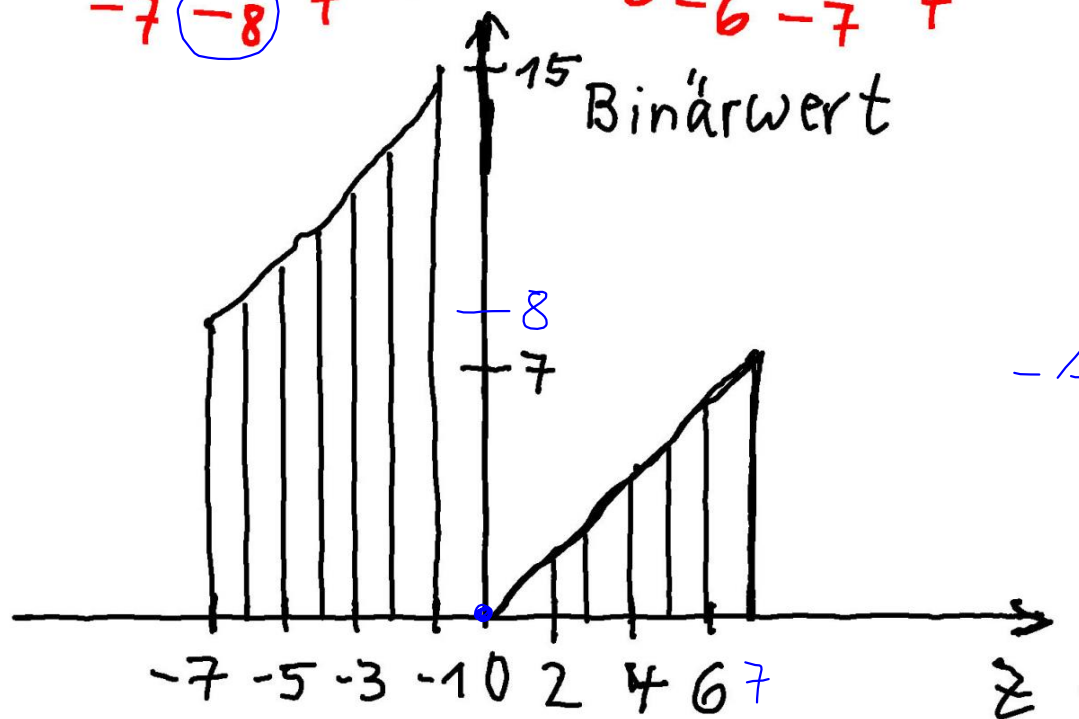
Komplement



Tetraden  
 $n=4$   
 $\downarrow$   
 16 Zahlen

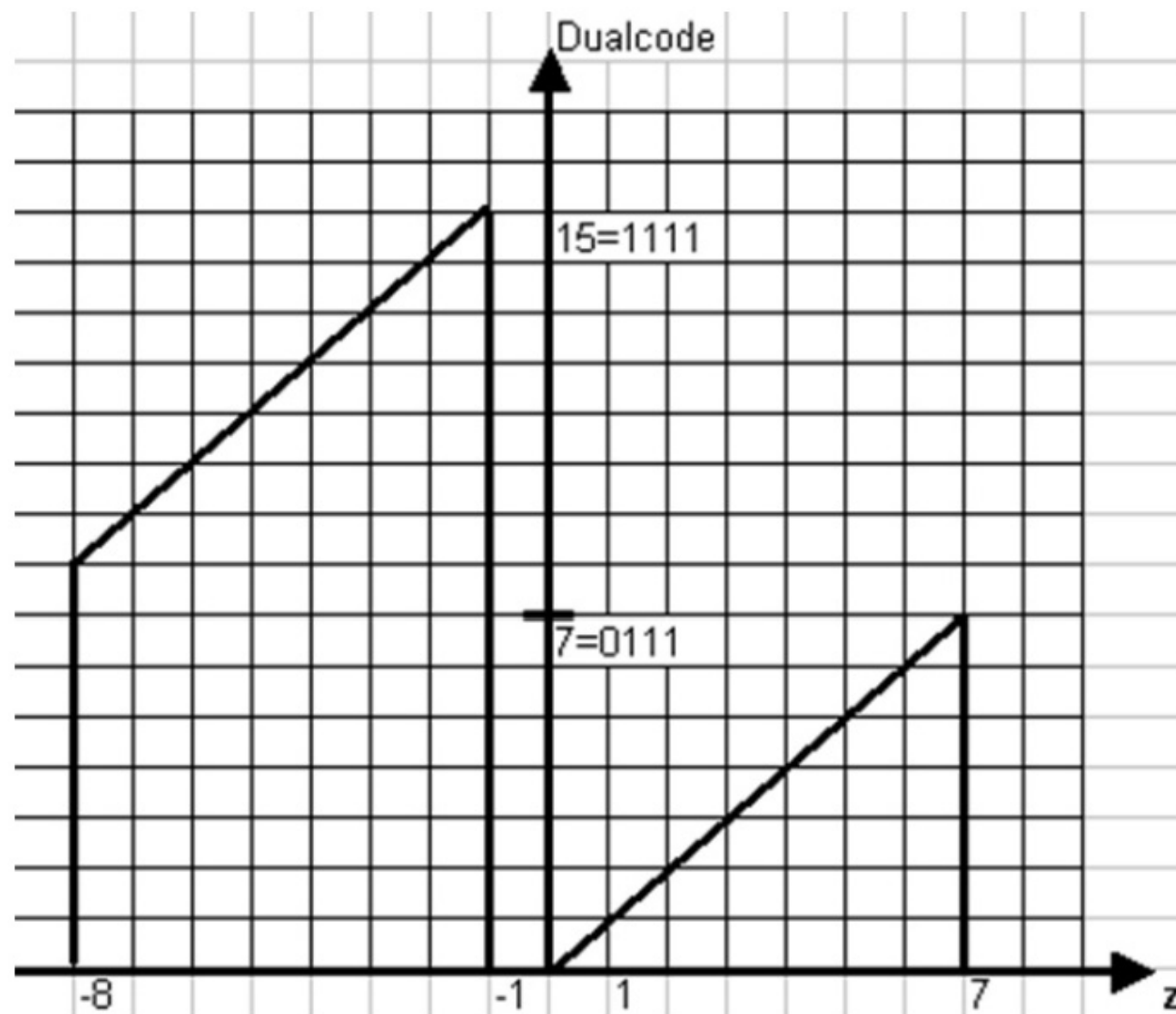
$-7 \dots -1 \quad 0 \dots 7$   
 $\quad \quad \quad -0$

$-4 \rightarrow 8$



**n=4 Bits**

0	0000	0
1	0001	1
2	0010	2
3	0011	3
4	0100	4
5	0101	5
6	0110	6
7	0111	7
8	1000	-8
9	1001	-7
A	1010	-6
B	1011	-5
C	1100	-4
D	1101	-3
E	1110	-2
F	1111	-1



Bsp:  $n=3 \rightarrow z \in \{-4..0..3\}$

