

# Binäre Zahlendarstellungen

Binäre Darstellung ("Bits" aus {0, 1})

$$b_n b_{n-1} \dots b_1 b_0$$

entspricht der Zahl  $b_n \cdot 2^n + \dots + b_1 \cdot 2 + b_0$

Beispiel: **101011** entspricht 43.

Niedrigstes Bit, Least Significant Bit (LSB)

Höchstes Bit, Most Significant Bit (MSB)

# 4. Zahlendarstellungen

Wertebereich der Typen int, float und double Gemischte Ausdrücke und Konversionen; Löcher im Wertebereich; Fließkommazahlensysteme; IEEE Standard; Grenzen der Fließkommaarithmetik; Fließkomma-Richtlinien;

# Binäre Zahlen: Zahlen der Computer?

Wahrheit: Computer rechnen mit Binärzahlen.



# Binäre Zahlen: Zahlen der Computer?

Klischee: Computer reden 0/1-Kauderwelsch.



## Rechentricks

- Abschätzung der Grössenordnung von Zweierpotenzen<sup>3</sup>:

$$\begin{aligned}2^{10} &= 1024 = 1\text{Ki} \approx 10^3. \\2^{20} &= 1\text{Mi} \approx 10^6, \\2^{30} &= 1\text{Gi} \approx 10^9, \\2^{32} &= 4 \cdot (1024)^3 = 4\text{Gi} \approx 4 \cdot 10^9. \\2^{64} &= 16\text{Ei} \approx 16 \cdot 10^{18}.\end{aligned}$$

<sup>3</sup>Dezimal vs. Binäre Einheiten: MB - Megabyte vs. MiB - Megabibyte (etc.)  
kilo (K, Ki) – mega (M, Mi) – giga (G, Gi) – tera(T, Ti) – peta(P, Pi) – exa (E, Ei)

126

## Wertebereich des Typs int

```
public class Main {
    public static void main(String[] args) {
        Out.print("Minimum int value is ");
        Out.println(Integer.MIN_VALUE);
        Out.print("Maximum int value is ");
        Out.println(Integer.MAX_VALUE);
    }
}
```

```
Minimum int value is -2147483648.
Maximum int value is 2147483647.
```

Woher kommen diese Zahlen?

127

## Wertebereich des Typs int

Repräsentation mit 32 Bits. Wertebereich umfasst die  $2^{32}$  ganzen Zahlen:

$$\{-2^{31}, -2^{31} + 1, \dots, -1, 0, 1, \dots, 2^{31} - 2, 2^{31} - 1\}$$

Woher kommt gerade diese Aufteilung?

128

## Rechnen mit Binärzahlen (4 Stellen)

Einfache Addition

2	0010
+3	+0011
-----	-----
5	0101

Einfache Subtraktion

5	0101
-3	-0011
-----	-----
2	0010

129

## Rechnen mit Binärzahlen (4 Stellen)

Addition mit Überlauf

$$\begin{array}{r} 7 \\ +9 \\ \hline 16 \end{array} \qquad \begin{array}{r} 0111 \\ +1001 \\ \hline (1)0000 \end{array}$$

Negative Zahlen?

$$\begin{array}{r} 5 \\ +(-5) \\ \hline 0 \end{array} \qquad \begin{array}{r} 0101 \\ +???? \\ \hline (1)0000 \end{array}$$

130

## Rechnen mit Binärzahlen (4 Stellen)

Einfacher: -1

$$\begin{array}{r} 1 \\ +(-1) \\ \hline 0 \end{array} \qquad \begin{array}{r} 0001 \\ 1111 \\ \hline (1)0000 \end{array}$$

Nutzen das aus:

$$\begin{array}{r} 3 \\ +? \\ \hline -1 \end{array} \qquad \begin{array}{r} 0011 \\ +???? \\ \hline 1111 \end{array}$$

131

## Rechnen mit Binärzahlen (4 Stellen)

Invertieren!

$$\begin{array}{r} 3 \\ +(-4) \\ \hline -1 \end{array} \qquad \begin{array}{r} 0011 \\ +1100 \\ \hline 1111 \hat{=} 2^B - 1 \end{array}$$

$$\begin{array}{r} a \\ +(-a-1) \\ \hline -1 \end{array} \qquad \begin{array}{r} a \\ \bar{a} \\ \hline 1111 \hat{=} 2^B - 1 \end{array}$$

