

# Grundlagen der Informatik

Prof. Dr. J. Schmidt

Fakultät für Informatik

GDI – WS 2020/21

Zahlendarstellung – Multiplikation,  
Division, Gleitkommaarithmetik

- Welche charakteristischen Merkmale weist die binäre Multiplikation und Division auf?
- Wie werden in digitalen Rechenanlagen „reelle“ Zahlen dargestellt?



- Rechenregeln für die Multiplikation zweier binärer Zahlen
  - $0 \cdot 0 = 0$
  - $0 \cdot 1 = 0$
  - $1 \cdot 0 = 0$
  - $1 \cdot 1 = 1$
- Multiplikation mehrstelliger Zahlen
  - Multiplikation des Multiplikanden mit den einzelnen Stellen des Multiplikators
  - Stellenrichtige Addition der Zwischenergebnisse

Identisch mit den Regeln der logischen UND-Verknüpfung



# Binäre Multiplikation (2)

## Kapitel 2: Zahlendarstellung

Beispiel:  $10 \cdot 13$

$$\begin{array}{r} 1010 \cdot 1101 \\ \hline 1010 \\ 10100 \\ 00000 \\ 101000 \\ \hline 10000010 \end{array}$$

**Ergebnis: 130**



# Binäre Multiplikation (3)

## Kapitel 2: Zahlendarstellung

Beispiel:  $17.375 \cdot 9.75$

$$\begin{array}{r} 1\ 0\ 0\ 0\ 1.0\ 1\ 1\ \cdot\ 1\ 0\ 0\ 1.1\ 1 \\ \hline \phantom{1\ 0\ 0\ 0\ 1.0\ 1\ 1\ \cdot\ 1\ 0\ 0\ 1.1\ 1} 1\ 0\ 0\ 0\ 1\ 0\ 1\ 1 \\ \phantom{1\ 0\ 0\ 0\ 1.0\ 1\ 1\ \cdot\ 1\ 0\ 0\ 1.1\ 1} 1\ 0\ 0\ 0\ 1\ 0\ 1\ 1 \\ \phantom{1\ 0\ 0\ 0\ 1.0\ 1\ 1\ \cdot\ 1\ 0\ 0\ 1.1\ 1} 1\ 0\ 0\ 0\ 1\ 0\ 1\ 1 \\ \phantom{1\ 0\ 0\ 0\ 1.0\ 1\ 1\ \cdot\ 1\ 0\ 0\ 1.1\ 1} 1\ 0\ 0\ 0\ 1\ 0\ 1\ 1 \\ \phantom{1\ 0\ 0\ 0\ 1.0\ 1\ 1\ \cdot\ 1\ 0\ 0\ 1.1\ 1} 1\ 0\ 0\ 0\ 1\ 0\ 1\ 1 \\ \hline \end{array}$$

**Zwischenergebnis**

1 0 1 0 1 0 0 1 . 0 1 1 0 1

Stellenrichtiges Einfügen des Kommas



**Ergebnis:  $169.40625_{(10)}$**



- Rechenregeln für die Division zweier binärer Zahlen
  - $0 : 0 =$  nicht definiert
  - $0 : 1 = 0$
  - $1 : 0 =$  nicht definiert
  - $1 : 1 = 1$
- Durchführung der binären Division analog dem im Zehnersystem üblichen Verfahren



# Binäre Division (2)

## Kapitel 2: Zahlendarstellung

Beispiel  $20 : 6$

$$1\ 0\ 1\ 0\ 0 : 1\ 1\ 0 = 1\ 1.0\ 1\ 0\ 1\ \dots$$

$$\begin{array}{r} 1\ 0\ 1\ 0\ 0 \\ -1\ 1\ 0 \\ \hline \end{array}$$

$$1\ 0\ 0\ 0$$

$$\begin{array}{r} 1\ 0\ 0\ 0 \\ -1\ 1\ 0 \\ \hline \end{array}$$

$$1\ 0\ 0\ 0$$

$$\begin{array}{r} 1\ 0\ 0\ 0 \\ -1\ 1\ 0 \\ \hline \end{array}$$

$$1\ 0\ 0\ 0$$

$$\begin{array}{r} 1\ 0\ 0\ 0 \\ -1\ 1\ 0 \\ \hline \end{array}$$

...

**Ergebnis:  $3.333\dots_{(10)}$**



- **Multiplikation** wird mittels **wiederholter Addition** durchgeführt
- **Division** wird mittels **wiederholter Subtraktion** durchgeführt
- **Sonderfall:**  
Multiplikator oder Divisor sind eine Zweierpotenz  $2^k$ 
  - Multiplikation bzw. Division kann einfacher und schneller erfolgen
  - Durch Verschiebung von entsprechend vielen Bits ( $k$ ) nach links bzw. rechts
  - Bei  $2^1$  um **1** Bit,  $2^2$  um **2** Bits, bei  $2^3$  um **3** Bits usw.





