



## Computing Tricks

- Estimate the orders of magnitude of powers of two.<sup>3</sup>:

$$2^{10} = 1024 = 1\text{Ki} \approx 10^3.$$

$$2^{20} = 1\text{Mi} \approx 10^6,$$

$$2^{30} = 1\text{Gi} \approx 10^9,$$

$$2^{32} = 4 \cdot (1024)^3 = 4\text{Gi} \approx 4 \cdot 10^9.$$

$$2^{64} = 16\text{Ei} \approx 16 \cdot 10^{18}.$$

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

## Domain of Type 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.

Where do these numbers come from?

## Domain of the Type int

Representation with 32 bits. Domain comprises the  $2^{32}$  integers:

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

Where does this partitioning come from?

## Computing with Binary Numbers (4 digits)

Simple Addition

2	0010
+3	+0011
5	0101

Simple Subtraction

5	0101
-3	-0011
2	0010

## Computing with Binary Numbers (4 digits)

Addition with Overflow

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

Negative Numbers?

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

## Computing with Binary Numbers (4 digits)

Simpler -1

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

Utilize this:

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

130

131

## Computing with Binary Numbers (4 digits)

Invert!

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

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

## Computing with Binary Numbers (4 digits)

■ Negation: inversion and addition of 1

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

■ Wrap around semantics (calculating modulo  $2^B$ )

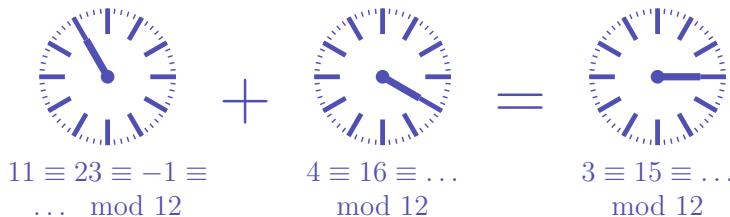
$$-a \hat{=} 2^B - a$$

132

133

## Why this works

Modulo arithmetics: Compute on a circle<sup>4</sup>



<sup>4</sup>The arithmetics also work with decimal numbers (and for multiplication).

## Negative Numbers (3 Digits)

	$a$	$-a$	
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		

The most significant bit decides about the sign.

134

135

## Two's Complement

### ■ Negation by bitwise negation and addition of 1

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

- Arithmetics of addition and subtraction *identical* to unsigned arithmetics

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

- Intuitive “wrap-around” conversion of negative numbers.

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

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

## Over- and Underflow

- Arithmetic operations ( $+, -, *$ ) can lead to numbers outside the valid domain.

- Results can be incorrect!

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

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

- There is *no error message!*

136

137

## “Proper Calculation”

```
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();  
        }  
    }  
}
```

## Types `float` and `double`

- are the fundamental types for floating point numbers
- approximate the field of real numbers ( $\mathbb{R}, +, \times$ ) from mathematics
- have a great domain, sufficient for many applications (`double` provides more places than `float`)
- are fast on many computers

138

139

## Fixpoint numbers

- fixed number of integer places (e.g. 7)
- fixed number of decimal places (e.g. 3)

$$0.0824 = 0000000.082 \leftarrow \text{third place truncated}$$

### Nachteile

- Domain is getting even smaller than for integers.
- If a number can be represented depends on the position of the comma.

## Floating point Numbers

- fixed number of significant places (e.g. 10)
- plus position of the comma

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

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

- Zahl ist  $Mantissa \times 10^{Exponent}$

140

141

## Domain

Integer Types:

- Over- and Underflow relatively frequent, but ...
- the domain is contiguous (no “holes”):  $\mathbb{Z}$  is “discrete”.

Floating point types:

- Overflow and Underflow seldom, but ...
- there are holes:  $\mathbb{R}$  is “continuous”.

## Holes in the Domain

```
public class Main {  
    public static void main(String[] args) {  
        Out.print("First number =? ");      input 1.1  
        float n1 = In.readFloat();  
  
        Out.print("Second number =? ");     input 1.0  
        float n2 = In.readFloat();  
  
        Out.print("Their difference =? ");  input 0.1  
        float d = In.readFloat();  
  
        Out.print("computed difference - input difference = ");  
        Out.println(n1-n2-d);  
    }  
}
```

What is going on here?