132

## Rechnen mit Binärzahlen (4 Stellen)

- Negation: Inversion und Addition von 1

$$-a \hat{=} \bar{a} + 1$$

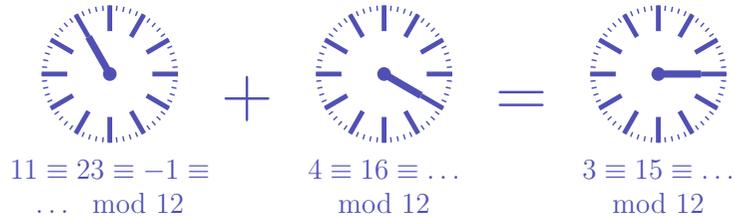
- Wrap-around Semantik (Rechnen modulo  $2^B$ )

$$-a \hat{=} 2^B - a$$

133

## Warum das funktioniert

Modulo-Arithmetik: Rechnen im Kreis<sup>4</sup>



<sup>4</sup>Die Arithmetik funktioniert auch mit Dezimalzahlen (und auch für die Multiplikation)

134

## Negative Zahlen (3 Stellen)

	$a$	$-a$	
0	000	000	0
1	001	111	-1
2	010	110	-2
3	011	101	-3
4	100	100	-4
5	101		
6	110		
7	111		

Das höchste Bit entscheidet über das Vorzeichen.

135

## Zweierkomplement

- Negation durch bitweise Negation und Addition von 1.

$$-2 = -[0010] = [1101] + [0001] = [1110]$$

- Arithmetik der Addition und Subtraktion *identisch* zur vorzeichenlosen Arithmetik.

$$3 - 2 = 3 + (-2) = [0011] + [1110] = [0001]$$

- Intuitive „Wrap-Around“ Konversion negativer Zahlen.

$$-n \rightarrow 2^B - n$$

- Wertebereich:  $-2^{B-1} \dots 2^{B-1} - 1$

136

## Überlauf und Unterlauf

- Arithmetische Operationen (+, -, \*) können aus dem Wertebereich herausführen.
- Ergebnisse können inkorrekt sein.

$$\text{power8: } 15^8 = -1732076671$$

$$\text{power20: } 3^{20} = -808182895$$

- Es gibt *keine Fehlermeldung!*

137

## „Richtig Rechnen“

```
public class Main {  
  
    public static void main(String[] args) {  
        Out.print("Celsius: ");  
        int celsius = In.readInt();  
        int fahrenheit = 9 * celsius / 5 + 32;  
        Out.print(celsius + " degrees Celsius are ");  
        Out.println(fahrenheit + " degrees Fahrenheit");  
    }  
}  
  
public class Main {  
  
    public static void main(String[] args) {  
        Out.print("Celsius: ");  
        float celsius = In.readInt();
```

138

## Typen float und double

- sind die fundamentalen Typen für Fließkommazahlen
- approximieren den Körper der reellen Zahlen ( $\mathbb{R}, +, \times$ ) in der Mathematik
- haben grossen Wertebereich, ausreichend für viele Anwendungen (double hat mehr Stellen als float)
- sind auf vielen Rechnern sehr schnell

139

## Fixkommazahlen

- feste Anzahl Vorkommastellen (z.B. 7)
- feste Anzahl Nachkommastellen (z.B. 3)

0.0824 = 0000000.082 ← dritte Stelle abgeschnitten

### Nachteile

- Wertebereich wird *noch* kleiner als bei ganzen Zahlen.
- Repräsentierbarkeit hängt von der Stelle des Kommas ab.

140

## Fließkommazahlen

- feste Anzahl signifikante Stellen (z.B. 10)
- plus Position des Kommas

$$82.4 = 824 \cdot 10^{-1}$$

$$0.0824 = 824 \cdot 10^{-4}$$

- Zahl ist  $\text{Signifikand} \times 10^{\text{Exponent}}$

141

## Wertebereich

Ganzzahlige Typen:

- Über- und Unterlauf häufig, aber ...
- Wertebereich ist zusammenhängend (keine „Löcher“):  $\mathbb{Z}$  ist „diskret“.

Fließkommamatypen:

- Über- und Unterlauf selten, aber ...
- es gibt Löcher:  $\mathbb{R}$  ist „kontinuierlich“.