# Realisierung im Rechner (2)

## Kapitel 2: Zahlendarstellung

### Beispiele

- $13 \cdot 4$

$$1101 \cdot 100 = 110100$$

				1	1	0	1
			1	1	0	1	0
		1	1	0	1	0	0

2 Bits nach links

= 52

- $20 \cdot 8$

$$10100 \cdot 1000 = 10100000$$

3 Bits nach links

- $20 : 4$

$$10100 : 100 = 101$$

2 Bits nach rechts

- $26 : 4$

$$11010 : 100 = 110 \quad (\text{Rest } 2) \quad (\rightarrow \text{evtl. Informationsverlust!})$$

2 Bits nach rechts



# Realisierung im Rechner (3)

## Kapitel 2: Zahlendarstellung

- Ablage der zu verschiebenden Zahl in einem dem Rechenwerk direkt zugeordneten Speicherplatz (**Register**, Akkumulator)
- Register besitzen meist Übertrags-Bit (**Carry**)
- Beispiel:  
Verschieben  $1101_{(2)}$  um 1 Stelle nach rechts

Register

0	0	0	0	1	1	0	1
---	---	---	---	---	---	---	---

0	0	0	0	0	1	1	0
---	---	---	---	---	---	---	---

Carry

0
---

1
---

Übertrags-Bit wird von  
0 auf 1 gesetzt  
(Setzen eines Flag)



- In Programmiersprachen wird anstatt des im Deutschen üblichen Kommas der **Punkt** verwendet, um ganzzahligen Teil vom gebrochenen Teil einer „reellen“ Zahl abzutrennen
- Zwei unterschiedliche Typen
  - Festkommazahlen (fixed point)
  - Gleitkommazahlen (floating point)



Punkt (das Komma) steht immer an einer bestimmten festgelegten Stelle

- wobei der Punkt nicht eigens mitgespeichert wird

$$Z = z_{n-1}z_{n-2} \cdots z_1z_0 . z_{-1}z_{-2} \cdots z_{-m} \quad (2)$$

$$Z = \sum_{i=-m}^{n-1} z_i \cdot 2^i$$

- $Z$  hat die Länge  $n + m$   
**n** Stellen vor dem und **m** Stellen nach dem Komma



- Nachteile der Festkomma-Darstellung
  - Mit einer bestimmten Anzahl von Bits kann nur ein beschränkter Wertebereich abgedeckt werden
  - Die Stelle des Punkts (Kommas) muss allgemein festgelegt werden

(Wo soll man diesen festlegen, wenn manchmal mit sehr kleinen, hochgenauen Werten und ein anderes Mal mit sehr großen Werten gearbeitet werden muss?)
- Festkommadarstellung wird nur in Rechnern verwendet, die für Spezialanwendungen benötigt werden
  - sonst Gleitkommadarstellung



- Jede reelle Zahl kann beispielsweise in folgender Form angegeben werden

$$2.3756 \cdot 10^3$$

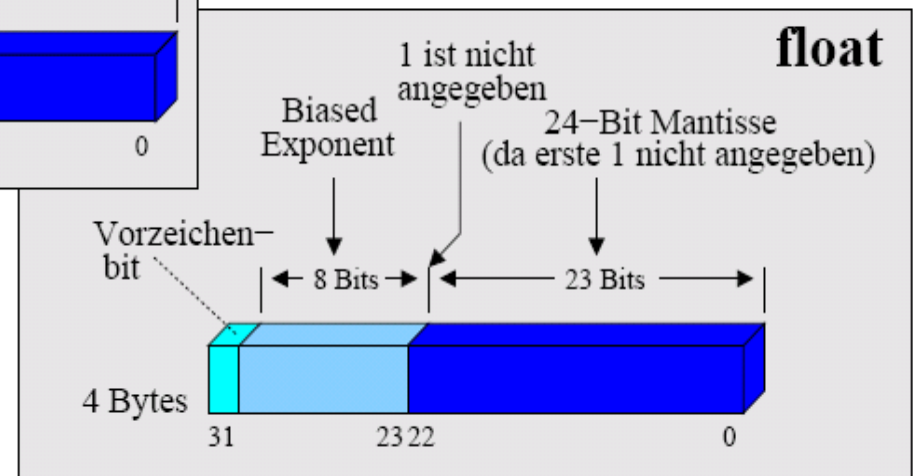
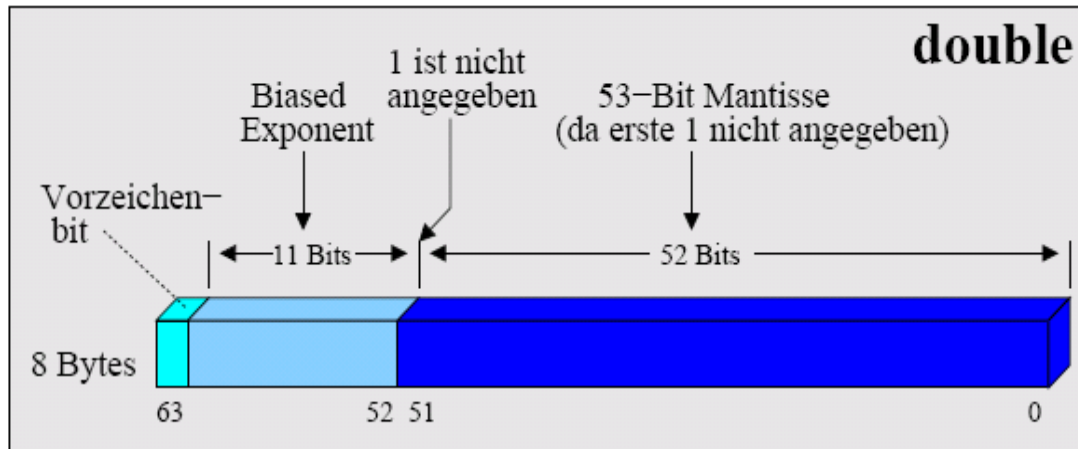
- Zwei Bestandteile
  - **Mantisse** (2.3756) und
  - **Exponent** (3), der ganzzahlig ist
- Wird in den meisten Rechnern verwendet
  - aber nicht Basis 10, sondern Basis 2
- Die für die Darstellung verwendete Anzahl von Bytes legt fest, ob man mit
  - **einfacher** (Datentyp **float**) oder mit
  - **doppelter Genauigkeit** (Datentyp **double**) arbeitet



