

Fliesskommazahlen für Schüler der Gymnasialstufe

amaximov

März 2020

Inhaltsverzeichnis

1 Einführung	1
2 Fliesskommazahlen	5
2.1 Kleinste und grösste positive Zahlen	5
2.2 Darstellbare Zahlen	9
3 Addition	12
4 Zusammenfassung	16
5 Beispiellösungen	17
5.1 Fliesskommazahlen	17
5.2 Addition	18

1 Einführung

Wie können Computer so viele unterschiedliche Dinge mit nur Nullen und Einsen darstellen? Auf dem Bildschirm sehen wir Texte, Bilder, Videos, die wir noch dazu verändern können.

Vorletzte Woche haben wir gesehen, wie Computer ganze Zahlen speichern und manipulieren. Sie stellen die Zahlen in Basis 2 dar und speichern eine feste Anzahl Bits. Letzte Woche haben wir angefangen, uns mit der Darstellung der reellen Zahlen zu beschäftigen, mit den Fliesskommazahlen. Heute machen wir damit weiter.

Die Fliesskommazahlen, wie wir gesehen haben, sind im Grunde genommen nichts anderes, als die Exponentialschreibweise, die wir schon aus Chemie und Physik kennen: Die signifikanten Stellen werden mit der Basis hoch einen Exponenten multipliziert. Im Unterschied zu Menschen, arbeitet der Computer meistens in der Basis 2 statt 10 und schränkt die Anzahl der möglichen signifikanten Stellen und Exponenten ein.

Um diese Einschränkungen sichtbar und anfassbar zu machen, hatten wir das "Kasten-Seil"-Modell eingeführt.

Stellen wir uns reelle Zahlen in Basis 2 vor:

$\dots 1001,00000\dots$
 $\dots 0000,10010\dots$
 $\dots 0001,00100\dots$

Die Signifikanten Stellen, die wir speichern wollen, tun wir in einen Kasten.

$\dots \boxed{1001} 00000\dots$
 $\dots 0000, \boxed{10010}\dots$
 $\dots 000 \boxed{10010} 0\dots$

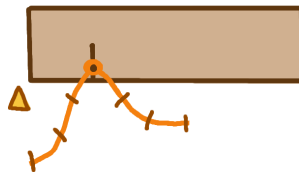
Um zu wissen, woher wir diese signifikanten Stellen rausgenommen haben, markieren wir auf einem Seil den Abstand bis zum Komma.

$\dots \boxed{1001} 00000\dots$
 $\dots 0000, \boxed{10010}\dots$
 $\dots 000 \boxed{10010} 0\dots$

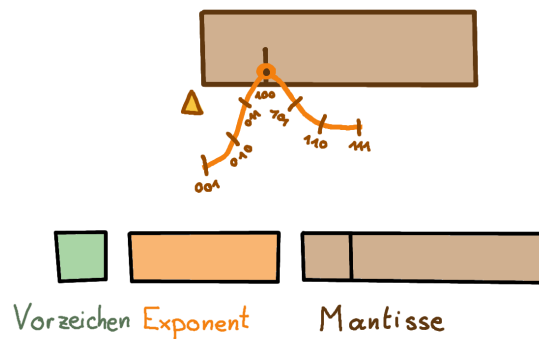
Jetzt haben wir alle Informationen gespeichert, die wir brauchen, um diese

Zahl wiederherzustellen. Damit diese Darstellung eindeutig ist, verlangen wir, dass das erste Bit der Mantisse eine Eins ist.

Folgende Elemente charakterisieren ein Fließkommazahlensystem: Die Grösse vom Kasten, d.h. die Mantissenlänge, und die Länge vom Seil, d.h. der Exponentenbereich. Die Bits im Kasten und die Markierung am Seil stellen eine Zahl dar.



Der Computer hat intern keine Kasten und keine Seile. Er arbeitet mit Bitmuster. Die Bits werden in 3 Bereichen aufgeteilt: Vorzeichen (grün auf dem Bild), Exponent (Orange auf dem Bild) und Mantisse (braun auf dem Bild). Im Mantissenteil werden die Bits aus dem Kasten gespeichert. Im Exponententeil wird die Markierung am Seil kodiert. Im Vorzeichenteil wird das Vorzeichen kodiert (0 für positive Zahlen und 1 für negative).

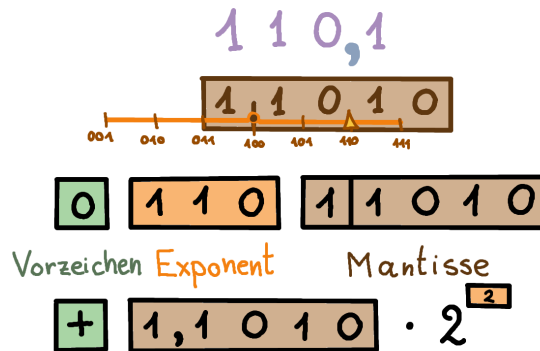


Beispiel 1.1. Wir werden jetzt zusammen die Zahl 6.5 im Fließkommazahlensystem mit Mantissenlänge 5 und Exponenten von -3 bis 3 darstellen.

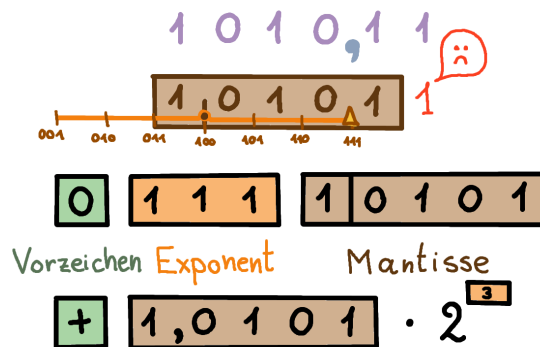
Die reelle Zahl in Basis 2 ist 110.1 .

Der Kasten hat 5 Plätze. Alle signifikanten Stellen haben dort Platz. Dann verbinden wir das Seil mit dem Komma und setzen eine Markierung. Das lässt

sich direkt ins Bitmuster übersetzen: Das Vorzeichen ist positiv, die Kodierung vom Exponenten lässt sich am Seil ablesen, die Mantisse speichert man direkt.



Beispiel 1.2. Nicht alle Zahlen lassen sich im Fließkommazahlensystem genau darstellen. Manche müssen gerundet werden. Zum Beispiel, die Zahl 10.75 sieht in Binär so aus: 1010.11. Sie hat 6 signifikante Stellen, aber nur 5 haben in der Mantisse Platz. Die letzte Eins kann nicht gespeichert werden und geht verloren.



