


```
float f = 0.1F;
printf(" f = %56.40f\n", f);
double d = 0.1;
printf(" d = %56.40lf\n", d);
```

liefert auf unserem Rechner die folgenden Werte:

```
f =          0.1000000014901161193847656250000000000000
d =          0.10000000000000000055511151231257827021182
```

Der binär gespeicherte Wert ist $(1, 1001100110011001101)_2 \cdot 2^{-4}$ in normalisierter Darstellung. Das erste Bit vor dem Komma ist das Hidden-Bit. Die hinteren, niederwertigen Stellen von $0,1 = (0, \overline{00011})_2$ wurden nicht einfach abgeschnitten, denn dann hätte sich ein Wert kleiner als $0,1$ ergeben, sondern es wurde gerundet.

Rundung Bei binären Rundungen muss laut IEEE 754 zur nächstgelegenen darstellbaren Zahl gerundet werden. Wenn diese nicht eindeutig definiert ist (genau in der Mitte zwischen zwei darstellbaren Zahlen), wird so gerundet, dass das niederwertigste Bit der Mantisse 0 wird. Statistisch wird dabei in 50% der Fälle auf-, in den anderen 50% der Fälle abgerundet. (laut wikipedia)

Die obige Ausgabe des C-Programms kann je nach Hardware, Betriebssystem und Compiler variieren, da C keine exakten Vorgaben für die Datentypen macht, sondern nur Mindestanforderungen definiert.

Arithmetische Operationen Eine weitere Ungenauigkeit tritt bei der Addition von Gleitkommazahlen auf, und zwar durch die Denormalisierung der kleineren Zahl. Das Rechnen mit solchen Zahlen wird auch in dem Buch *Eine Einführung in die Mathematik an Beispielen aus der Informatik* von Goebels und Rethmann beschrieben. Hier ein kurzer Auszug daraus:

Jetzt rechnen wir mit Gleitpunktzahlen. Schauen wir uns zunächst die Multiplikation und Division für Zahlen $x = m_x \cdot 2^{d_x}$ und $y = m_y \cdot 2^{d_y}$ an. Aus den Rechenregeln für Potenzen reeller Zahlen ergibt sich:

- Bei der Multiplikation werden die Mantissen multipliziert und die Exponenten addiert.

$$x \cdot y = (m_x \cdot m_y) \cdot 2^{d_x + d_y}$$

- Bei der Division werden die Mantissen dividiert und die Exponenten subtrahiert.

$$x : y = (m_x : m_y) \cdot 2^{d_x - d_y}$$

Betrachten wir auch dazu wieder ein Beispiel, zunächst für die Multiplikation $121,8 \cdot 0,37$ im Dezimalsystem, weil uns das vertrauter ist:

$$(1,218 \cdot 10^2) \cdot (3,7 \cdot 10^{-1}) = 4,5066 \cdot 10^1 = 45,066$$

Und nun ein Beispiel für die Division $450,66 : 37$ im Dezimalsystem:

$$(4,5066 \cdot 10^2) : (3,7 \cdot 10^1) = 1,218 \cdot 10^1 = 12,18.$$

Die Multiplikation und Division der Zahlen erfordert keine angepasste Darstellung. Das ist bei der Addition und Subtraktion der Zahlen $x = m_x \cdot 2^{d_x}$ und $y = m_y \cdot 2^{d_y}$ anders. Wenn wir schriftlich addieren und subtrahieren, dann muss das Komma bzw. der Dezimalpunkt an der gleichen Stelle stehen. Bei Gleitpunktzahlen wird die Position des Kommas durch den Exponenten bestimmt, so dass wir diesen angleichen müssen. Rechnerisch sieht das dann so aus:

$$x \pm y = \begin{cases} (m_x \cdot 2^{d_x-d_y} \pm m_y) \cdot 2^{d_y}, & \text{falls } d_x \leq d_y \\ (m_x \pm m_y \cdot 2^{d_y-d_x}) \cdot 2^{d_x}, & \text{falls } d_x > d_y. \end{cases}$$

