

Rechnerarithmetik: Fließpunktzahlen

TechGl 2 - SoSe 2014

Genauigkeit und Formate

$$\text{Bias} = 2^{(\#e)-1} - 1$$

Mini-Float 16b	<table><tr><td>VZ</td><td>e</td><td>fraction</td></tr><tr><td>1b</td><td>5b</td><td>10b</td></tr></table>	VZ	e	fraction	1b	5b	10b	Bias = $2^{5-1} - 1 = 15$
VZ	e	fraction						
1b	5b	10b						

Singleprecision-Float 32b	<table><tr><td>VZ</td><td>e</td><td>fraction</td></tr><tr><td>1b</td><td>8b</td><td>23b</td></tr></table>	VZ	e	fraction	1b	8b	23b	Bias = $2^{8-1} - 1 = 127$
VZ	e	fraction						
1b	8b	23b						

Doubleprecision-Float 64b	<table><tr><td>VZ</td><td>e</td><td>fraction</td></tr><tr><td>1b</td><td>11b</td><td>52b</td></tr></table>	VZ	e	fraction	1b	11b	52b	Bias = $2^{11-1} - 1 = 1023$
VZ	e	fraction						
1b	11b	52b						

Berechnung von Fließpunktzahlen

Allgemein	<table><tr><td>VZ</td><td>e</td><td>fraction</td></tr></table>	VZ	e	fraction	$e = \text{exp} + \text{Bias}$
VZ	e	fraction			

$$n = (-1)^{\text{VZ}} \cdot (1.\text{fraction}) \cdot 2^{\text{exp}}$$

Mini-Float: $x = 1100111011001011_2$

→

1	10011	1011001011
---	-------	------------

 $10011 = \text{exp} + 15 \rightarrow \text{exp} = 4$

$$\begin{aligned} \rightarrow n &= (-1)^1 \cdot (1.1011001011) \cdot 2^4 \\ &= -(11011.001011) = -(11011 + 0.001011) \\ &= -(27 + (2^{-3} + 2^{-5} + 2^{-6})) = -27.171875_{10} \end{aligned}$$

Mini-Float: $n = 14.125_{10}$

$$\begin{aligned} &= 14_{10} + 0.125_{10} = 1110_2 + 0.001_2 \\ &= 1110.001_2 \\ &= (-1)^0 \cdot 1.110001_2 \cdot 2^3 \\ &\rightarrow \text{exp} = 3 \rightarrow \text{exp} + 15 = 18 = 10010 \\ &\rightarrow \text{VZ} = 0 \\ &\rightarrow \text{fraction} = 110001_2 = 1100010000_2 \end{aligned}$$

→

0	10010	1100010000
---	-------	------------

= 0100111100010000