2 Fließkommazahlen

2.1 Kleinste und grösste positive Zahlen

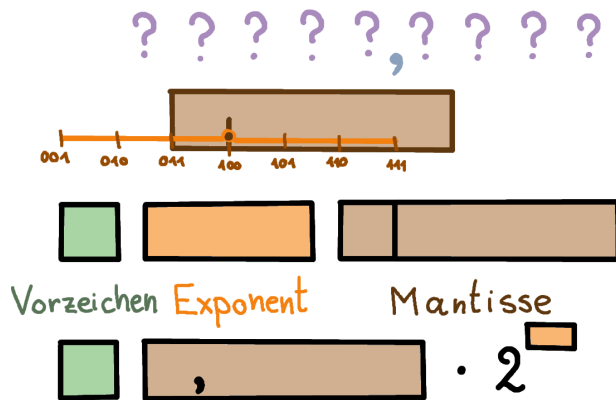
Es gibt unendlich viele reelle Zahlen. Wir können aber nur endlich viele davon in einem Fließkommazahlensystem darstellen. Im Kasten, welcher uns die Mantisse veranschaulicht, finden nur endlich viele Bits platz. Das Seil, welches den Kasten an den Komma bindet und welches uns den Exponenten veranschaulicht, hat eine endliche Länge.

Weil es endlich viele darstellbare Zahlen gibt, muss es eine kleinste und eine grösste Zahl geben. In diesem Abschnitt werden wir die kleinste und die grösste positive darstellbare Zahlen finden.

Beispiel 2.1. Wir konstruieren die grösste positive Zahl im Fließkommazahlensystem mit Mantissenlänge 5 und Exponenten von -3 bis 3 .

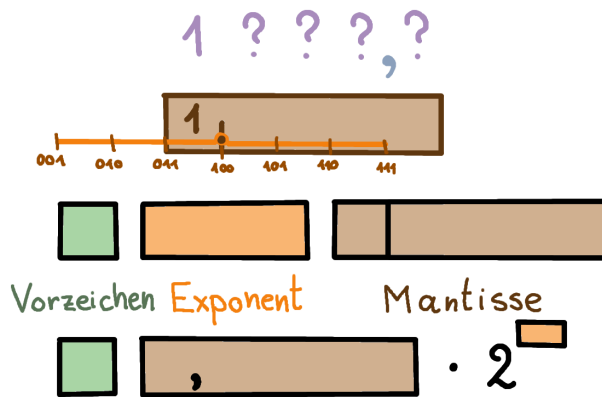
In violett wird die reelle Zahl in Basis 2 aufgeschrieben, in der zweiten Zeite kommt die "Kasten-und-Seil"-Darstellung aus der Einführung, in der dritten Zeite das Bitmuster und als letztes die Exponentialschreibweise. Darstellungsübergreifend ist die Mantisse in braun markiert, der Exponent in Orange und das Vorzeichen in grün.

Als erstes platzieren wir den Kasten. Damit die Zahl möglichst gross wird, muss der Kasten nach links möglichst weit weg vom Komma stehen. Wir haben aber eine Einschränkung: Das Seil muss immer mit dem Komma verbunden bleiben.

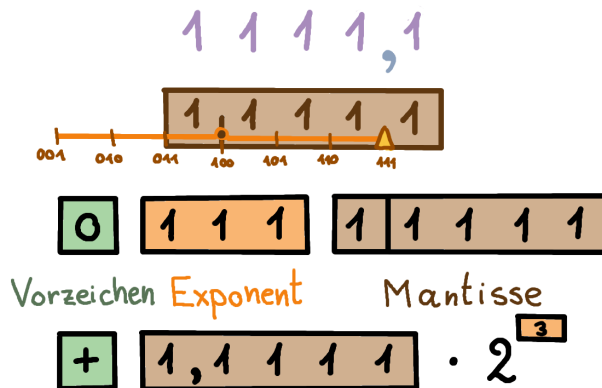


Der Exponent muss also möglichst gross sein.

Was ist mit der Mantisse? Sicher muss eine Eins an der ersten Stelle stehen.



Damit die Mantisse möglichst gross wird, muss sie aus lauter Einser bestehen.

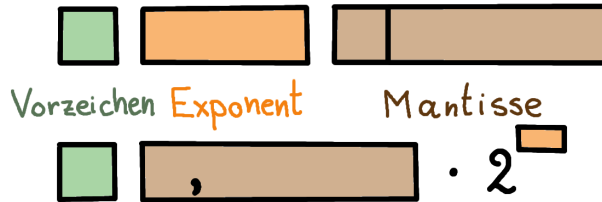
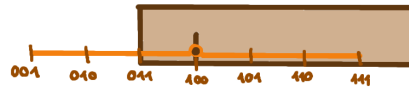


Die grösste darstellbare Zahl in diesem Fließkommazahlensystem ist also 15.5.

Beispiel 2.2. Wir konstruieren die kleinste positive Zahl im Fließkommazahlensystem mit Mantissenlänge 5 und Exponenten von -3 bis 3 .

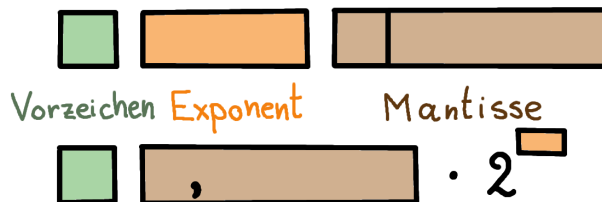
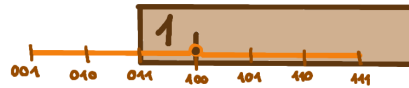
Als erstes platzieren wir den Kasten. Damit die Zahl möglichst klein wird, muss der Kasten nach rechts möglichst weit weg vom Komma stehen. Wir haben aber eine Einschränkung: Das Seil muss immer mit dem Komma verbunden bleiben.

? ? ? ? ? , ? ? ? ?

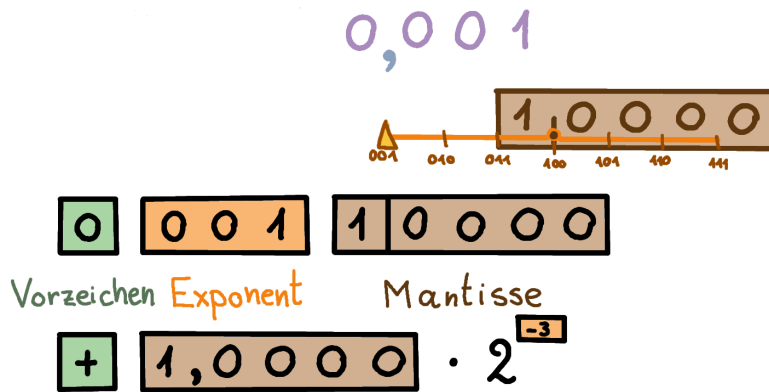


Der Exponent muss also möglichst klein sein.
Was ist mit der Mantisse? Sicher muss eine Eins an der ersten Stelle stehen.

? ? ? ? ? , ? ? 1 ?



Damit die Mantisse möglichst klein wird, müssen wir so viele Stellen wie möglich auf Null setzen.



Die kleinste darstellbare Zahl in diesem Fließkommazahlensystem ist also $1/8$.

Aufgabe 2.1. Betrachte das Fließkommazahlensystem mit Mantissenlänge 3 und Exponent von -1 bis 1 . Was sind die grösste und die kleinste darstellbare Zahlen in diesem Fließkommazahlensystem? Gib den Dezimalwert der Zahl an und stelle sie als Bitmuster und in der Exponentialschreibweise dar.