$z = -1$   
 $z_{2er-kompl.} = 2^3 - |-1| = 7$

$z = -4$   
 $z_{2er\ kompl} = 2^3 - |-4| = 4$

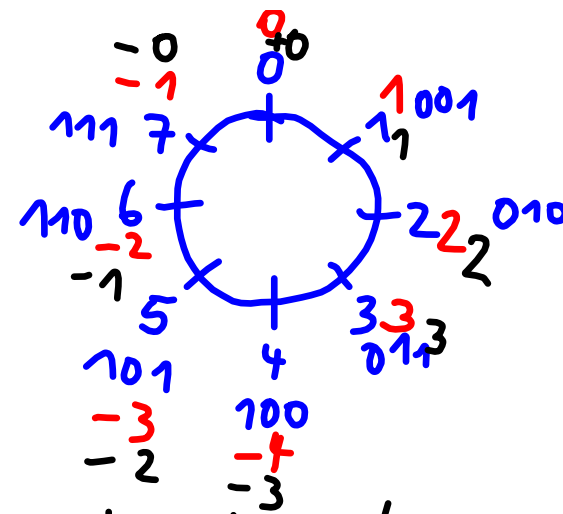
$= NOT |z| + 1$   
 $= |-100|$

$011 + 001 = 100 = 4$

Zahlenbereichsüberschreitung / Overflow / Überläufe

Add.  $1+2 = \begin{array}{r} 001 \\ 010 \\ \hline 011 \end{array} \rightarrow = 3$  aber  $1+4 = \begin{array}{r} 001 \\ 100 \\ \hline 101 = 5 \end{array} \rightarrow = -3$

