

**Base case**  
 $N = d_n R^n + d_1 R^1 + d_0 R^0$   
 the  $d$  specifies the Number system  $\rightarrow d_2 ==$  binary  
 can also be written as  $R_2$   
 This can also be used to expand numbers:  
 $N_{10} 255 = 2 * 10^2 + 5 * 10^1 + 5 * 10^0$   
 $N_2 110 = 1 * 2^2 + 1 * 2^1 + 0 * 2^0 \Rightarrow N_{10} 6$

**Quantities:**  
 $N \rightarrow$  natural numbers ||  $Z \rightarrow$  full numbers  
 $Q \rightarrow$  rational numbers ||  $R \rightarrow$  real numbers

**Common number systems:**  
 Decimal:  $N_{10} = n * 10^n \dots 0 * 10^0$   
 Binary:  $N_2 = n * 2^n \dots 0 * 2^0$   
 $2^{10} = 1024, 2^9 = 512, 2^8 = 256, 2^7 = 128, 2^6 = 64,$   
 $2^5 = 32, 2^4 = 16, 2^3 = 8, 2^2 = 4, 2^1 = 2, 2^0 = 1$   
 Hexadecimal:  $N_{16} = n * 16^n \dots 0 * 16^0$   
 notation: 0 1 2 3 4 5 6 7 8 9 A B C D E F  
 $16^5 = 1048576, 16^4 = 65536, 16^3 = 4096, 16^2 = 256,$   
 $16^1 = 16, 16^0 = 1$

**Modulo**  
 $8 \bmod 4 = (8) \rightarrow 0, 8 \bmod 3 = (6) \rightarrow 2, 8 \bmod 5 = (5) \rightarrow 3$   
 if  $x < y$  in  $x \bmod y$  then the result will always be  $x!$   
 any negative numbers can be considered as NOTnegative  
 aka only absolute values! modulo deals with  $|x|$   
 many programming languages actually do not follow this!  
 they have their own implementation of modulo.  
 $5 \equiv 3 \bmod 2 \rightarrow$  as  $5 \bmod 2 = 1$  and  $3 \bmod 2 = 1$

**Codeword length**  
 Byte = 8 bit || Word = 16 or 32 bit  
 TCP packet = 1024 bit

**Cyclic group**  
**Es sei  $F(a) = a^3 + a + 1 = 0$ ,**  
 ■ Dann können wir zunächst festhalten  
 ■  $a = a$   
 ■  $a^2 = a^2$  aber  
 ■  $a^3 = a + 1$   
 ■  $a^4 = a(a + 1) = a^2 + a$   
 ■  $a^5 = a(a^2 + a) = a^3 + a^2 = a^2 + a + 1$   
 ■  $a^6 = a(a^2 + a + 1) = a^3 + a^2 + a = a + 1 + a^2 + a = a^2 + 1$   
 ■  $a^7 = a(a^2 + 1) = a^3 + a = a + 1 + a = a + 1$   
 ■  $a^8 = a$  : der Zyklus beginnt von vorne!  
 ■  $\{0, 1, a, a^2, a + 1, a^2 + a, a^2 + a + 1, a^2 + 1\}$   
 ■  $\{000, 001, 010, 100, 011, 110, 111, 101\}$

**WHAT THE FUCK**  
**Result Quantity** the result of all possible outcomes  
 it is denoted with:  $\Omega$   
 A single element of the result list is:  $\omega \rightarrow \omega \in \Omega$   
 The list of results is  $|\Omega|$   
 Example Dice roll:  $\Omega = \{1, 2, 3, 4, 5, 6\}$

Probability:  $P(A) = \frac{\text{best results}}{\text{all results}} = \frac{|A|}{|\Omega|} = \frac{|A|}{n}$   
 So what is the probability of rolling a 6?  
 $P(\text{desired number to roll}) = \frac{\text{only 1 good result!}}{6 \text{ possible results}} = \frac{1}{6}$   
 hence the chance is 1 in 6  
 Why this complicated method? You can modify desired results!  
 just change the A in P(A)!