Für den Exponenten im Bitmuster verwende folgende Kodierung: 01 kodiert -1 , 10 kodiert 0 und 11 kodiert 1 .

Aufgabe 2.2. Betrachte das Fließkommazahlensystem mit Mantissenlänge 4 und Exponent von -1 bis 1 . Was sind die grösste und die kleinste darstellbare Zahlen in diesem Fließkommazahlensystem? Gib den Dezimalwert der Zahl an und stelle sie als Bitmuster und in der Exponentialschreibweise dar.

Für den Exponenten im Bitmuster verwende folgende Kodierung: 01 kodiert -1 , 10 kodiert 0 und 11 kodiert 1 .

Im Allgemeinen für einen Fließkommazahlensystem mit Mantissenlänge m und Exponenten zwischen e_{\min} und e_{\max} findet man die grösste und kleinste positive Zahlen wie folgt.

Für die grösste positive Zahl wählt man den grösstmöglichen Exponenten e_{\max} und die grösstmögliche Mantisse $1.111\dots111$. In der Exponentialschreibweise ist die grösste Zahl also

$$1.111111\dots111 \cdot 2^{e_{\max}}$$

und hat das Bitmuster 0 1111...111 (1)111111...111.

Für die kleinste positive Zahl wählt man den kleinsten möglichen Exponenten e_{\min} und die kleinste mögliche Mantisse. Beachte, dass die Mantisse immer mit einer Eins starten muss. Die kleinste mögliche Mantisse ist deswegen $1.0000\dots000$. In der Exponentialschreibweise ist die kleinste Zahl also

$$1.0000000\dots000 \cdot 2^{e_{\min}}$$

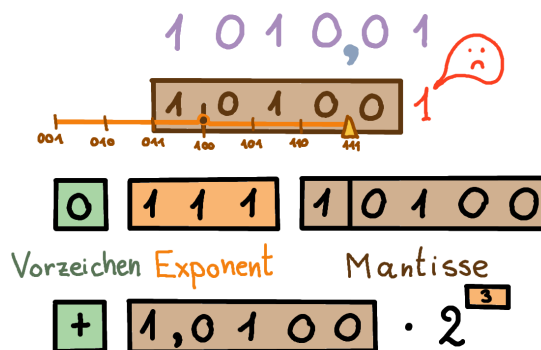
und hat das Bitmuster 0 0000...001 (1)00000000...000.

2.2 Darstellbare Zahlen

Wir wissen, dass es eine kleinste und eine grösste Zahl gibt. Dass man nicht alle reelle Zahlen zwischen diesen zwei Schranken darstellen kann, können wir uns denken. Die Frage ist nun, welche Zahlen sich darstellen lassen und wie sich der Abstand zwischen darstellbaren Zahlen verhält.

Hier und in den folgenden Kapiteln, falls nicht speziell vermerkt, werden wir mit Mantissenlänge 5 und Exponentenbereich von -3 bis 3 arbeiten.

Beispiel 2.3. Nehmen wir eine Zahl zwischen $1/8$ und 15.5 (die kleinste und grösste darstellbare Zahlen in diesem Fließkommazahlensystem), zum Beispiel 10.25 . Lässt sich diese Zahl darstellen?

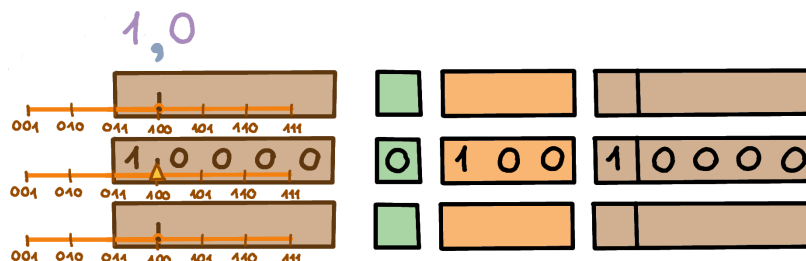


Diese Zahl lässt sich in gegebenem System nicht exakt darstellen. Für die letzte 1 gibt es in der Mantisse kein Platz. Deswegen wird 10.25 als 10.0 dargestellt.

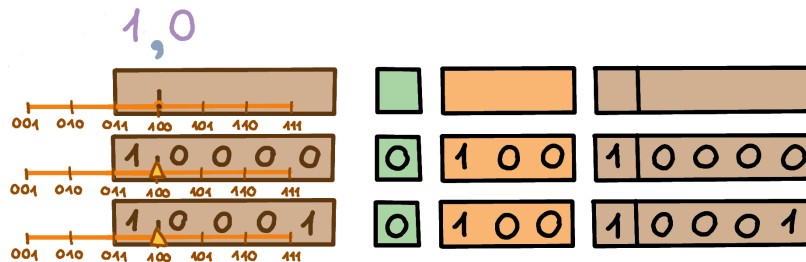
Wir haben gesehen, dass reelle Zahlen sich nur dann exakt darstellen lassen, wenn alle signifikante Stellen in der Mantisse Platz haben.

Beispiel 2.4. Betrachten wir die Zahl 1 . Was ist die nächstgrösste darstellbare Zahl? Und die nächstkleinste?

Im ersten Schritt werden wir die Zahl 1 im gegebenem Fließkommazahlensystem darstellen.

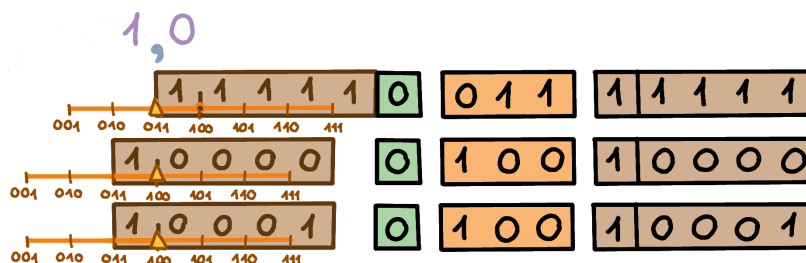


Die nächste darstellbare Zahl finden wir, indem wir in der Mantisse ganz rechts eine Eins addieren.



Die nächste darstellbare Zahl ist also $17/16$.

Die vorherige darstellbare Zahl finden wir, indem wir versuchen die Mantisse kleiner zu machen. Da die Mantisse von 1 die kleinste mögliche Mantisse ist, müssen wir den Exponenten um Eins zurücksetzen und die grösstmögliche Mantisse wählen.



Die vorherige Zahl ist also $31/32$.

Beachte, dass der Abstand zur nächsten und vorherigen darstellbaren Zahlen in diesem Fall nicht symmetrisch ist: die nächste Zahl ist $1/16$ entfernt, während die vorherige nur $1/32$.

Aufgabe 2.3. Finde die nächste und die vorherige darstellbare Zahlen von folgenden Zahlen. Schreibe die Werte in der Dezimaldarstellung auf und stelle alle Zahlen als Bitmuster und in der Exponentialschreibweise dar. Sind alle Nachbarn gleich entfernt?