142

143

## Floating Point Number Systems

A Floating Point Number System is defined by the four natural numbers:

- $\beta \geq 2$ , the Basis,
- $p \geq 1$ , the precision (number of places),
- $e_{\min}$ , the smallest possible exponent,
- $e_{\max}$ , the largest possible exponent.

Notation:

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

## Floating Point Number Systems

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

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

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

represented with Basis  $\beta$ :

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

144

145

## Floating Point Number Systems

Example

■  $\beta = 10$

Representations of the decimal number 0.1

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

## Normalized Representation

Normalized Number:

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

### Bemerkung 1

The normalized representation is unique and therefore preferred.

### Remark 2

The number 0 (and all numbers smaller than  $\beta^{e_{\min}}$ ) have no normalized representation (we will deal with this later)!

146

147

## Set of Normalized Numbers

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

## Normalized Representation

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

(only positive numbers)

$d_0.d_1d_2$	$e = -2$	$e = -1$	$e = 0$	$e = 1$	$e = 2$
$1.00_2$	0.25	0.5	1	2	4
$1.01_2$	0.3125	0.625	1.25	2.5	5
$1.10_2$	0.375	0.75	1.5	3	6
$1.11_2$	0.4375	0.875	1.75	3.5	7



148

149

## Binary and Decimal Systems

- Internally the computer computes with  $\beta = 2$  (binary system)
- Literals and inputs have  $\beta = 10$  (decimal system)
- Inputs have to be converted!

## Binary representation of 1.1

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

$\Rightarrow 1.0\overline{0011}$ , periodic, not finite

## Conversion Decimal $\rightarrow$ Binary

Angenommen,  $0 < x < 2$ .

- Hence:  $x' = b_{-1} \bullet b_{-2} b_{-3} b_{-4} \dots = 2 \cdot (x - b_0)$
- Step 1 (for  $x$ ): Compute  $b_0$ :

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

- Step 2 (for  $x$ ): Compute  $b_{-1}, b_{-2}, \dots$ :  
Go to step 1 (for  $x' = 2 \cdot (x - b_0)$ )

150

152

## Binary Number Representations of 1.1 and 0.1

- are not finite, there are errors when converting into a (finite) binary floating point system.
- 1.1f and 0.1f do not equal 1.1 and 0.1, but slightly inaccurate approximation of these numbers.

$$1.1 = 1.1000000000000000888178\dots$$

$$1.1f = 1.1000000238418\dots$$

153

154

## Computing with Floating Point Numbers

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

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

1. adjust exponents by denormalizing of one number
2. binary addition of the mantissa
3. renormalize
4. round to  $p$  significant places, if necessary

## The IEEE Standard 754

- defines floating point number systems and their rounding behavior
- is used nearly everywhere
- Single precision (`float`) numbers:  
 $F^*(2, 24, -126, 127)$  plus  $0, \infty, \dots$
- Double precision (`double`) numbers:  
 $F^*(2, 53, -1022, 1023)$  plus  $0, \infty, \dots$
- All arithmetic operations round the exact result to the next representable number

155

156

## 32-bit Representation of a Floating Point Number



± Exponent

Mantisse

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

1.00000000000000000000000000  
1.11111111111111111111111111

## The IEEE Standard 754

Why

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

- 1 sign bit
- 23 bit for the mantissa (leading bit is 1 and is not stored)
- 8 bit for the exponent (256 possible values) (254 possible exponents, 2 special values:  $0, \infty, \dots$ )

⇒ 32 bit overall.

157

158

Why

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

- 1 sign bit
- 52 bit for the mantissa (leading bit is 1 and is not stored)
- 11 bit for the exponent (2046 possible exponents, 2 special values: 0,  $\infty, \dots$ )

$\Rightarrow$  64 bit overall.

## Rule 1

Do not test rounded floating point numbers for equality.

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

endless loop because i never becomes exactly 1

159

160

## Floating Point Rules

## Rule 2

Do not add two numbers providing very different orders of magnitude!

$$\begin{aligned}
 & 1.000 \cdot 2^5 \\
 & + 1.000 \cdot 2^0 \\
 & = 1.00001 \cdot 2^5 \\
 & \text{“=” } 1.000 \cdot 2^5 \text{ (Rounding on 4 places)}
 \end{aligned}$$

Addition of 1 does not provide any effect!