**Inverse Probability:**  $P(\text{inverse}) = 1 - P(A)$   
 dice  $\rightarrow 1 - \frac{1}{6} = \frac{5}{6}$

**Addition rule:**  
 $P(A \cup B) = P(A) + P(B) - P(A \cap B)$   
 !!The last part is needed, as otherwise the number  
 would exceed the possible states!!  
 $P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$

**Amount of possibilities:**  
**ordered probes with replication:**  
 2 coins, head and tail, possibilities?  $k = \text{head/tail} = 2$   $n = \text{coins} = 2$   
 $\Omega = n^k = 2^2$   
**ordered probes without replication:**  
 5 dices. How many combinations?  
 dice numbers =  $n = 6$  (1-6), dice amount =  $k = 5$   
 $\text{possibilities} = \Omega = \frac{n!}{(n-k)!} = \Omega = \frac{6!}{(6-5)!} = 720$   
 Or this:  
 $\Omega = \Pi_{n-k+1}^n = \Pi_6^{6-5+1} 6 = 2 * 3 \dots 5 * 6 = 720$   
**unordered probes without replication:**  
 25 players, each should only play once with the other.  
 $\Omega = \frac{n!}{k!(n-k)!} \rightarrow \frac{25!}{2!(25-2)!} \rightarrow \frac{\text{too big}}{\text{too big}} = 300$   
 as you can see the bottom is a BIG calculation, so  
 $\Omega = \frac{\Pi_{n-k+1}^n}{k!} \rightarrow \frac{\Pi_{25-2+1}^{25}}{2!} \rightarrow \frac{24 * 25}{2} = 300$

Note that  $k$  can also be defined as the  
 length of the tuple we want to receive.  
 $\rightarrow (\text{Player}, \text{Player}) \rightarrow 2$

**Source to Sink Information**

Nachricht (Darstellung & Bedeutung)	redundant	nicht-redundant
irrelevant	Zeichenvorrat bei Quelle und Senke verschieden	
relevant	vorhersagbar	Information

**Entropy**  
 information content  
 this essentially just us how many bits are needed  
 $k$  is base state count  $\rightarrow$  bit = 2  
 and  $N$  is the full number of states  
 example: list True, False, True, False 4 states total, base 2.  
 $H_0 = \log_k(N)[k] \rightarrow H_0 = \log_2(4)[\text{bit}] = 2$

**information flow**  
 essentially information content over time  
 $H_0^* = \frac{\log_2(N)}{\tau} [\frac{\text{bit}}{s}]$   
 information quantity / Surprise  
 $I(x_k) = -\log_2(P(x_k))[\text{bit}]$   
 Entropy (Surprise per element)  
 0 means no symbols. 1 means perfect balance 50-50  
 $H(X) = \sum_{k=1}^N P(x_k) * I(x_k) [\frac{\text{bit}}{\text{symbol}}]$   
 where  $X$  is the list of symbols  
 Sink Redundance / Code Redundance  
 $R_Q = H_0 - H(X) [\frac{\text{bit}}{\text{symbol}}]$   
 $R_c = L - H(X) [\frac{\text{bit}}{\text{symbol}}]$   
 Code Word Length  
 $L(x_k) = \text{rounded}(I(x_k))[\text{bit}]$   
 Median Code Word Length  
 $L = \sum_{k=1}^N P(x_k) * L(x_k) [\frac{\text{bit}}{\text{symbol}}]$   
 Entropy of the entire Code  
 $H_c(X) = \sum_{k=1}^N P(x_k) * L(x_k) [\frac{\text{bit}}{\text{symbol}}]$   
 $H_c$  can be a real number  $\rightarrow H_c \in \mathbb{R}$

**Für jede beliebige zugehörige Binärcodierung mit Präfixeigenschaft ist die mittlere Codewortlänge nicht kleiner als die Entropie  $H(X)$ :**  
 $H(X) \leq L$

**Für jede beliebige Quelle kann eine Binärcodierung gefunden werden, so dass die folgende Ungleichung gilt:**  
 $H(X) \leq L \leq H(X) + 1$

Sink without memory  
 $P(x_k, y_k) = P(x_k) + P(y_k)$   
 Sink with memory  
 $P(x_k, y_i) = P(x_k) + P(x_k | y_i)$   
 Entropy without memory / Combined Entropy  
 $H(H, Y) = \sum_{x_k} \sum_{y_i} P(x_k, y_i) * (-\log_2(P(x_k, y_i)))$   
 or:  $H(X, Y) = H(X) + H(Y)$   
 Entropy with memory  
 $H(H, Y) = \sum_{x_k} \sum_{y_i} P(x_k) * P(x_k, y_i) * (-\log_2(P(x_k) * P(x_k | y_i)))$

**Encoding of Symbols**  
 • Ordne die Zeichen gemäss ihrer Auftretswahrscheinlichkeit  
 • Die beiden Zeichen mit der kleinsten Auftretswahrscheinlichkeit haben die gleiche CW-Länge  $L_k$   
 • Sei  $L_k$  die mittlere CW-Länge für eine Quelle mit  $N$  Zeichen und  $L_{N-1}$  die mittlere CW-Länge für den Fall, dass die beiden letzten zu einem einzigen Zeichen zusammengefasst werden, dann gilt:  
 $L_N - (p(x_{N-1}) + p(x_N)) : L(x_N) = L_{N-1} - (p(x_{N-1}) + p(x_N)) : (L(x_N) - 1)$   
 $\Rightarrow L_N = L_{N-1} + p(x_{N-1}) + p(x_N)$