142

## Löcher im Wertebereich

```
public class Main {
    public static void main(String[] args) {
        Out.print("First number =? ");
        float n1 = In.readFloat();

        Out.print("Second number =? ");
        float n2 = In.readFloat();

        Out.print("Their difference =? ");
        float d = In.readFloat();

        Out.print("computed difference - input difference = ");
        Out.println(n1-n2-d);
    }
}
```

Eingabe 1.1

Eingabe 1.0

Eingabe 0.1

Ausgabe 2.2351742E-8

Ja was ist denn hier los?

143

## Fließkommazahlensysteme

Ein Fließkommazahlensystem ist durch vier natürliche Zahlen definiert:

- $\beta \geq 2$ , die Basis,
- $p \geq 1$ , die Präzision (Stellenzahl),
- $e_{\min}$ , der kleinste Exponent,
- $e_{\max}$ , der grösste Exponent.

Bezeichnung:

$$F(\beta, p, e_{\min}, e_{\max})$$

144

## Fließkommazahlensysteme

$F(\beta, p, e_{\min}, e_{\max})$  enthält die Zahlen

$$\pm \sum_{i=0}^{p-1} d_i \beta^{-i} \cdot \beta^e,$$

$$d_i \in \{0, \dots, \beta - 1\}, \quad e \in \{e_{\min}, \dots, e_{\max}\}.$$

In Basis- $\beta$ -Darstellung:

$$\pm d_0 \bullet d_1 \dots d_{p-1} \times \beta^e,$$

145

## Fließkommazahlensysteme

Beispiel

■  $\beta = 10$

Darstellungen der Dezimalzahl 0.1

$$1.0 \cdot 10^{-1}, \quad 0.1 \cdot 10^0, \quad 0.01 \cdot 10^1, \quad \dots$$

146

## Normalisierte Darstellung

Normalisierte Zahl:

$$\pm d_0 \bullet d_1 \dots d_{p-1} \times \beta^e, \quad d_0 \neq 0$$

### Bemerkung 1

Die normalisierte Darstellung ist eindeutig und deshalb zu bevorzugen.

### Bemerkung 2

Die Zahl 0 (und alle Zahlen kleiner als  $\beta^{e_{\min}}$ ) haben keine normalisierte Darstellung (werden wir später beheben)!

147

## Menge der normalisierten Zahlen

$$F^*(\beta, p, e_{\min}, e_{\max})$$

148

## Normalisierte Darstellung

Beispiel  $F^*(2, 3, -2, 2)$

(nur positive Zahlen)

$d_0 \bullet d_1 d_2$	$e = -2$	$e = -1$	$e = 0$	$e = 1$	$e = 2$
1.00 <sub>2</sub>	0.25	0.5	1	2	4
1.01 <sub>2</sub>	0.3125	0.625	1.25	2.5	5
1.10 <sub>2</sub>	0.375	0.75	1.5	3	6
1.11 <sub>2</sub>	0.4375	0.875	1.75	3.5	7



149

## Binäre und dezimale Systeme

- Intern rechnet der Computer mit  $\beta = 2$  (binäres System)
- Literale und Eingaben haben  $\beta = 10$  (dezimales System)
- Eingaben müssen umgerechnet werden!

## Umrechnung dezimal $\rightarrow$ binär

Angenommen,  $0 < x < 2$ .

- Also:  $x' = b_{-1}b_{-2}b_{-3}b_{-4}\dots = 2 \cdot (x - b_0)$
- Schritt 1 (für  $x$ ): Berechnen von  $b_0$ :

$$b_0 = \begin{cases} 1, & \text{falls } x \geq 1 \\ 0, & \text{sonst} \end{cases}$$

- Schritt 2 (für  $x$ ): Berechnen von  $b_{-1}, b_{-2}, \dots$ :  
Gehe zu Schritt 1 (für  $x' = 2 \cdot (x - b_0)$ )

150

152

## Binärdarstellung von 1.1

$x$	$b_i$	$x - b_i$	$2(x - b_i)$
1.1	$b_0 = \mathbf{1}$	0.1	0.2
0.2	$b_{-1} = \mathbf{0}$	0.2	0.4
0.4	$b_{-2} = \mathbf{0}$	0.4	0.8
0.8	$b_{-3} = \mathbf{0}$	0.8	1.6
1.6	$b_{-4} = \mathbf{1}$	0.6	1.2
1.2	$b_{-5} = \mathbf{1}$	0.2	0.4

$\Rightarrow 1.\overline{00011}$ , periodisch, *nicht* endlich

153

## Binärdarstellungen von 1.1 und 0.1

- sind nicht endlich, also gibt es Fehler bei der Konversion in ein (endliches) binäres Fließkommazahlensystem.
- 1.1f und 0.1f sind nicht 1.1 und 0.1, sondern geringfügig fehlerhafte Approximationen dieser Zahlen.