- (a) 2
- (b) 3
- (c) 4

Teste dich selber

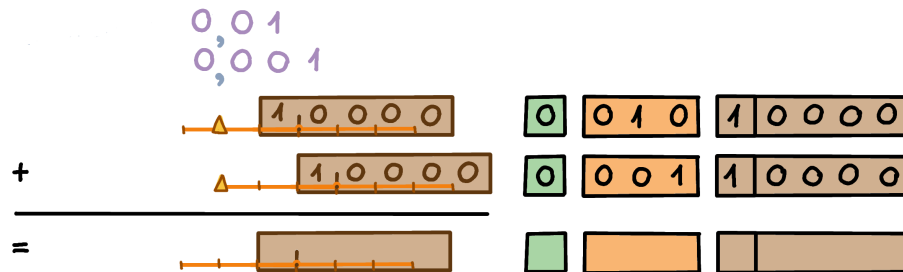
Aufgabe 2.4. Beantworte folgende Fragen:

- (a) *Kann man im Fliesskommazahlensystem alle reelle Zahlen darstellen? Wieso?*
- (b) *Gibt es eine grösste Zahl im Fliesskommazahlensystem? Falls nein, warum? Falls ja, wie findet man sie?*
- (c) *Gibt es eine kleinste Zahl im Fliesskommazahlensystem? Falls nein, warum? Falls ja, wie findet man sie?*
- (d) *Gib eine Zahl zwischen $1/2$ und 3.5 an, die im Fliesskommazahlensystem mit Mantissenlänge 3 und Exponenten von -1 bis 1 nicht darstellbar ist.*
- (e) *Sind alle Zahlen im Fliesskommazahlensystem gleichverteilt? Falls nicht, welche Zahlen stehen dichter beieinander, die kleineren oder die grösseren?*

3 Addition

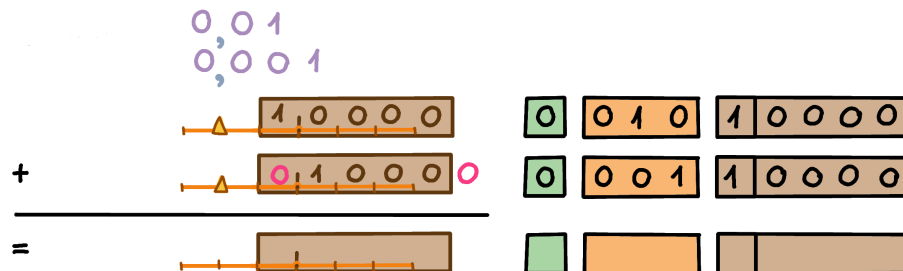
Im vorherigen Kapitel haben wir gesehen, welche Zahlen in einem Fließkommazahlensystem dargestellt werden können, das heisst welche Zahlen exakt in einem Computer gespeichert werden können. Computer werden aber nicht nur zum Speichern von Zahlen verwendet, sondern auch für Berechnungen. Auch bei Berechnungen verhalten sich Fließkommazahlen nicht ganz wie reelle Zahlen. In diesem Kapitel werden wir dies am Beispiel der Addition erfahren.

Beispiel 3.1. Wir möchten $1/4 + 1/8$ ausrechnen. Der erste Schritt ist beide Zahlen aufzuschreiben. Wie in den vorherigen Kapiteln, sind in violett die reelle Zahlen in Basis 2 aufgeschrieben und braune "Kasten" mit orangenem "Seil" verwendet, um Mantisse und Exponent zu veranschaulichen. Rechts wird das Bitmuster in der gewöhnlichen Form angegeben.

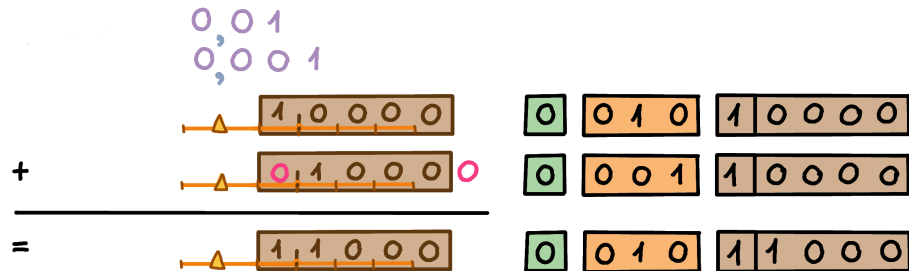


Damit wir die Bits der Mantisse stellenweise addieren können, wie wir das von den ganzen Zahlen kennen, müssen wir die zwei "Kasten" so verschieben, dass sie sich untereinander befinden. Da alles, was ausserhalb vom "Kasten" landet, verloren geht, werden wir den Kasten von der kleineren Zahl unter den Kasten von der grösseren Zahl schieben. So werden wir die Stellen mit dem niedrigsten Wert verlieren. In diesem Fall verlieren wir eine Null, der Wert der Zahl verändert sich also nicht.

Beachte, dass wenn der Kasten verschoben wird, verschiebt sich auch die Markierung bezüglich des "Seils", das heisst der Exponent verändert sich. Die Markierung am "Seil" bleibt immer unter dem Komma.

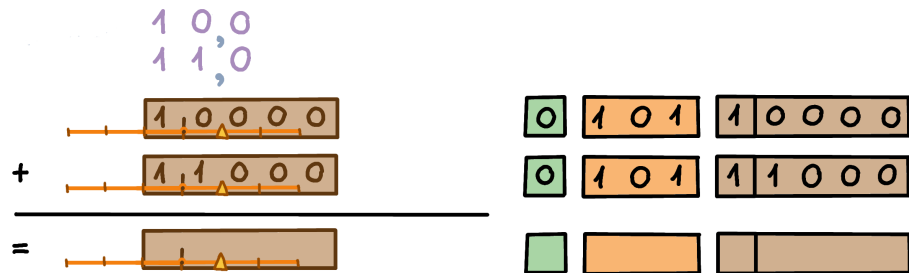


Wenn die Kasten untereinander sind, können wir die Bits in den Kasten wie gewöhnlich addieren, wie bei den Integers.

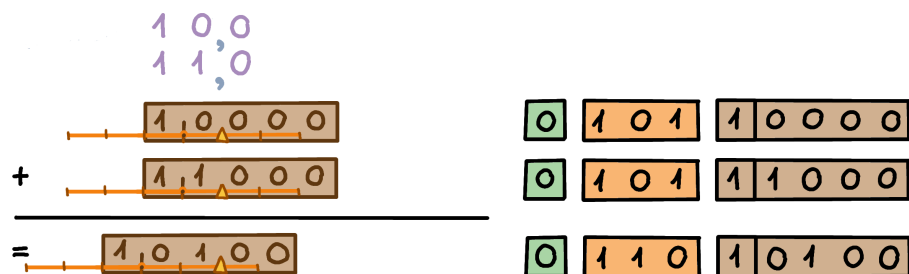


Wir haben ausgerechnet, dass $1/4 + 1/8 = 3/8$, in der Exponentialschreibweise $1.1000 \cdot 2^{-2}$.

Beispiel 3.2. Wir möchten $2 + 3$ ausrechnen. Im ersten Schritt schreiben wir beide Zahlen auf.