1	2	3	4	5	6	7	8	9
0.22	0.19	0.15	0.12	0.08	0.07	0.07	0.06	0.04
1	2	3	4	8	9	5	6	7
				0	1			
0.22	0.19	0.15	0.12	0.1	0.08	0.07	0.07	
1	2	3	6	7	4	8	9	5
			0	1		0	1	
0.22	0.19	0.15	0.14	0.12	0.1	0.08		
1	2	8	9	5	3	6	7	4
		00	01	1		0	1	
0.22	0.19	0.18	0.15	0.14	0.12			

continue this pattern until every symbol has a code  
 note the extra 0 on every step

**Run Length Encoding RLE/RLC**  
 □ Quelltext  $w$ : A3g2beh3f4g  $\Rightarrow |w| = 15$   
 □ Codiert  $w_k$ : A3g2beh3f4g  $\Rightarrow |w_k| = 11$   
 A+3xg+2xb+e+h+3xf+4xg  
 shortening of length  
 by compressing repetition.

**Encoder and Decoder**  
 You need to either choose 1 or 0 as the starting bit. After that the decoder can print out the correct code.

**Chiffre text**  
 You can "encrypt" your data by shifting the codes by a certain amount.  
 In the caesar chiffre this is done with the number 4.  $a \rightarrow e$   
 Please do not use this, use RSA or other algorithms.

**Errors**

$p(x_1) = 0.5$   $x_1$   $\begin{matrix} p \\ 1-p \end{matrix}$   $\begin{matrix} \rightarrow y_1 \\ \rightarrow y_2 \end{matrix}$   $p(y_1) = p(x_1) \cdot p + p(x_2) \cdot (1-q)$   
 $p(x_2) = 0.5$   $x_2$   $\begin{matrix} 1-p \\ q \end{matrix}$   $\begin{matrix} \rightarrow y_1 \\ \rightarrow y_2 \end{matrix}$   $p(y_2) = p(x_1) \cdot (1-p) + p(x_2) \cdot q$

$p(y|x) = \begin{bmatrix} p & 1-p \\ 1-q & q \end{bmatrix} \rightarrow \begin{bmatrix} \sum_{x=1}^2 1 \\ \sum_{x=1}^2 1 \end{bmatrix}$

1-p and 1-q are the chance for error. Which we of course have to take into account.

$p(x_1) = 0.5$   $x_1$   $\begin{matrix} 0.95 \\ 0.025 \end{matrix}$   $\begin{matrix} \rightarrow y_1 \\ \rightarrow y_2 \end{matrix}$   $\begin{bmatrix} p(x_1) & p(x_2) & p(x_3) & p(x_4) & p(x_5) & p(x_6) & p(x_7) & p(x_8) \\ p(x_1) & p(x_2) & p(x_3) & p(x_4) & p(x_5) & p(x_6) & p(x_7) & p(x_8) \\ p(x_1) & p(x_2) & p(x_3) & p(x_4) & p(x_5) & p(x_6) & p(x_7) & p(x_8) \\ p(x_1) & p(x_2) & p(x_3) & p(x_4) & p(x_5) & p(x_6) & p(x_7) & p(x_8) \end{bmatrix}$   
 $p(x_2) = 0.25$   $x_2$   $\begin{matrix} 0.95 \\ 0.025 \end{matrix}$   $\begin{matrix} \rightarrow y_1 \\ \rightarrow y_2 \end{matrix}$   
 $p(x_3) = 0.25$   $x_3$   $\begin{matrix} 0.95 \\ 0.025 \end{matrix}$   $\begin{matrix} \rightarrow y_1 \\ \rightarrow y_2 \end{matrix}$

$p(y|x) = \begin{bmatrix} 0.95 & 0.025 & 0.025 \\ 0.025 & 0.95 & 0.025 \\ 0.025 & 0.025 & 0.95 \end{bmatrix}$   $\begin{bmatrix} 0.4875 & 0.5095 & 0.25025 & 0.25025 \\ 0.25625 & 0.50025 & 0.25095 & 0.25025 \\ 0.25625 & 0.50025 & 0.25025 & 0.25095 \end{bmatrix}$

**Conditional Entropy  $\rightarrow$  Entropy of Y given X**  
 $H(Y|X) = \sum_{k=1}^N \sum_{i=1}^N P(x_k, y_i) * (-\log_2(\frac{P(x_k, y_i)}{P(x_k)}))$