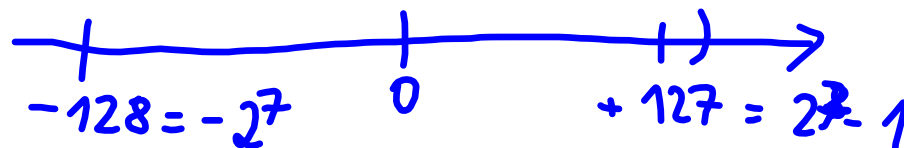
$-2-3 = (-2)+(-3) = \begin{array}{r} 110 \\ +101 \\ \hline 011 = 3 \end{array} \rightarrow = +3$



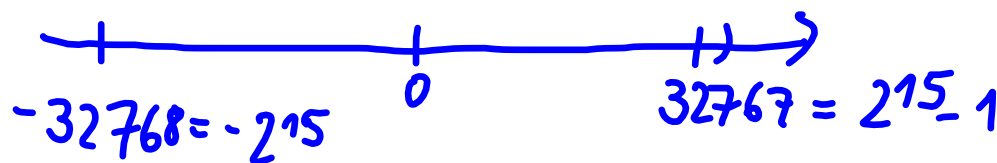
Einerkomplement

$z = -3 \rightarrow z_{1er-k.} = NOT |-z|$   
 $= NOT |-011|$   
 $= 100 = 4_{10}$

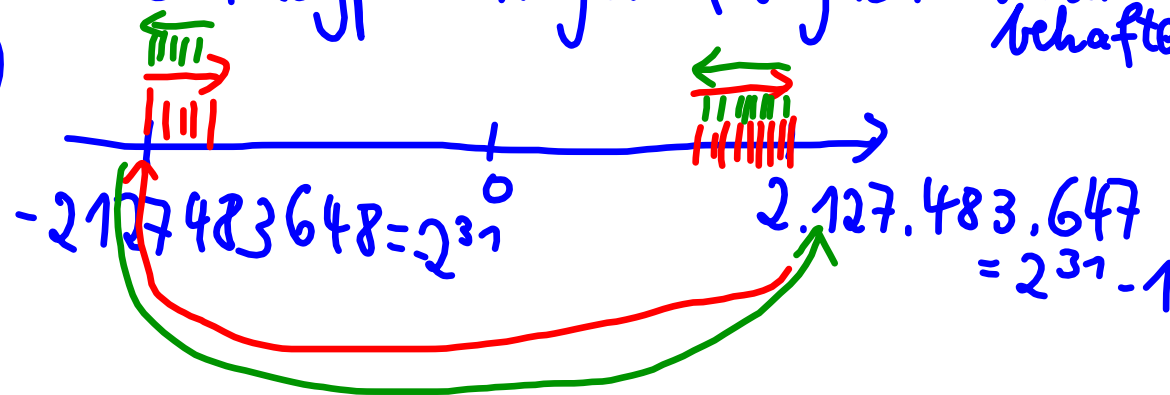
Bsp:  $n = 8$   
 (Zweier-  
 Komplement)



$n = 16$  Datentyp shortint ( $65536 = 2^{16}$ )



$n = 32$   
 (4 Bytes) Datentyp integer (signed = zeichen-  
behaftet)



# Gleitkomma-Zahldarstellung

(nach IEEE 754 / 854)

Gleitkommazahl  
(Floating point)  $z = \pm m \cdot B^e$   $m = \text{Mantisse (Genauigkeit)}$   
 $e = \text{Exponent (Größenordnung)}$

Festkommazahl  $z = z_n \dots z_0, z_{-1} \dots z_{-m}$

normierte Gleitkommazahl: mit  $1 < |m| < B$  (außer Null)