Die Kasten befinden sich schon untereinander. Wir müssen also nichts verschieben und können sofort losrechnen.



Beachte, dass der Kasten vom Ergebnis bezüglich den Kasten der Summanden verschoben ist, um die neue signifikante Stelle zu enthalten.

Wir haben ausgerechnet, dass $2 + 3 = 5$, in der Exponentialschreibweise $1.0100 \cdot 2^2$.

Aufgabe 3.1. *Rechne folgende Summen aus. Die Mantissenlänge beträgt 5 Bits, der Exponent geht von -3 bis 3 . Gebe bitte das Bitmuster und die Exponentialdarstellung des Resultats an.*

(a) $5/8 + 3/4$

(b) $10 + 2.25$

(c) $17/16 + 2$

Die Addition im Fliesskommazahlensystem ist wie gewöhnlich kommutativ, weil wir immer die kleinste Zahl so verschieben, dass ihr Kasten unter dem Kasten der grösseren Zahl steht und dann die Bits in beiden Kästen stellenweise zusammen addieren. In der folgenden Aufgabe werden wir prüfen, ob sie auch assoziativ ist.

Lernaufgabe. *Du wirst jetzt herausfinden, ob die Addition bei den Fliesskommazahlen assoziativ ist. Berechne dazu zwei Mal die gleiche Summe in einem Fliesskommazahlensystem mit Mantissenlänge 5 und Exponenten von -3 bis 3 : Das erste Mal als $1/8 + 2/8 + 3/8 + 4/8 + 5/8 + 6/8 + 7/8 + 8/8$ und das zweite Mal als $8/8 + 7/8 + 6/8 + 5/8 + 4/8 + 3/8 + 2/8 + 1/8$.*

Welchen Resultat erwartest du? Sind die zwei Summen gleich oder unterschiedlich? Kannst du daraus folgern, ob Addition assoziativ ist? Nimm dir Zeit und rechne die zwei Summen tatsächlich aus.

Die zwei Summen, die du ausgerechnet hast, liefern unterschiedliche Ergebnisse. Die erste liefert den exakten Wert 4.5, während bei der zweiten Summe kriegen wir im Fliesskommazahlensystem nur 4.25, und das obwohl der exakte Wert dargestellt werden kann. Das passiert, weil man bei den Fliesskommazahlen nur Zahlen der ähnlichen Grössenordnung exakt addieren kann. In der ersten Summe addieren wir die kleineren Summanden am Anfang, wenn die kumulative Summe noch nicht zu gross ist. In der zweiten Summe wächst die kumulative Summe sehr schnell, und irgendwann sind die Summanden zu klein bezüglich der kumulativen Summe, um einen Unterschied zu machen.

Daraus können wir folgern, dass die Addition bei den Fliesskommazahlen nicht assoziativ ist.

Aufgabe 3.2. *Betrachten wir die Summe $1/8 + 1/8 + 1/8 + \dots + 1/8$. Bei den reellen Zahlen können wir mit solchen Summen auf beliebig grossen Zahlen kommen. Bei den Fliesskommazahlen kann das nicht gehen, weil, wie wir im vorherigen Kapitel gesehen haben, es eine grösste Fliesskommazahl gibt. Aber können wir diese Zahl auch tatsächlich erreichen?*

In einem Fliesskommazahlensystem mit Mantissenlänge 5 und Exponentenbereich von -3 bis 3 , was ist die grösste Zahl, die wir erreichen können, wenn wir beliebig viele $1/8$ zusammen rechnen? Wie viele Summanden brauchen wir, um diese Zahl zu erreichen?

Teste dich selber

Aufgabe 3.3. Beantworte folgende Fragen:

- (a) Warum kann man im Allgemeinen zwei Mantissen nicht stellenweise zusammen addieren?
- (b) Gregory behauptet, dass der Kasten vom Ergebnis sich immer genau unter dem Kasten der grössten Zahl befindet. Hat er recht? Argumentiere.
- (c) Hannah behauptet, dass die Addition bei den Fließkommazahlen nicht kommutativ und nicht assoziativ ist. Hat sie recht? Argumentiere.

Aufgabe 3.4. Die Ameisenkönigin möchte ausrechnen, wie viele Ameisen braucht sie, um 10 Reiskörnchen zu transportieren. Sie weiss, dass eine Ameise allein $1/4$ Reiskorn transportiert. Die Ameisenkönigin hat dazu folgendes Programm geschrieben.

```
1 def nof_ameisen():
2     sum = 0.0
3     i = 0
4     while node != 10.0:
5         i += 1
6         sum += 0.25
7     return i
```

Listing 1: Programm von der Ameisenkönigin

Die Ameisencomputer arbeiten mit Fließkommazahlen mit Mantissenlänge 5 und Exponenten zwischen -3 und 3 . Kann die Ameisenkönigin mit diesem Programm die gewünschte Anzahl Ameisen herausfinden? Falls ja, wie viele Ameisen braucht sie, um 10 Reiskörnchen zu transportieren laut diesem Programm? Falls nein, was ist die maximale Summe, die das Programm erreichen kann?

4 Zusammenfassung

Wir haben gesehen, wie man im Computer reelle Zahlen durch **Fliesskommazahlen** approximieren kann. Da wir eine endliche Darstellung verwenden, gibt es eine endliche Anzahl Zahlen, die wir darstellen können. Es gibt also eine grösste und eine kleinste positive Zahl.

Die **grösste positive Zahl** kriegen wir, wenn wir die grösstmögliche Mantisse mit dem grösstmöglichen Exponenten kombinieren, d.h. die Mantisse besteht aus lauter Einser und der Exponent ist maximal.

Die **kleinste positive Zahl** kriegen wir, wenn wir die kleinstmögliche Mantisse mit dem kleinstmöglichen Exponenten kombinieren. Da wir in der Darstellung verlangen, dass das erste Bit der Mantisse eine Eins ist, besteht die kleinstmögliche Mantisse aus einer führenden Eins und vielen Nullen. Der Exponent ist minimal.

Wir haben gelernt, die nächste und die vorherige darstellbare Zahl zu bestimmen. So haben wir gesehen, dass darstellbare Zahlen nicht gleichverteilt auf der Zahlengerade auftreten, sondern dass der **Abstand** zwischen benachbarten darstellbaren Zahlen wächst, wenn die Zahlen grösser werden.

Mit den Fliesskommazahlen kann man auch rechnen. Wir haben insbesondere die **Addition** kennengelernt.

Wenn man zwei Fliesskommazahlen zusammen addieren möchte, muss man sie zuerst zum gleichen Exponenten bringen, dann kann man die neuen Mantissen wie ganze Zahlen addieren. Am Schluss muss man sicherstellen, dass der Exponent und die Mantisse gültig sind: der Exponent muss zwischen dem minimalen und dem maximalen Exponenten sein, die Mantisse muss eine führende Eins haben.