**Chain Rule**  
 $H(Y|X) = H(X, Y) - H(X) \parallel H(Y \setminus X)$

**Bayes Rule**  
 $H(Y|X) = H(X|Y) - H(X) + H(Y) \parallel H(Y \setminus X)$

**H(X)** **H(Y)** **H(X,Y)** **H(X|Y)** **H(Y|X)** **I(X;Y)**

**Transinformation**  
 likelihood of information being correct at arrival.  
 $T = H(X) - H(X|Y) \parallel H(Y) - H(Y|X)$   
 or:  $|I(X;Y)|$

**Hamming distance / distance to next valid codeword**  
 $h = \text{Min}_{i,j} (d(x_i, x_j)) \rightarrow h = \text{message points} + 1$   
 error detection distance  
 the amount of bits that differ from input to output  
 $e^* = h - 1$   
 error correction distance for h even  
 $h = 2e + 2 \rightarrow e = \frac{h-2}{2}$   
 error correction distance for h uneven  
 $h = 2e + 1 \rightarrow e = \frac{h-1}{2}$   
 Consider the valid input either 111 or 000.  
 The Hamming distance  $h$  is therefore 3 bits.  
 The detection distance  $e^*$  is 3 - 1  
 Due to h being uneven, the correction distance  $e$  is  $\frac{h-1}{2}$   
 which results in 1.

**tightly packed coderoom**  
 $n =$  dimension of code  
 $m =$  dimension of messages  $2^m * \sum_{w=0}^e \binom{n}{w} \leq 2^n$   
 $k =$  dimension of control  $\rightarrow n = m + k$   
 The code is considered to be tightly packed  
 if the equation has the result 2. aka == not smaller.

$x_1$	$x_2$	$x_3$
0	0	0
0	1	1
1	0	1
1	1	0
0	0	1
0	1	0
1	0	0
1	1	1

$\theta = (\theta + \theta) \bmod 2 = 0$  OK  
 $1 = (\theta + 1) \bmod 2 = 1$  OK  
 $1 = (1 + \theta) \bmod 2 = 1$  OK  
 $\theta = (1 + 1) \bmod 2 = 0$  OK  
 $x_3 = (x_1 + x_2) \bmod 2$   
 $1 = (\theta + \theta) \bmod 2 = 0$  NOT OK  
 $\theta = (\theta + 1) \bmod 2 = 1$  NOT OK  
 $\theta = (1 + \theta) \bmod 2 = 1$  NOT OK  
 $1 = (1 + 1) \bmod 2 = 0$  NOT OK

**Hamming Codes**  
 The hamming code is very easy to implement  
 $\Sigma_i x_i * \vec{P}_i \equiv \vec{0} \bmod 2$   
 The syndrome  $\vec{Z} = \Sigma_i x_i * \vec{P}_i \bmod 2$   
 1,2,4,8,16...  $2^x$  are parity checks

parity bit	0	1	000	001		
2	3	010	011			
4	5	100	101			
6	7	110	111			

$\blacksquare = Q_1$   
 $\blacksquare = Q_2$   
 $\blacksquare = Q_3$

**example for code 1001**

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

$\vec{Z} = 000 = \text{no error}$   $\vec{Z} = 101 = \text{error at } 101 = 5$   
 note that the 001 010 100 of the parity checks are simply the unit vector  $\vec{0} \parallel$  !!!

parity checks needed:  
 $par = \log_2(\text{bit amount of code})$   
1101 = 4 bits -> 3 parity checks  
as 4 can be displayed by 3 bits -> 100

OR we can use XOR

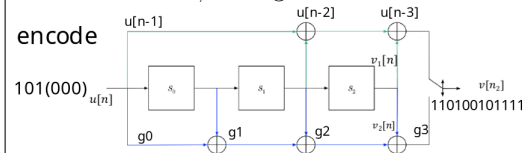
write down all the positions with 1  
in binary:

001 001  
011 011  
101 **error**

errors can be represented by  
either removing one bit  
from the calculation -> setting  
it to 0 or adding it, setting it to 1

$\oplus 111$   
 $\oplus 111$   
 $000$   $101$

convolutional code / folding code

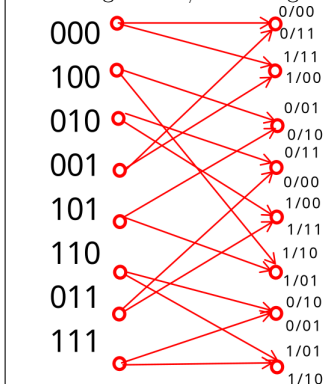


u[n]	u[n-1]	u[n-2]	u[n-3]	v1[n]	v2[n]
1	0	0	0	1	1
0	1	0	0	0	1
0	1	0	0	0	1
1	0	1	0	0	0
0	1	0	1	1	0
0	0	1	0	1	1
0	0	1	0	1	1
0	0	0	0	-	-

this algorithm is heavily dependent on the configuration, aka the more states you have the more bits you get as the output.