$$\begin{aligned} 1.1 &= \underline{1.1000000000000000}888178\dots \\ 1.1f &= \underline{1.1000000}238418\dots \end{aligned}$$

154

## Rechnen mit Fließkommazahlen

Beispiel ( $\beta = 2, p = 4$ ):

$$\begin{array}{r} 1.111 \cdot 2^{-2} \\ + 1.011 \cdot 2^{-1} \\ \hline = 1.001 \cdot 2^0 \end{array}$$

1. Exponenten anpassen durch Denormalisieren einer Zahl 2. Binäre Addition der Signifikanden 3. Renormalisierung 4. Runden auf  $p$  signifikante Stellen, falls nötig

155

## Der IEEE Standard 754

- legt Fließkommazahlensysteme und deren Rundungsverhalten fest
- wird fast überall benutzt
- Single precision (`float`) Zahlen:

$$F^*(2, 24, -126, 127) \quad \text{plus } 0, \infty, \dots$$

- Double precision (`double`) Zahlen:

$$F^*(2, 53, -1022, 1023) \quad \text{plus } 0, \infty, \dots$$

- Alle arithmetischen Operationen runden das *exakte* Ergebnis auf die nächste darstellbare Zahl

156

## 32-bit Darstellung einer Fließkommazahl



± Exponent

Mantisse

±  $2^{-126}, \dots, 2^{127}$   
0,  $\infty, \dots$

1.00000000000000000000000  
1.11111111111111111111111

## Der IEEE Standard 754

Warum

$$F^*(2, 24, -126, 127)?$$

- 1 Bit für das Vorzeichen
- 23 Bit für den Signifikanden (führendes Bit ist 1 und wird nicht gespeichert)
- 8 Bit für den Exponenten (256 mögliche Werte)(254 mögliche Exponenten, 2 Spezialwerte: 0,  $\infty, \dots$ )

⇒ insgesamt 32 Bit.

157

158

## Der IEEE Standard 754

Warum

$$F^*(2, 53, -1022, 1023)?$$

- 1 Bit für das Vorzeichen
- 52 Bit für den Signifikanden (führendes Bit ist 1 und wird nicht gespeichert)
- 11 Bit für den Exponenten (2046 mögliche Exponenten, 2 Spezialwerte: 0,  $\infty$ , ...)

⇒ insgesamt 64 Bit.

159

## Fliesskomma-Richtlinien

### Regel 1

#### Regel 1

Teste keine gerundeten Fliesskommazahlen auf Gleichheit!

```
for (float i = 0.1; i != 1.0; i += 0.1){  
    Out.println(i);  
}
```

Endlosschleife, weil i niemals exakt 1 ist!

160

## Fliesskomma-Richtlinien

### Regel 2

#### Regel 2

Addiere keine zwei Zahlen sehr unterschiedlicher Grösse!

$$\begin{aligned} & 1.000 \cdot 2^5 \\ & + 1.000 \cdot 2^0 \\ & = 1.00001 \cdot 2^5 \\ & \text{"=" } 1.000 \cdot 2^5 \text{ (Rundung auf 4 Stellen)} \end{aligned}$$

Addition von 1 hat keinen Effekt!

161

## Fliesskomma-Richtlinien

### Regel 3

#### Regel 3

Subtrahiere keine zwei Zahlen sehr ähnlicher Grösse!

Auslöschungsproblematik (ohne weitere Erklärung).

162

## 5. Wahrheitswerte

Boolesche Funktionen; der Typ `boolean`; logische und relationale Operatoren; Kurzschlussauswertung

### Wo wollen wir hin?

```
int a = In.readInt();
if (a % 2 == 0){
    Out.println("even");
} else {
    Out.println("odd");
}
```

Verhalten hängt ab vom Wert eines **Booleschen Ausdrucks**

163

164

### Boolesche Werte in der Mathematik

Boolesche Ausdrücke können zwei mögliche Werte annehmen:

$F$  oder  $T$

- $F$  entspricht „*falsch*“
- $T$  entspricht „*wahr*“

### Der Typ `boolean` in Java

- Repräsentiert *Wahrheitswerte*
- Literale `false` und `true`
- Wertebereich `{false, true}`

```
boolean b = true; // Variable mit Wert true (wahr)
```

165

166

## Relationale Operatoren

- $a < b$  (kleiner als)
- $a >= b$  (grösser gleich)
- $a == b$  (gleich)
- $a != b$  (ungleich)

Zahlentyp  $\times$  Zahlentyp  $\rightarrow$  boolean

167

## Boolesche Funktionen in der Mathematik

- Boolesche Funktion

$$f : \{F, T\}^2 \rightarrow \{F, T\}$$

- $F$  entspricht „falsch“.
- $T$  entspricht „wahr“.