→  $B=2$   $m = 1, \dots$  (außer bei Null)

→  $B=10$   $m = \overset{1}{\underset{9}{\vdots}}, \dots$

allg. Zahldarstellung im HS (Codierung):  $VCM \hat{=} \text{Bitanzahl}$



# Gleitkomma-Zahldarstellung *in FPU*

(nach IEEE 754 / 854)

- a) kurze Gleitkommazahlen SHORTREAL / Zahlen mit einfacher Genauigkeit  
single precision

Gesamtlänge 4 Bytes = 32 Bits mit VCM = 1/8/23 Bits  
mit  $C \in \{0 \dots 127 \ 128 \dots 255\}$   
für  $e \in \{-127 \dots +128\}$ , also  $K = 127$

mit norm. Mantisse:  $m_0, m_1, m_2, \dots, m_{23}$  ( $m_0 = 1$ )

Rundungsfehler  $< 2^{-23} \sim \frac{1}{2} \cdot 2^{-23} = 2^{-24} \approx 2,4 \cdot 10^{-8}$  ← etwa in 8. Dez. Stelle

Datentyp in C: float; in FORTRAN REAL\*4

- b) lange Gleitkommazahlen LONGREAL / Zahlen mit doppelter Genauigkeit  
double precision

mit VCM = 1/11/52 also 8 Bytes = 64 Bits  
 $C \in \{-1023 \dots +1024\}$   
char  $\in \{0 \dots 2047\} \Rightarrow K = 1023$ ,

Rundungsfehler  $< 2^{-52} \approx 10^{-17}$  ← Fehler etwa in der 17. Dezimalstelle hinter dem Komma

Datentyp: in C double  
in FORTRAN REAL\*8

c) INTEL FPU (kein IEEE 754) mit 80 Bit-Register (10 Bytes)

$$VCM = 1/15/64 \Rightarrow K = 2^{15-1} - 1 = 16383$$

$$\text{Rundungsfehler} < 2^{-64} \sim 10^{-21}$$

Datentyp in C: double extended

Bsp:  $-38,625 = -100110, \overset{0,5}{1} \overset{0,25}{0} \overset{0,125}{1} = -1, \underbrace{00110101}_{\text{norm. Mantisse}} \cdot 2^5$

$$\text{Charakteristik} = \overset{k}{127} + \overset{e}{5} = 132 = 128 + 4 = 10000100$$



C 2 | 1 A | 8 0 | 0 0 = 0xC21A8000

# Gleitkommazahlen

## Ausnahmen (exceptions)

(wird angezeigt durch Signalisierungsbit im Stausregister der FPU (floating point unit)  
→ bewirken Behandlungsroutine

- 1) ungenaues Ergebnis → Rundung
- 2) Division durch Null d.h.  $x/0$  mit  $x \neq 0$   
→ Ergebnis =  $+\infty$
- 3) Unterlauf (underflow) → tiny numbers  
(kleiner als kleinste darstellbare Zahl)
- 4) Überlauf (overflow)  
Überschreiten des Zahlbereiches,  
Ergebnis =  $+\infty$  bzw. Sonderzahl
- 5) unzulässige Operation/invalid operation  
z.B.  $+\infty - \infty$  oder  $0 * \infty$  oder  $0/0$  oder  $\infty/\infty$   
oder SQRT(-x)

## Sonderzahlen

- 1) vorzeichenbehaftete Null (signed null)  $+0 = -0$
- 2) vorzeichenbehaftetes unendlich (signed infinity)
- 3) Nichtzahlen (not-a-numbers = NaN)
  - a) SNaN signaling NaNs  
z.B. Anzeige nicht initialisierter Variablen  
(anzeigende Nichtzahlen)
  - b) QNaN quiet NaNs bei ungültigen bzw.  
nichtverfügbaren Operanden  
(nichtanzeigende Nichtzahlen)