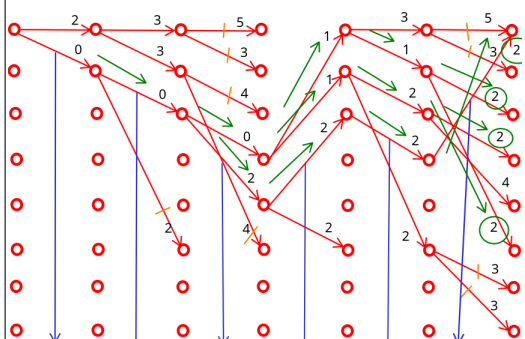
$\{u[n]\} = \{u_1, u_2, \dots, u_n\}$   
 $\{v[n]\} = \{g^0 u[n] + g^1 u[n-1] + g^2 u[n-2] + g^3 u[n-3], g^0 u[n+1-m], \dots\}$   
 $\{v[n]\} = \sum_{m=0}^M g^m u[n-m]$

Viterbi algorithm / Decoding convolutional code



now check what the output was with the values of the changes above the 1/10... 0/00...

consider the output from above with 2 errors: 110111011111  
the numbers on the line represent the difference to the code given  
->Hamming Distance



1/11 0/01 **1/00** 0/10 0/11 0/11  
!!corrected error!!

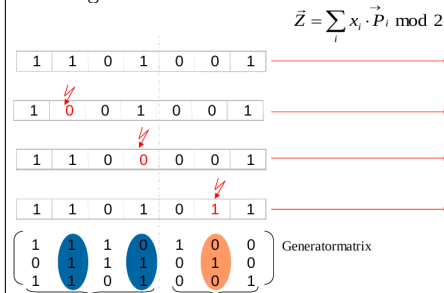
-> 101000

There is a big problem with this algorithm, as you can see there were multiple ways of solving it, and the others would NOT have given the same code. So multi error can still be tricky.

Input - Output ratio for convolutional codes

$$R = \frac{\text{input}}{\text{output}} \text{ in example } R = \frac{1}{2}$$

hamming code extension



message with code check

1101 001  
1 x 101 -> 1 x 101  
1 x 111 -> 1 x 111  
0 x 110 -> 1 x 011  
1 x 011 -> **1 x 001**  
0 x 100 -> 000  
0 x 010 -> 000  
1 x 001 -> **correct**

message with code check

1111 001  
1 x 101 -> 1 x 101  
1 x 111 -> 1 x 111  
**1 x 110** -> **1 x 110**  
1 x 011 -> 1 x 011  
0 x 100 -> 1 x 001  
0 x 010 -> **110**  
1 x 001 -> **wrong**  
**error at 110**

RSA: !!Check Github for another PDF!!

RINKEL DONE

due to space constraints, something form the last page:

Floating Point Comparison

Due to the limitation of floating point numbers, checking for equality the regular way doesn't make sense  
for (int = 0; i < n; ++i) { a += f; }  
b = ((float) n) \* f;  
if (a == b) { ... }

This checks bit pattern!!

Instead we check if it is approximately equal:

if (Math.abs(a - b) <= eps) { ... }

Problem with this: it's the absolute error instead we can choose the relative error:

float aa = Math.abs(a);

float bb = Math.abs(b);

if (Math.abs(a - b) <= eps \* Math.max(a, b))

The Max makes sure the eps is based on the bigger number.

There are still certain problems: calculations close to 0 might result in the relative error of 1, which is actually not an error once again rounding errors...

BE AWARE OF FLOATING ERRORS!!

potence system

decimal binary and hexadecimal are all built the same way

$$W_a(z) = \sum_{i=0}^{n-1} z_i * a^i = z_{n-1} * a^{n-1} + z_{n-2} * a^{n-2} + \dots + z_1 * a^1 + z_0 * a^0$$

Terms in the binary form

Bit = 1 set bit

Bit = 0 cleared bit

LSB Least Significant Bit, rightmost bit, bit 0

MSB Most Significant Bit, leftmost bit, bit (n-1)

!!Note that this is for little-endian!!

big-endian is the opposite!!

Nibble a 4 bit binary number

bigger numbers get grouped into nibbles -> 1101 0101

Octect a binary number with 8 bit

byte 8 bit

Note, in binary the leading 0s will still be written, in order to indicate how big the memory allocation is.

indices in bits in order to change single bits in a number, we can use the array notation. b=1010 b[3]=1 b=1011

amount of numbers to display

binary:  $2^n$  numbers -> n=bits

The amount of numbers to double every increase of potence.