168

## AND( $x, y$ )

$$x \wedge y$$

- “Logisches Und”

$$f : \{F, T\}^2 \rightarrow \{F, T\}$$

- $F$  entspricht „falsch“.
- $T$  entspricht „wahr“.

$x$	$y$	AND( $x, y$ )
$F$	$F$	$F$
$F$	$T$	$F$
$T$	$F$	$F$
$T$	$T$	$T$

## Logischer Operator &&

$a \ \&\& \ b$  (logisches Und)

boolean  $\times$  boolean  $\rightarrow$  boolean

```
int n = -1;
int p = 3;
boolean b = (n < 0) && (0 < p); // b = true (wahr)
```

169

170

## OR( $x, y$ )

$x \vee y$

- “Logisches Oder”

$$f : \{F, T\}^2 \rightarrow \{F, T\}$$

- $F$  entspricht „falsch”.
- $T$  entspricht „wahr”.

$x$	$y$	OR( $x, y$ )
$F$	$F$	$F$
$F$	$T$	$T$
$T$	$F$	$T$
$T$	$T$	$T$

## Logischer Operator ||

$a \ || \ b$  (logisches Oder)

`boolean × boolean → boolean`

```
int n = 1;
int p = 0;
boolean b = (n < 0) || (0 < p); // b = false (falsch)
```

171

172

## NOT( $x$ )

$\neg x$

- “Logisches Nicht”

$$f : \{F, T\} \rightarrow \{F, T\}$$

- $F$  entspricht „falsch”.
- $T$  entspricht „wahr”.

$x$	NOT( $x$ )
$F$	$T$
$T$	$F$

## Logischer Operator !

$!b$  (logisches Nicht)

`boolean → boolean`

```
int n = 1;
boolean b = !(n < 0); // b = true (wahr)
```

173

174

## Präzedenzen

$$\begin{array}{c} !b \ \&\& \ a \\ \Updownarrow \\ (!b) \ \&\& \ a \\ \\ a \ \&\& \ b \ || \ c \ \&\& \ d \\ \Updownarrow \\ (a \ \&\& \ b) \ || \ (c \ \&\& \ d) \\ \\ a \ || \ b \ \ \ \ \ \&\& \ \ \ \ \ \ c \ || \ d \\ \Updownarrow \\ a \ || \ (b \ \&\& \ c) \ || \ d \end{array}$$

175

## Präzedenzen

*Der unäre logische* Operator !  
bindet stärker als  
*binäre arithmetische* Operatoren. Diese  
binden stärker als  
*relationale* Operatoren,  
und diese binden stärker als  
*binäre logische* Operatoren.

$$\begin{array}{l} 7 + x < y \ \&\& \ y \ != \ 3 * z \ || \ ! \ b \\ 7 + x < y \ \&\& \ y \ != \ 3 * z \ || \ (!b) \end{array}$$

176

## DeMorgansche Regeln

- $!(a \ \&\& \ b) == (!a \ || \ !b)$
- $!(a \ || \ b) == (!a \ \&\& \ !b)$

! (reich *und* schön) == (arm *oder* hässlich)

177

## Anwendung: Entweder ... Oder (XOR)

$(x \    \ y)$	$\ \&\& \ !(x \ \&\& \ y)$	x oder y, und nicht beide
$(x \    \ y)$	$\ \&\& \ (!x \    \ !y)$	x oder y, und eines nicht
$!(x \ \&\& \ !y)$	$\ \&\& \ !(x \ \&\& \ y)$	nicht keines, und nicht beide
$!(x \ \&\& \ !y \    \ x \ \&\& \ y)$		nicht: keines oder beide

178

## Kurzschlussauswertung

- Logische Operatoren `&&` und `||` werten den *linken Operanden* *zuerst* aus.
- Falls das Ergebnis dann schon feststeht, wird der rechte Operand *nicht mehr* ausgewertet.

```
x != 0 && z / x > y
```

⇒ Keine Division durch 0