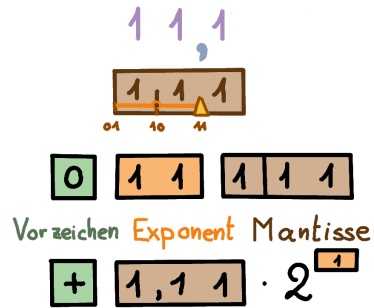
Wir haben gesehen, dass die Addition im Fliesskommazahlensystem Unterschiede zur Addition bei den reellen Zahlen aufweist. Erstens, nicht alle darstellbare Zahlen lassen sich exakt addieren. Zweitens, die Addition bei den Fliesskommazahlen ist **nicht assoziativ**. Es kann einen Unterschied machen, welche Teilsummen man zuerst berechnet.

Die Fliesskommazahlen sind eine mächtige Darstellung, die mit wenig Bits sehr unterschiedliche Zahlen speichern kann. Das hat aber auch seine Grenzen. Wir müssen in Kauf nehmen, dass die Resultate der Berechnungen nicht immer exakt sind.

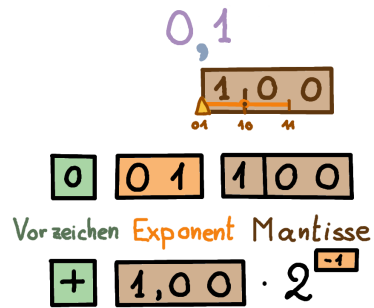
5 Beispiellösungen

5.1 Fließkommazahlen

Aufgabe 2.1 Die grösste Zahl in diesem System ist 3.5.

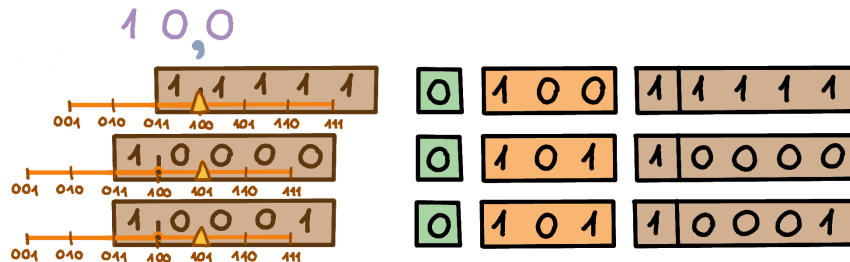


Die kleinste Zahl in diesem System ist 0.5.

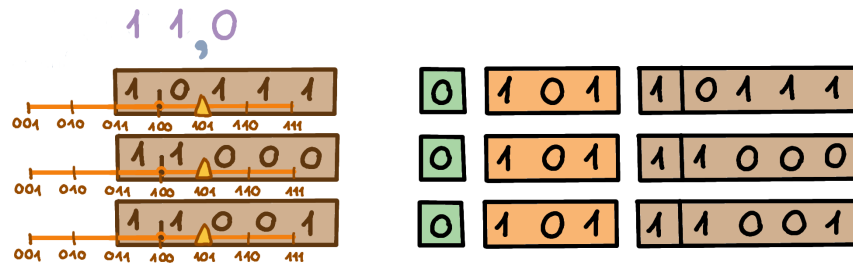


Aufgabe 2.3

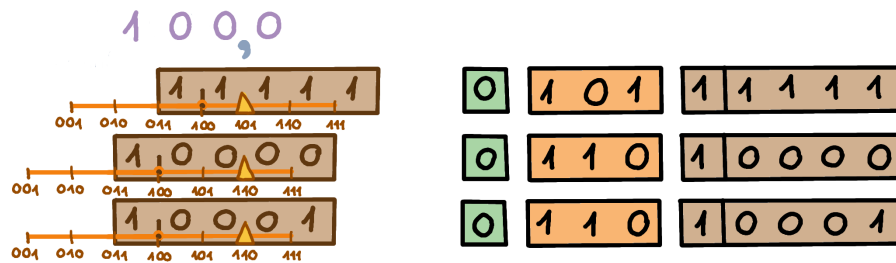
(a) Die Nachbarn von 2 sind $31/16$ und $17/8$.



(b) Die Nachbarn von 3 sind $23/8$ und $25/8$.



(c) Die Nachbarn von 4 sind $31/8$ und $17/4$.



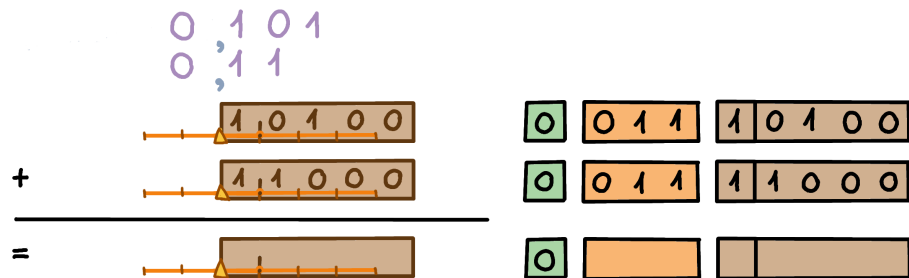
Aufgabe 2.4

- Nein, es gibt unendlich viele reelle Zahlen und endlich viele Fließkommazahlen.
- Ja, die grösste Zahl ist $1.1111 \dots 111 \cdot 2^{e_{max}}$.
- Ja, die kleinste Zahl ist $1.0000 \dots 000 \cdot 2^{e_{min}}$.
- Zum Beispiel, die Zahl 2.25 lässt sich in diesem System nicht exakt darstellen.
- Nein, die darstellbare Fließkommazahlen sind nicht gleichverteilt. Die kleineren stehen dichter beieinander, weil bei kleineren Zahlen die letzte Stelle der Mantisse weniger Wert ist.

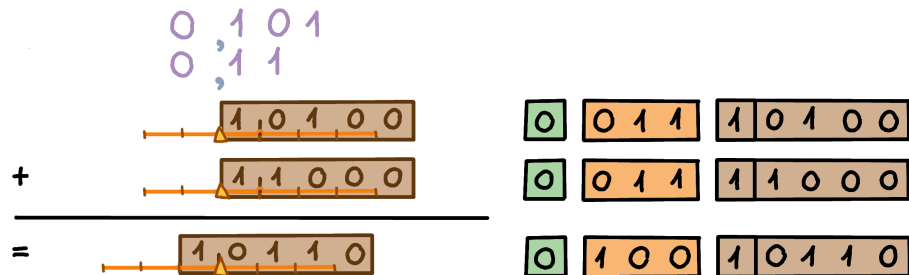
5.2 Addition

Aufgabe 3.1

- $5/8 + 3/4 = 11/8$, in der Exponentialschreibweise $1.0110 \cdot 2^0$
Im ersten Schritt schreiben wir die Zahlen auf.



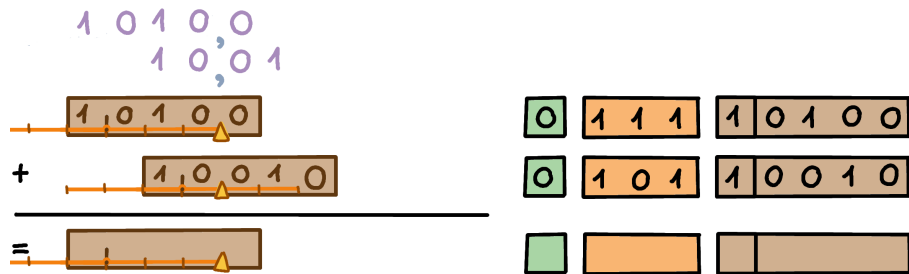
Da die zwei Kästen schon übereinander liegen, müssen wir sie nicht verschieben und können die Bits stellenweise zusammen addieren.