Wie man hier sieht, wird der kleine Exponent dem großen Exponenten angeglichen. Es macht keinen Sinn, den großen Exponenten an den kleinen anzupassen, da dadurch eventuell die höchstwertigen Stellen nicht mehr dargestellt werden könnten. Man nimmt lieber in Kauf, niederwertigere Stellen zu verlieren, und akzeptiert damit Rundungsungenauigkeiten. Bei einer Mantisse von 4 Stellen (ohne Vorzeichen) zur Basis 10 seien x, y, z wie folgt gegeben:

$$x := +1,235 \cdot 10^2, \quad y := +5,512 \cdot 10^4, \quad z := -5,511 \cdot 10^4.$$

Dann erhalten wir mit der Rundungsungenauigkeit $1,235 \cdot 10^{-2} \approx 0,012$:

$$\begin{aligned} x + y &= +1,235 \cdot 10^2 + 5,512 \cdot 10^4 = (+1,235 \cdot 10^{-2} + 5,512) \cdot 10^4 \\ &\approx (+0,012 + 5,512) \cdot 10^4 = +5,524 \cdot 10^4, \end{aligned}$$

$$(x + y) + z \approx +5,524 \cdot 10^4 - 5,511 \cdot 10^4 = +0,013 \cdot 10^4 = +1,300 \cdot 10^2.$$

Andererseits tritt bei der folgenden Rechnung kein Rundungsfehler auf:

$$y + z = +5,512 \cdot 10^4 - 5,511 \cdot 10^4 = +0,001 \cdot 10^4 = +1,000 \cdot 10^1,$$

$$\begin{aligned} x + (y + z) &= +1,235 \cdot 10^2 + 1,000 \cdot 10^1 = +1,235 \cdot 10^2 + 0,100 \cdot 10^2 \\ &= +1,335 \cdot 10^2 \neq +1,300 \cdot 10^2 = (x + y) + z. \end{aligned}$$

Das Assoziativgesetz gilt also nicht bei Gleitpunktzahlen mit endlicher Stellenzahl. Es ist sinnvoll, Zahlen mit ähnlich großen Exponenten zu addieren oder subtrahieren.

Ein anderes Beispiel, bei dem die beschränkte Stellenzahl der Mantisse Auswirkungen hat, erhalten wir bei Berechnung der linken Seite von

$$\frac{(x + y)^2 - x^2 - 2xy}{y^2} = \frac{x^2 + 2xy + y^2 - x^2 - 2xy}{y^2} = 1$$

auf einem 32 Bit Computer mit einfacher Genauigkeit (`float`) für die Zahlen $x = 1000$ und $y = 0,03125 = \frac{1}{32}$. Obwohl beide Zahlen exakt dargestellt werden können, bekommen wir das Ergebnis 0. Das liegt daran, dass $(x + y)^2$ eine große Zahl und y^2 eine im Vergleich so kleine Zahl ist, dass in der Zahlendarstellung aufgrund der Stellenzahl der Mantisse genau der y^2 -Anteil von $(x + y)^2 = x^2 + 2xy + y^2$ nicht mehr gespeichert werden kann. Daher entsteht bereits im Zähler die Null. Um solche groben Fehler zu vermeiden, sollten Zahlen bei Addition und Subtraktion ungefähr die gleiche Größenordnung haben.

Wir können festhalten, dass Ungenauigkeiten im Umgang mit Gleitpunktzahlen sowohl bei der Umwandlung vom Dezimal- ins Dualsystem als auch bei den arithmetischen Operationen auftreten. (Ende des Auszugs aus dem oben genannten Buch.)

Schauen wir uns einmal an, was bei einer fortgesetzten Multiplikation mit 10 passiert. Dazu schreiben wir zunächst ein kleines Programm:

Schauen wir uns nun an, wie die Ungenauigkeit bei der letzten Multiplikation mit $10 = (1010)_2$ zustande kommt. Es gilt $3007812,5 = (1,01101111001010100010010)_2 \cdot 2^{21}$. Das kann bspw. unter <https://www.matheretter.de/rechner/dezimalbinar> geprüft werden.