Note: There are no broken binary numbers. It's always  $2^n$

with n being an integer!!

prefixes

		K	Kilo	Ki	Kibi
$2^{10}$	$1.024 \cdot 10^3$	M	Mega	Mi	Mebi
$2^{20}$	$1.049 \cdot 10^6$	G	Giga	Gi	Gibi
$2^{30}$	$1.074 \cdot 10^9$	T	Tera	Ti	Tebi
$2^{40}$	$1.100 \cdot 10^{12}$	P	Peta	Pi	Pebi
$2^{50}$	$1.126 \cdot 10^{15}$	E	Exa	Ei	Exbi
$2^{60}$	$1.153 \cdot 10^{18}$				

suffixes: **binary:1011<sub>b</sub>, hex:1011<sub>h</sub>**

Konvention Dezimal Hexadez Binär

Intel Manual 1011 1011H 1011B

Wikipedia 1011<sub>10</sub> 1011<sub>16</sub> 1011<sub>2</sub>

C, C++, Java, Python, PHP 1011 0x1011 0b1011

Pascal, Delphi 1011 \$1011 %1011

Ada 1011 16#1011 2#1011

Visual Basic 1011 &H1011 &B1011

Hex

1 hex decimal is a Nibble, 4 bit

2 hex decimals are a byte

this is why bits are grouped into nibbles and bytes.

Notation: 12 34 -> 0001 0010 0011 0100

just like with binary, we write the leading 0s for the memory size.

biggest number possible with n bit:

	9 Bit	1FF	13 Bit	1FFF
:	10 Bit	3FF	14 Bit	3FFF
7 Bit	7F	11 Bit	7FF	:
8 Bit	FF	12 Bit	FFF	:

Logical Operations

It is a function from n bits to 1 bit -> bool

$$f: \mathbb{B}_n \rightarrow \mathbb{B} \parallel f :: [a] \rightarrow Bool$$

parameter: variable to use inside the function

argument: specific value used in a specific use can also be described as fixed values. Note f(a,b) = a + g(a,b) in this function the a and b in f are parameters, but in g they are arguments !!

Unary function: a function with 1 parameter

Binary function: a function with 2 parameters

Ternary function: a function with 3 parameters

n-ary function: a function with n parameters

Nullary function: function with zero parameters

There are exactly 4 unary functions (with bits!)

F(X) = 0 -> False || T(X) = 1 -> True

id(x) = x, just like haskell id...

not(x) =  $\neg x$

There are 16 binary functions

Disjunction: OR  $x \vee v$

Conjunction: AND  $x \wedge v$

The others are made with unary, nullary or the 2 above

x	y	0	xy	xy	x	xy	y	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	1
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
0	1	0	0	0	0	1	1	1	0	0	0	1	1	1	1
1	0	0	0	1	1	0	0	1	0	0	1	1	0	0	1
1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1

x	y	0	xy	x	xy	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	$\neg x \vee y$	1
0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
0	1	0	0	0	1	1	1	1	0	0	0	1	1	1	1
1	0	0	0	1	0	0	1	1	0	0	1	1	0	0	1
1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1

note that some are actually unary or nullary, aka they don't actually do anything with the parameters.

AND: the basis of modern computers  
it is easy to create with transistors.  
It only results in False if both inputs are true,  
hence the name Not-AND.  
mathematical notation:  $NAND = x|y = \overline{x \wedge y}$   
all base operations can be made with NAND  
 $\overline{x} = x \wedge \overline{x} = x | x$   
 $x \wedge y = \overline{x | y} = (x | y) | (x | y)$   
 $x \vee y = \overline{x \wedge \overline{y}} = (x | x) | (y | y)$

NOR  
NOR is only true when neither of the inputs are true.  
aka NOR = 1 if x==0 && y==0  
 $\overline{x \vee y}$   
Just like with NAND, all functions  
can be made with NOR  
XOR: exclusive or  
The XOR is true if only one input is true.  
 $x \oplus y$

Addition of bits.

x	y	x + y	x ∧ y	x ⊕ y
0	0	00	0	0
0	1	01	0	1
1	0	01	0	1
1	1	10	1	0

AND signifies the overflow of bits. 1 0 <- 1 1  
XOR signifies the addition of bits. 0 OR 1  
Literal: variable or negation of variable:  $x_1, \overline{x_2}$   
Conjunction-term: conjunction of literals  $x_1 x_2 = x_1 \wedge x_2$   
Disjunction-term: disjunction of literals  $x_1 \vee x_2$   
Minterm: conjunction with ALL parameters of a function  
Maxterm: disjunction with ALL parameters of a function