Null:  $z = 0,0 * 2^0$



unendlich:  $z = \pm 0,0 * 2^{255} = \pm 0,0 * 2^{2047}$

normalisiert:  $z = \pm 1, x * 2^{e-127} = \pm 1, x * 2^{e-1023}$

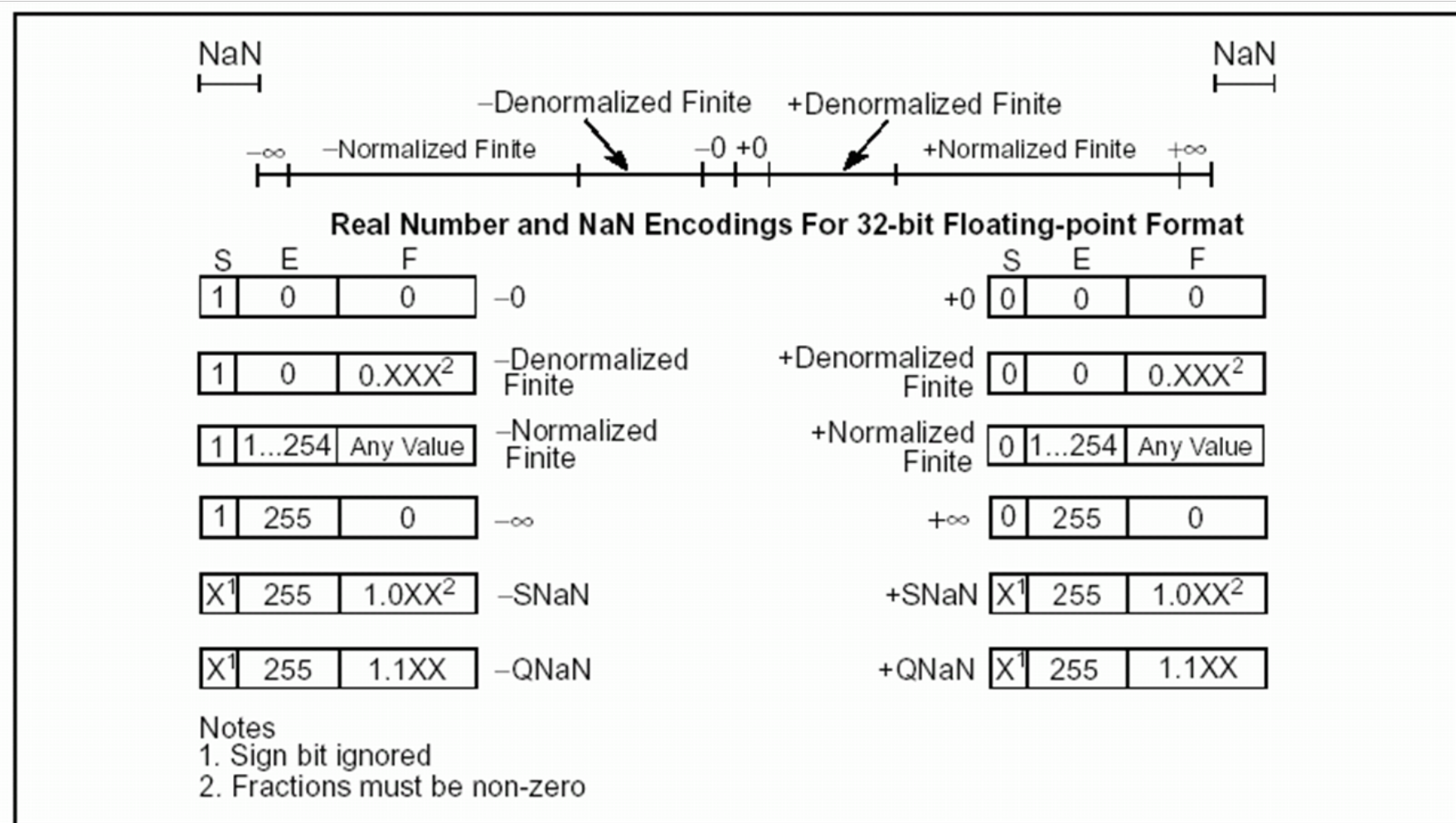
( $0 < e < 255$ )      ( $0 < e < 2047$ )