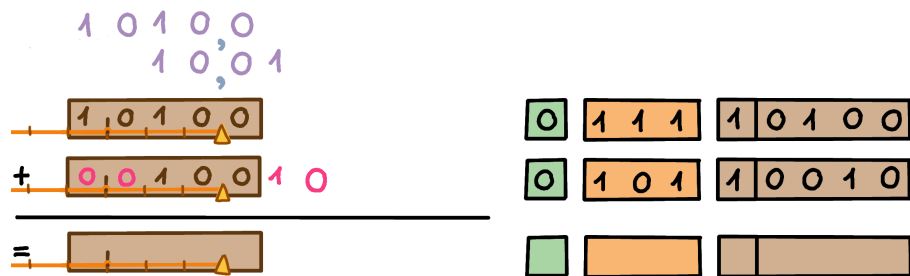
Der Kasten vom Ergebnis ist verschoben bezüglich den Kästen der Summanden.

- (b) $10 + 2.25 = 12$, in der Exponentialschreibweise $1.1000 \cdot 2^3$

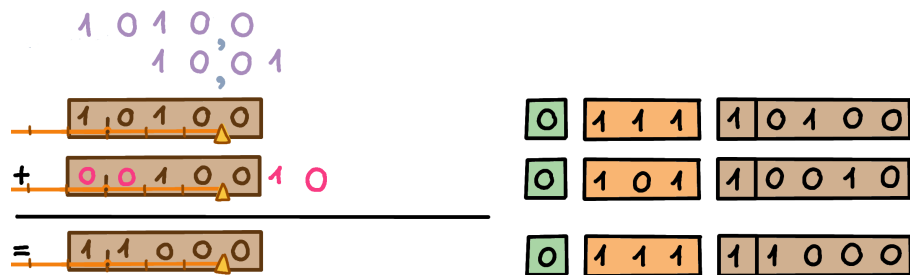
Im ersten Schritt schreiben wir die Zahlen auf.



Im zweiten Schritt schieben wir den Kasten von der kleinsten Zahl unter den Kasten der grössten Zahl. Dabei gehen zwei Stellen verloren, eine davon ist eine Eins.

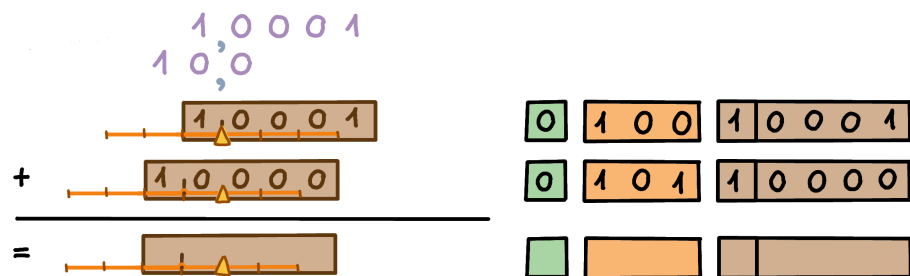


Nun können wir die Bits stellenweise zusammenrechnen.

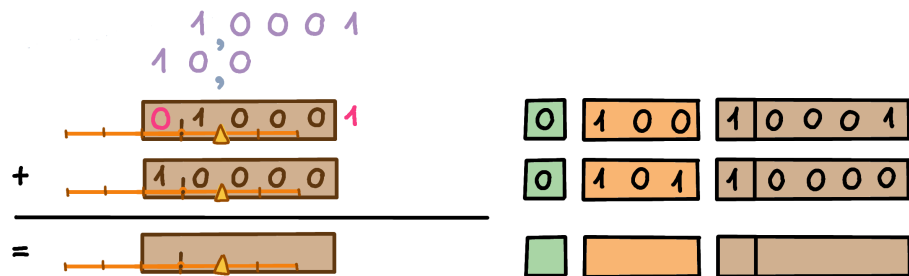


(c) $17/16 + 2 = 3$, in der Exponentialschreibweise $1.1000 \cdot 2^1$.

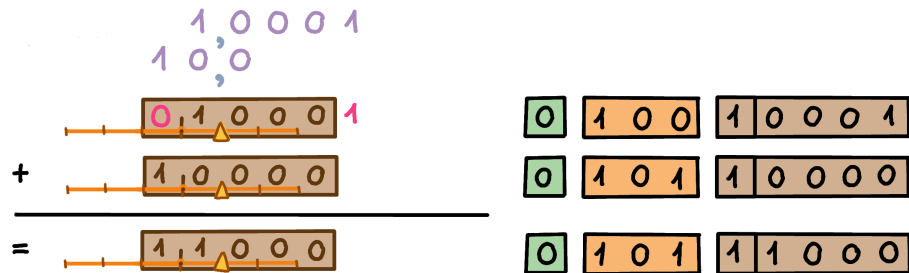
Im ersten Schritt schreiben wir die Zahlen auf.



Im zweiten Schritt schieben wir den Kasten von der kleinsten Zahl unter den Kasten der grössten Zahl. Dabei geht eine Stelle verloren.

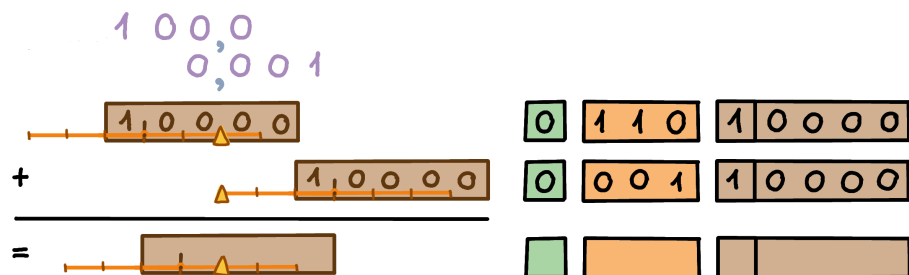


Nun können wir die Bits stellenweise zusammenrechnen.

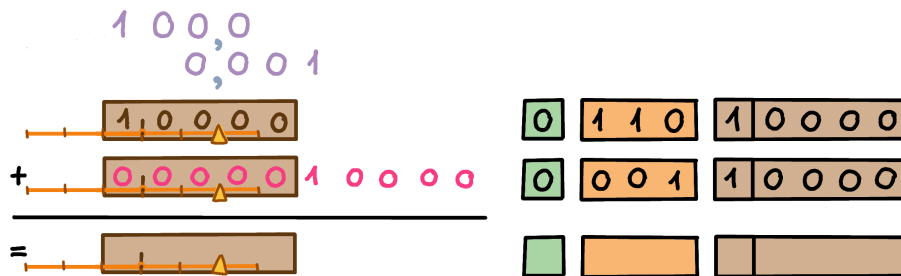


Aufgabe 3.2 Die maximale Zahl, die wir erreichen können, wenn wir $1/8 + 1/8 + \dots + 1/8$ zusammen rechnen, ist 4.0.