Disjunctive Normalform DNF  
a disjunction of conjunctions.  $x_1 x_2 \vee \overline{x_1} \overline{x_2}$   
Functions are often displayed as DNF, since this format  
requires only 3 symbols:  $\vee, \wedge, \neg$   
The canonical DNF C-DNF: a DNF with all parameters  
the canonical DNF is often used to display the true false table.  
The C-DNF is then often simplified to get the end result

x	y	$x \oplus y = \overline{x}y \vee x\overline{y}$	$x y = \overline{x\overline{y} \vee \overline{x}y \vee x\overline{y}}$
0	0	0	1
0	1	$\overline{x}y$	1
1	0	$x\overline{y}$	1
1	1	0	0

$\overline{x\overline{y} \vee \overline{x}y \vee x\overline{y}} = \overline{x\overline{y}} \vee \overline{\overline{x}y} = \overline{x} \vee y = \overline{x} \vee \overline{\overline{y}}$   
 $x \oplus y = \overline{x}y \vee x\overline{y}$

note that DNF has nothing to do with the ferrari engine.

logical function only return 1 or 0. In order to  
return an entire number, we would have to map this function.  
luckily there are several predefined mapped functions  
- nor(x) invert all bits in x  
- and(x,y) check the individual bits of x,y with and  
- or(x,y) check the individual bits of x,y with or  
- nor(x,y) check the individual bits of x,y with nor

Bitwise operations in java...

NOT:  $z = \neg q$   $z_i \leftarrow \neg q_i$   
AND:  $z = q \& p$   $z_i \leftarrow q_i \wedge p_i$   
OR:  $z = q | p$   $z_i \leftarrow q_i \vee p_i$   
XOR:  $z = q \wedge p$   $z_i \leftarrow q_i \oplus p_i$

|| and && are evaluations  
note that with || and &&, if the first evaluation is enough  
to determine the result, the second one won't be executed.  
a=False, b=True -> a && b -> a is false, therefore result false.

Variable sizes in java  
int = 32 bit  
short = 16 bit  
byte = 8 bit  
long 64 bit  
Please note that these should only be used when,  
you either save significant memory, or the integer isn't big enough.

multiplying a binary number  
you have to multiply every single bit and add the corresponding 0s

$2^2 \cdot 101011$   
 $4 \cdot 101011$   
 $100 \cdot 101011$   
 $000000$   
 $000000$   
 $10101100$

Multiplication of binary number  
 $b = 2^{n-1} \cdot b_{n-1} + \dots + 2^0 \cdot b_0$   
Multiplication with potences. -> also binary!  
 $c = 2^m \cdot (2^{n-1} \cdot b_{n-1} + \dots + 2^0 \cdot b_0)$   
 $= 2^m \cdot 2^{n-1} \cdot b_{n-1} + \dots + 2^m \cdot 2^0 \cdot b_0$   
 $= 2^{m+n-1} \cdot b_{n-1} + \dots + 2^m \cdot b_0$   
 $= 2^{m+n-1} \cdot c_{m+n-1} + \dots + 2^m \cdot c_m + 2^{m-1} \cdot 0 + 2^0 \cdot 0$   
 $c_{m+n-1} = b_{n-1}, \dots, c_m = b_0, c_{m-1} = \dots = c_0 = 0$

Multiplication with potences is just a leftshift!!  
 $2^4 \cdot 101 = \text{add 4 0s to } 101 \rightarrow 101'0000$

with potences is a rightshift!!  
 $1'0111 = \text{remove 3 bits from } 1'0111 \rightarrow 101$   
 $\overline{23}$

Left & Rightshift java  
logical right:  $a \gg x \rightarrow 101 \gg 2 \rightarrow 1 \rightarrow 010$   
logical left:  $a \ll x \rightarrow 101 \ll 2 \rightarrow 1'0100$   
arithmetic right:  $a \gg x \rightarrow 101 \gg 1 \rightarrow 110$   
arithmetic left:  $a \ll x \rightarrow 101 \ll 2 \rightarrow 1'0111$   
the >> and << add the MSB value instead of 0  
The reason for this is unsigned and signed!!