unnormalisiert:  $z = \pm 0, x * 2^{e-126} = \pm 0, x * 2^{e-1022}$

*SPrec.*      *DPrec.*

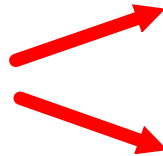
	Vorzeichen	Charakteristik	Mantisse
NaN	0	11.....11	 11 ..... 11 ⋮ 00 ..... 01
pos./neg. normalisierte Zahlen	0/1	11 ..... 10 ⋮ 00 ..... 01	11 ..... 11 ⋮ 00 ..... 00
$+\infty$	0	11 ..... 11	00 ..... 00
$\pm 0$	0/1	00 ..... 00	00 ..... 00
$-\infty$	1	11 ..... 11	00 ..... 00
NaN	1	11.....11	 SNaN 000 ..... 01 ⋮ QNaN 111 ..... 11

# Gleitkommazahlen



## alternative Zahldarstellungen

a) Zahlen als Strings/Zeichenkette



Nullterminiert  
(String endet beim Zeichen 0x00 = \00)

im ersten Byte wird eine Längenangabe  $L < 256$   
codiert, danach folgen L-BCD-codierte  
Dezimalziffern

b) Bibliotheken für beliebig genaues Rechnen (bel. große Exponenten und Mantissen)

in PERL:

```
use Math:BigInt; # lange Integerarithmetik  
use Math:BigRat; # Arithmetik mit langen rationalen Zahlen
```

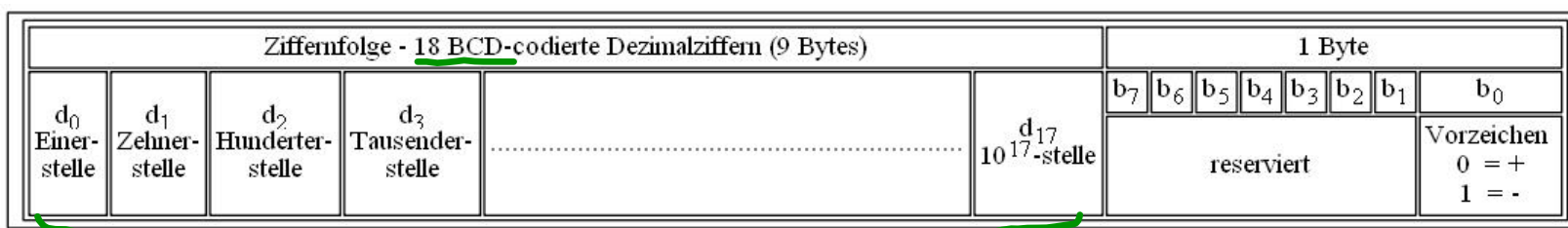
c) Zahlen durch Intervalle ersetzen

(Intervallararithmetik zum Einschließen der Rechenergebnisse in Intervalle)

d) BCD-Darstellung / BCD-Arithmetik

# BCD-Zahldarstellung (Binary Coded Decimals)

x87-BCD-Zahlenformat (10 Bytes)



BCD-Arithmetik:

Bsp:

$$\begin{array}{r} 87 \\ + 49 \\ \hline 136 \end{array}$$

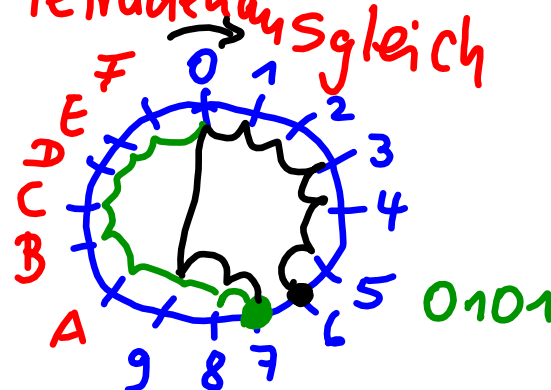
18 Dez.ziffern

$$\begin{array}{r} 1000 \quad 0111 \\ 0100 \quad 1001 \\ \hline 1101 \quad 0000 \\ \textcircled{13} = \\ \text{0110} \quad 0110 = 6 \\ \hline 0011 \quad 0110 \\ \text{"3} \quad \text{"6} \end{array}$$