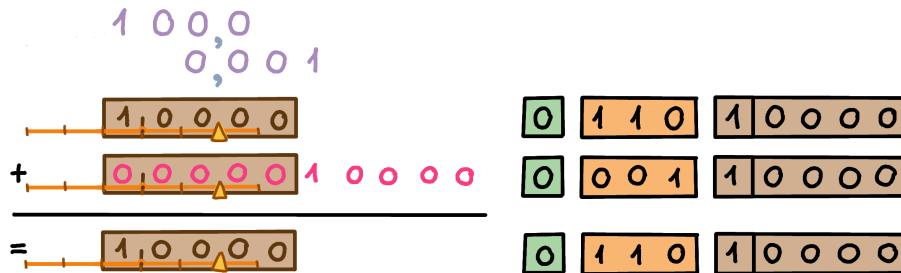
Zum einen, wenn wir die 4.0 erreicht haben, kommen wir nicht mehr weiter. Das sehen wir, wenn wir $4.0 + 1/8$ ausrechnen. Wie gewöhnlich schreiben wir zuerst die Summanden untereinander.



Wenn wir den Kasten von $1/8$ unter den Kasten von 4.0 verschieben, sehen wir, dass alle signifikanten Stellen von $1/8$ verloren gehen, auch die führende Eins.



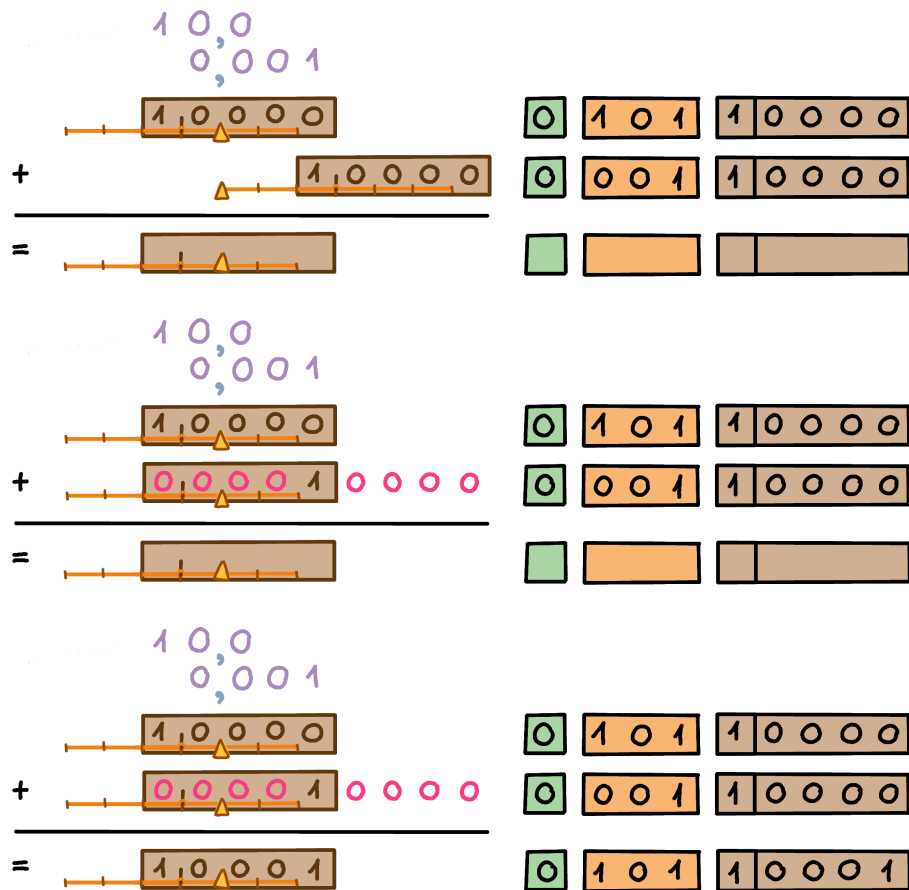
Deswegen, wenn wir $4.0 + 1/8$ ausrechnen, kriegen wir 4.0.



Egal wie viele $1/8$ rechnen wir zusammen, bleiben wir immer bei 4.0.

Jetzt bleibt uns zu zeigen, dass wir die 4.0 auch tatsächlich erreichen können. Das Problem bei der 4.0 ist, dass alle signifikanten Stellen von $1/8$ verloren gehen. Das passiert, weil der Unterschied zwischen dem Exponenten von 4.0 und dem Exponenten von $1/8$ die ganze Mantissenlänge beträgt. Das passiert bei einem kleineren Exponenten nicht. Zum Beispiel, wenn wir $2.0 + 1/8$ ausrechnen, sehen wir, dass das Ergebnis wie erwartet $17/8$ ist.

Um zu zeigen, dass das Problem erst bei $4.0 + 1/8$ auftritt, rechnen wir $2.0 + 1/8$. Das Ergebnis ist wie erwartet $17/8$.



Wir erreichen also die 4.0 nach 32 Summanden und kommen dann nicht mehr weiter.

Aufgabe 3.3

- Der Wert der Bits in der Mantisse hängt vom Exponenten ab. Zum Beispiel, dieselbe Mantisse 1.0000 mit unterschiedlichen Exponenten kann 4, 2, 1, $1/2$, $1/4$ und $1/8$ darstellen. Wir wollen nicht, dass $1 + 2$ das gleiche Ergebnis liefert die $1 + 1/4$. Wir wollen nur Bits mit dem gleichen Wert zusammen addieren. Deswegen müssen wir vor der Addition sicherstellen, dass die Kästen der beiden Summanden exakt untereinander stehen.
- Die Aussage von Gregory ist falsch. Der Kasten vom Ergebnis kann sich bewegen bezüglich des Kastens vom grössten Summanden. Dies passiert, zum Beispiel, wenn man $2.5 + 1.75$ ausrechnet.

- (c) Hannah hat teilweise recht. Die Addition bei den Fließkommazahlen ist kommutativ aber nicht assoziativ.

Wenn wir zwei Zahlen zusammen addieren und diese zwei Zahlen vertauschen, kriegen wir das gleiche Ergebnis auch bei Fließkommazahlen.

Wenn wir aber die Reihenfolge verändern, in welcher die Zahlen zusammengerechnet werden, können wir unterschiedlich Ergebnisse bekommen. Das passiert, weil wir nur dann den exakten Wert ausrechnen können, wenn die Größenordnung der Teilsummanden ähnlich ist.

Aufgabe 3.4 Nein, das Programm der Ameisenkönigin wird unendlich lange laufen und die Anzahl Ameisen, die es braucht, um 10 Reiskörnchen zu transportieren, nie ausgeben. Das Problem ist analog zu dem, was wir in Aufgabe 3.2 gesehen haben. Das Programm läuft wie erwartet bis wir die 8.0 erreichen. Wenn wir aber $1/4$ dazu rechnen, dann verlieren wir alle signifikanten Stellen von $1/4$ und die 8.0 bleibt unverändert.