# Gleitkommazahlen (2)

## Kapitel 2: Zahlendarstellung

- IEEE-Format (binär, Norm IEEE 754-2019)
  - C/C++- und Java-Datentypen verwenden vier Bytes für **float** und acht Bytes für **double**
  - in der Norm werden auch Darstellungen für halbe (16 Bit) und vierfache Genauigkeit (128 Bit) definiert



- IEEE-Format

- geht von **normalisierten Gleitkommazahlen** aus

- Normalisierung

- der Exponent wird so verändert, dass
  - der gedachte Dezimalpunkt immer rechts von der ersten Nicht-Null-Ziffer liegt  
(im Binärsystem ist dies immer eine 1)





### Beispiel

- $17.625_{(10)}$

$$= 16 + 1 + \frac{1}{2} + \frac{1}{8}$$

$$= 10001.101_{(2)}$$

$$= 10001.101 \cdot 2^0$$

- Normalisierte Form

- Schiebe Dezimalpunkt hinter die erste signifikante Ziffer
- Passe Exponent entsprechend an

$$= 1.0001101 \cdot 2^4$$



Aus Normalisierung (IEEE-Format) ergibt sich

- In der **Mantisse** steht das höchstwertige „Einser-Bit“ immer links vom gedachten Dezimalpunkt
  - außer für 0.0 und einige andere Sonderfälle
  - dieses Bit wird nicht gespeichert
- **Exponent** ist eine ganze Zahl, welche (nach Addition eines **Bias**) ohne Vorzeichen dargestellt wird
  - Durch **Bias**-Addition wird für den Exponent keine Vorzeichenrechnung benötigt (immer positiv)
  - Wert von **Bias** hängt vom Genauigkeitsgrad ab
    - float (mit 4 Bytes, 8 Bits für Exponent): Bias = 127
    - double (mit 8 Bytes, 11 Bits für Exponent): Bias = 1023
- **Vorzeichenbit** zeigt das Vorzeichen der Mantisse
  - Mantisse immer als Betragswert
  - auch im negativen Fall nicht als Komplement



### Beispiel

- 17.625 ( $1.0001101 \cdot 2^4$ )

3 1	3 0	2 9	2 8	2 7	2 6	2 5	2 4	2 3	2 2	2 1	2 0	1 9	1 8	1 7	1 6	1 5	1 4	1 3	1 2	1 1	1 0	9	8	7	6	5	4	3	2	1	0
0	1	0	0	0	0	0	1	1	0	0	0	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

- Biased Exponent: 10000011 = 131
- Bias: 01111111 = 127
- Wirklicher Exponent: 00000100 = 4



Nach IEEE gilt für **float** (einfach)  
und **double** (doppelt)

	Einfach	Doppelt
Vorzeichen Bits	1	1
Exponenten Bits	8	11
Mantissen Bits	23	52
Bit insgesamt	32	64
BIAS	127	1023
Exponentenbereich	$[-126, 127]$	$[-1022, 1023]$