Tetraden

$$\begin{array}{l} 0000 \quad 0 \\ \vdots \\ 1001 \quad 9 \end{array}$$

10 10 } keine Dez.ziffern  
Pseudotetraden  
Tetradenvergleich



# Arithmetik

$$1+1 = 2 = 10_2$$

$$1+1+1+1 = 4 = \underline{100}_2$$

Addition im Dual-, Oktal- und Hexadezimalsystem

	Dez	Oktal	Hexadez.	Dual
Summand	107,4	357	47,8	10111,11
+ Summand	47,7	+ 732	+ 27,A	1110,10 1111,01 1101,11
Summe	155,1	1311	6F,2	101101,10

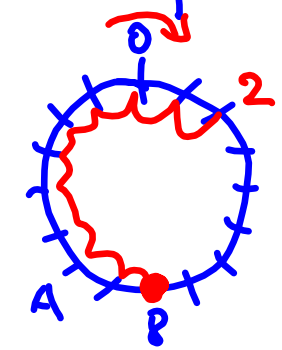
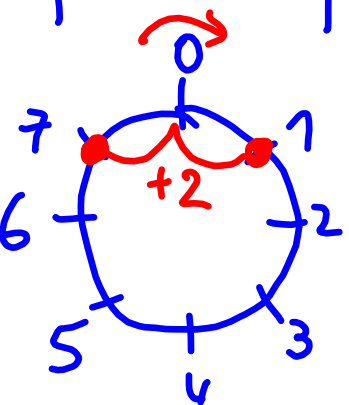
$$11 = 1 \cdot 8 + 3$$

$$7 + 2 = 9$$

$$= 1 \cdot 8 + 1$$

$$= 11_8$$

$$8 = 10_8$$



$$8 + A = 18$$

$$= 1 \cdot 16 + 2$$

$$= 12H$$

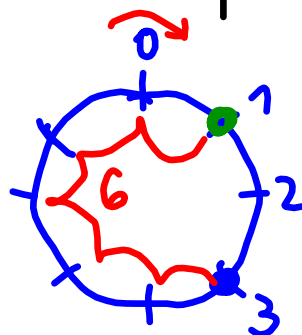
# Arithmetik

Subtraktion im Dual-, Oktal- und Hexadezimalsystem

	Dez	Okt.	Hex.	Dual
Minuend	193	301	C1	2
- Subtrahend	- 57	- 103	- 43	
Differenz	136	176	7E	

3-7  
13-7

Borgen  
einer 1,  
d.h. 10



$9 \cdot 3 = 11 \cdot 3 = 6$

$11H-3 = 17-3 = 14 = E$

$12-5 = 7$

$10_8 - 1 = 7$

# Arithmetik

Multiplikation im Dual-, Oktal- und Hexdezimalsystem

Handwritten multiplication of 78 (hex) by A6 (hex) in hexadecimal:

$$\begin{array}{r}
 78 \cdot A6 \\
 \hline
 150 \quad \vdots \\
 46 \quad \vdots \\
 \quad 30 \\
 \quad 2A \\
 \hline
 4DD0
 \end{array}$$

Handwritten multiplication of 78 (hex) by A6 (hex) in decimal:

$$\begin{array}{r}
 78 \cdot A6 \\
 \hline
 430 \\
 200 \\
 \hline
 4DD0
 \end{array}$$

$$\begin{aligned}
 A \cdot 8 &= 80_{10} \\
 &= 5 \cdot 16 + 0 \\
 &= 50H \\
 A \cdot 7 &= 70_{10} \\
 &= 4 \cdot 16 + 6 \\
 &= 46H \\
 6 \cdot 8 &= 30H \\
 6 \cdot 7 &= 42 = 32 + 10 \\
 &= 2A
 \end{aligned}$$

Duale  
kleine 1x1

	0	1	...	8	9	A	B	C	D	E	F
0											
1											
2											
...											
8											
9											
A											
B											
C											
D											
E											
F											

Handwritten annotations in the table:  
 - '30' is written in the row for '3' and column for '0'.  
 - '50' is circled in the row for '5' and column for '0'.  
 - The column headers 'A', 'B', 'C', 'D', 'E', 'F' are written in blue.

	0	1
0	0	0
1	0	1