$$\begin{array}{r}
 1, 0110 \ 1111 \ 0010 \ 1010 \ 0010 \ 010 \quad [\cdot 2^{21}] \quad \cdot (1010)_2 \\
 \hline
 = 1, 0110 \ 1111 \ 0010 \ 1010 \ 0010 \ 010 \quad [\cdot 2^{22}] \\
 + 1, 0110 \ 1111 \ 0010 \ 1010 \ 0010 \ 010 \quad [\cdot 2^{24}] \\
 \hline
 = 0, 0101 \ 1011 \ 1100 \ 1010 \ 1000 \ 10010 \quad [\cdot 2^{24}] \quad (\text{Denormalisieren}) \\
 + 1, 0110 \ 1111 \ 0010 \ 1010 \ 0010 \ 010 \quad [\cdot 2^{24}] \\
 \hline
 = 1, 1100 \ 1010 \ 1111 \ 0100 \ 1010 \ 11010 \quad [\cdot 2^{24}] \quad (\text{Abrunden}) \\
 = 1, 1100 \ 1010 \ 1111 \ 0100 \ 1010 \ 110 \quad [\cdot 2^{24}] \\
 = 30.078.124
 \end{array}$$

Problem hier ist, dass die Zahlen zu groß werden und mit einer 23-bit großen Mantisse nicht dargestellt werden können. Die Mantisse besteht nur aus 23 Bit, aber die Mantisse wird mit 2^{24} multipliziert, das Komma also um 24 Stellen nach rechts verschoben. Ungerade Zahlen können so nicht mehr dargestellt werden, es fehlt die entsprechende Stelle in der Mantisse. Es entstehen also „Lücken“ im Zahlenbereich. Man sieht das auch an dem folgenden Programm:

```

#include <stdio.h>

int main(void) {
    int x = 16777212;
    float f = 16777212.0;

    printf("      x      |      f\n");
    printf("-----+-----\n");
    for (int i = 0; i < 10; i++) {
        printf(" %9d | %12.f\n", x, f);
        x += 1;
        f += 1.0f;
    }
    return 0;
}

```

Hier die Ausgabe des Programms:

```

      x      |      f
-----+-----
 16777212 | 16777212.0
 16777213 | 16777213.0
 16777214 | 16777214.0
 16777215 | 16777215.0
 16777216 | 16777216.0
 16777217 | 16777216.0
 16777218 | 16777216.0
 16777219 | 16777216.0
 16777220 | 16777216.0
 16777221 | 16777216.0

```

Gleitpunktzahlen vom Typ `float` haben Einbußen hinsichtlich der Genauigkeit im Vergleich zu ganzen Zahlen vom Typ `int`, da die Zahlenbereiche unterschiedlich groß sind, aber die selbe Anzahl Bit zur Codierung genutzt werden:

Wir erhalten also den exakten Bruch, so wie gewünscht. Leider hat auch dieses Verfahren einen Nachteil. Betrachten wir die Ausgabe des Programms für den Wert 0,1:

```

0.1000000014901161193847656
0.2000000029802322387695312
0.4000000059604644775390625
0.8000000119209289550781250
1.6000000238418579101562500
3.2000000476837158203125000
6.4000000953674316406250000
12.8000001907348632812500000
25.6000003814697265625000000
51.2000007629394531250000000
102.4000015258789062500000000
204.8000030517578125000000000
409.6000061035156250000000000
819.2000122070312500000000000
1638.4000244140625000000000000
3276.8000488281250000000000000
6553.6000976562500000000000000
13107.2001953125000000000000000
26214.4003906250000000000000000
52428.8007812500000000000000000
104857.6015625000000000000000000
209715.2031250000000000000000000
419430.4062500000000000000000000
838860.8125000000000000000000000
1677721.6250000000000000000000000
3355443.2500000000000000000000000
6710886.5000000000000000000000000
13421773.0000000000000000000000000
0.1 = 13421773 / 134217728

```

Es ergeben sich außer der initialen Ungenauigkeit bei der Umwandlung ins Binärsystem keine weiteren Ungenauigkeiten. Leider, denn der so ermittelte Bruch ist nur eine Näherung für den Wert 0,1.

Schauen wir uns noch ein weiteres Beispiel an, bei dem durch Rechenungenauigkeiten bei der Multiplikation mit 10 ein exakter Bruch entsteht. Auch der Wert 2,0815 kann binär nicht exakt dargestellt werden. Allerdings verschwinden auch hier irgendwann die „überflüssigen“ Nachkommastellen.

```

2.0815000534057617187500000
20.8150005340576171875000000
208.1500091552734375000000000
2081.5000000000000000000000000
20815.0000000000000000000000000

```

Es gilt $2,0815 = (10,000101001101110100101110\dots)_2$, was bei einer 23 bit großen Mantisse und einem Hidden-Bit zur Zahl $2,08150005340576171875 = (1,00001010011011101001100)_2 \cdot 2^1$