### Null

- positive (+0.0) und negative (−0.0) Null
  - bei Vergleichsoperationen sind beide gleich
- Verwendung
  - Darstellung der Null
  - Rundung auf  $\pm 0.0$  bei Unterlauf (Loch um die Null)
- Darstellung
  - (biased) Exponent = 0
  - Mantisse = 0



### Unendlich (INF)

- plus ( $+\infty$ ) und minus ( $-\infty$ ) Unendlich
- Verwendung
  - Zahlen, deren Betrag zu groß ist um noch dargestellt werden zu können (Überlauf)
  - Rechnungen, die per Definition Unendlich ergeben (z.B. Division einer Zahl  $z \neq 0$  durch Null:  $z / 0.0 = \infty$ )
- Darstellung
  - (biased) Exponent = 111....1
  - Mantisse = 0



### keine Zahl (NaN)

- Verwendung

- Darstellung ungültiger Werte
- Rechnungen, die undefinierte Ergebnisse liefern
  - $0.0 / 0.0 = \text{NaN}$
  - $\infty / \infty = \text{NaN}$
  - $\sqrt{-3} = \text{NaN}$

- Darstellung

- (biased) Exponent  $= 111\dots 1$
- Mantisse  $> 0$



Geben Sie für die Dezimalzahl 125.875 die `float`-Darstellung (nach IEEE-Format) an.

1. Bestimmung Bias =  $127_{(10)}$
2. Umwandlung in duale Festkommazahl ohne Vorzeichen
  - Anteil vor Komma:  $125_{(10)} = 1111101_{(2)}$
  - Anteil nach Komma:  $0.875_{(10)} = 0.111_{(2)}$   
 $0.875 \cdot 2 = 1.75 \rightarrow 1$   
 $0.75 \cdot 2 = 1.5 \rightarrow 1$   
 $0.5 \cdot 2 = 1.0 \rightarrow 1$
3. Normieren  
 $1111101.111 \cdot 2^0 = 1.111101111 \cdot 2^6$
4. Berechnung dualer Exponent  
 $2^6 \rightarrow 6_{(10)} + \text{Bias} = 133_{(10)} = 10000101_{(2)}$
5. Bestimmung Vorzeichen Bit  
Positiv  $\rightarrow 0$
6. Zusammensetzen

V	Exponent	Mantisse
0	10000101	111101111000000000000000





# Ungenauigkeit von Gleitkommazahlen (1)

## Kapitel 2: Zahlendarstellung

- Gleitkommazahlen, die im Dezimalsystem genau darstellbar sind, können im Dualsystem nicht immer genau dargestellt werden
  - Daraus resultieren **Ungenauigkeiten**
  - Merke: `float`- und `double`-Werte sollte man niemals auf Gleichheit prüfen!
- Beispiel
  - Ausgeben der Zahlen 0.1 bis 1.0 mit einer Schrittweite von 0.1
  - `float`- oder `double`-Variable als Laufvariable bei Iterationen



# Ungenauigkeit von Gleitkommazahlen (2)

## Kapitel 2: Zahlendarstellung

Endlosschleife, da die Abbruchbedingung nie erreicht wird

```
#include <stdio.h>

int main(void)
{
    float i = 0.1;

    while (i != 1.0)
    {
        printf("%.10f\n", i);
        i = i + 0.1;
    }
    return 0;
}
```



## Gewünschtes Ergebnis wird erreicht

```
#include <stdio.h>

const float EPSILON = 1e-6;

int main(void)
{
    float i = 0.1;

    while (i <= 1.0+EPSILON)
    {
        printf("%.10f\n", i);
        i = i + 0.1;
    }
    return 0;
}
```



# Ungenauigkeit von Gleitkommazahlen (4)

Bessere Methode:  
Umgehen der Ungenauigkeit durch Arbeiten mit  
ganzzahligen Laufvariablen

```
#include <stdio.h>

int main(void) {
    int i = 1;

    while (i <= 10)
    {
        printf("%.10f\n", (float)i/10);
        i = i + 1;
    }
    return 0;
}
```



# Ungenauigkeit von Gleitkommazahlen (5)

## Kapitel 2: Zahlendarstellung

- Gleitkommaarithmetik ist nicht assoziativ!

- $$\begin{array}{lcl} (u + v) + w & \neq & u + (v + w) \\ (u * v) * w & \neq & u * (v * w) \end{array}$$

- ...und auch nicht distributiv

- $$u * (v + w) \neq (u * v) + (u * w)$$

- Beispiel (dezimal, 8 Stellen genau)

- $$\begin{array}{lcl} (11111113. + (-11111111.)) + 7.51111111 & = & \\ 2.0000000 & + 7.51111111 & = 9.51111111 \end{array}$$
- $$\begin{array}{lcl} 11111113. + (-11111111. + 7.51111111) & = & \\ 11111113. + (-111111103.) & = & 10.000000 \end{array}$$