## Rule 2

## Floating Point Guidelines

## Rule 3

## Rule 4

Do not subtract two numbers with a very similar value.

Cancellation problems (without further explanations).

161

162

## Our Goal

# 5. Logical Values

Boolean Functions; the Type `boolean`; logical and relational operators; shortcut evaluation

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

Behavior depends on the value of a **Boolean expression**

163

164

## Boolean Values in Mathematics

Boolean expressions can take on one of two values:

*F* or *T*

- *F* corresponds to “*wrong*”
- *T* corresponds to “*true*”

## The Type `boolean` in Java

- represents *logical values*
- Literals `false` and `true`
- Domain {`false`, `true`}

```
boolean b = true; // Variable with value true
```

165

166

## Relational Operators

- a < b (smaller than)
- a >= b (greater than)
- a == b (equals)
- a != b (unequal)

number type  $\times$  number type  $\rightarrow$  boolean

## AND( $x, y$ )

- “logical and”

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

- $F$  corresponds to “false”.
- $T$  corresponds to “true”.

## $x \wedge y$

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

## Logischer Operator $\&\&$

a  $\&\&$  b (logical and)

boolean  $\times$  boolean  $\rightarrow$  boolean

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

## Boolean Functions in Mathematics

- Boolean function

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

- $F$  corresponds to “false”.
- $T$  corresponds to “true”.

$\text{OR}(x, y)$

$x \vee y$

**Logical Operator ||**

- “logical or”

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

- $\textcolor{red}{F}$  corresponds to “false”.

- $\textcolor{blue}{T}$  corresponds to “true”.

$x$	$y$	$\text{OR}(x, y)$
$\textcolor{red}{F}$	$\textcolor{red}{F}$	$\textcolor{red}{F}$
$\textcolor{red}{F}$	$\textcolor{blue}{T}$	$\textcolor{blue}{T}$
$\textcolor{blue}{T}$	$\textcolor{red}{F}$	$\textcolor{blue}{T}$
$\textcolor{blue}{T}$	$\textcolor{blue}{T}$	$\textcolor{blue}{T}$

$a \text{ || } b$  (logical or)

`boolean × boolean → boolean`

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

171

172

$\text{NOT}(x)$

$\neg x$

**Logical Operator !**

- “logical not”

$$f : \{\textcolor{red}{F}, \textcolor{blue}{T}\} \rightarrow \{\textcolor{red}{F}, \textcolor{blue}{T}\}$$

$x$	$\text{NOT}(x)$
$\textcolor{red}{F}$	$\textcolor{blue}{T}$
$\textcolor{blue}{T}$	$\textcolor{red}{F}$

$!b$  (logical not)

`boolean → boolean`

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

173

174

## Precedences

$$\begin{array}{l} !b \And a \\ \Downarrow \\ (!b) \And a \end{array}$$
$$\begin{array}{l} a \And b \Or c \And d \\ \Downarrow \\ (a \And b) \Or (c \And d) \end{array}$$
$$\begin{array}{l} a \Or b \And c \Or d \\ \Downarrow \\ a \Or (b \And c) \Or d \end{array}$$

## Precedences

The *unary logical* operator ! provides a stronger binding than *binary arithmetic* operators. These bind stronger than *relational* operators, and these bind stronger than *binary logical* operators.

```
7 + x < y && y != 3 * z || !b  
7 + x < y && y != 3 * z || (!b)
```

175

176

## DeMorgan Rules

- $!(a \And b) == (!a \Or !b)$
- $!(a \Or b) == (!a \And !b)$

$!(\text{rich and beautiful}) == (\text{poor or ugly})$

## Application: either ... or (XOR)

$(x \Or y) \And !(x \And y)$  x or y, and not both

$(x \Or y) \And !(x \Or y)$  x or y, and one of them not

$!(!x \And !y) \And !(x \And y)$  not none and not both

$!(!x \And !y) \Or x \And y$  not: both or none

177

178

## Shortcut Evaluation

- Logical operators `&&` and `||` evaluate the *left operand first*.
- If the result is then known, the right operand will *not be* evaluated.

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

⇒ No division by 0