Reading a bit  
 $b \wedge m = b_7 b_6 b_5 b_4 b_3 b_2 b_1 b_0 \wedge m_7 m_6 m_5 m_4 m_3 m_2 m_1 m_0$   
 $= b_7 b_6 b_5 b_4 b_3 b_2 b_1 b_0 \wedge 0100'0000$   
 $= 0b0'00'0000$   
Java b & 0b0100\_0000

Setting a bit  
 $b \vee m = b_7 b_6 b_5 b_4 b_3 b_2 b_1 b_0 \vee m_7 m_6 m_5 m_4 m_3 m_2 m_1 m_0$   
 $= b_7 b_6 b_5 b_4 b_3 b_2 b_1 b_0 \vee 0100'0000$   
 $= b_7 1 b_6 b_5 b_4 b_3 b_2 b_1 b_0$   
Java b v 0b0100\_0000

Deleting a bit  
 $b \wedge m' = b_7 b_6 b_5 b_4 b_3 b_2 b_1 b_0 \wedge m'_7 m'_6 m'_5 m'_4 m'_3 m'_2 m'_1 m'_0$   
 $= b_7 b_6 b_5 b_4 b_3 b_2 b_1 b_0 \wedge 1011'1111$   
 $= b_7 0 b_6 b_5 b_4 b_3 b_2 b_1 b_0$   
Java b & 0b1011\_1111

Combine these with right and left shift!  
read:  $b \& (1 \ll n) \gg n$  000(b AND 1111...)  
set:  $(b | (1 \ll n)) \& 0b1111...$   
delete:  $b \& (1 \ll n)$  b AND not 11111..  
given a binary number c and a bit amount p. shift c  
by p.  
1. create mask, 2. read n bits from c into bb  
3. shift n bits from c to cc in order to create room  
for p bits, 4. set bb in cc

	m = (1 << p) - 1;	//
1	bb = b & ~(m << n);	//
2	cc = (c   m) << n;	//
3	r = bb   cc;	//

Addition and Subtraction  
Addition:

	1	1	1	1
+1	0 <sub>1</sub>	0 <sub>1</sub>	0 <sub>1</sub>	1
(1)	0	0	0	0

if the subtraction would  
result in 0, then we  
add another 1 to  
the number above!!

Subtraction: (Zweierkomplement)

	1	1	0	1
-	0	0 <sub>1</sub>	1	0
(1)	0	1	1	1

(1) 0 0 0 0 0  
(0) 1 0 1 0 1 0 1 1  
(0) 1 1 1 1 1  
N(1) =  $2^n - 1 = \frac{1 \dots 1}{n}$   
N(2<sup>n-1</sup> - 1) =  $2^n - (2^{n-1} - 1)$   
 $= 2 \cdot 2^{n-1} - 2^{n-1} + 1$   
 $= 2^{n-1} + 1$   
 $= 10 \dots 01$   
 $n-2$

Signed & Unsigned  
Unsigned: only positive integers!  
Signed: MSB 0 = positive, MSB 1 = negative  
the rest is a regular binary number.  
note: This is something you simply need to know.  
It isn't included in some encoding!!

	signed	unsigned
4 Bit	-8 ... 7	0 ... 15
8 Bit	-128 ... 127	0 ... 255
16 Bit	-32K ... 32K - 1	0 ... 64K - 1
32 Bit	-2G ... 2G - 1	0 ... 4G - 1

note the difference in positive and negative  
in the signed category! size constraints!

	0000'0000 <sub>b</sub>	...	0111'1111 <sub>b</sub>	1000'0000 <sub>b</sub>	...	1111'1111 <sub>b</sub>
unsigned	0	...	127	128	...	255
signed	0	...	127	-128	...	-1

Note -1 is 1111'1111. Negative numbers are  
calculated: 0 - number  
This means an overflow on unsigned ints will lead  
to it being 0 again.  
On signed ints, it will drop to negative maximum.  
Special cases in signed:  
 $2^{n-1}$  max -> always negative as MSB = 1  
Max : 100000000...  
0 -> can't be negative as this bit is used for  $2^{n-1}$

Note that when increasing memory for signed values,  
you need to use the > operators to copy the MSB.  
increasing memory for -1 -> 4 bit to 8 bit  
0000'1111 = 15 !! WRONG !! -> 1111'1111 = -1!!

For signed left shift, check if you have spare memory  
left shift without checking might result in loss of MSB!  
4 bit max: 1001 << 2 = 0100 !! prefix changed !!

multiplication: series of left shifts.  
1101 \* 110 = 11'0100 + 1'1010 = 100'1110  
10 -> add 1 zero, 100 -> add 2 zeros, get the sum of both  
Size increase: max double ->  $x^2$   
110 \* 110 -> 1100 + 1'1000 = 11'1000 (3 to 6 bits)

$1111'1111_b = 1111'1111_b = 1111'1111'0000'0001_b = -1$   
signed Multiplication ≠ unsigned Multiplication  
check if both operands are negative, if so invert them to positive.  
do unsigned multiplication  
if only one operand was negative, take the negative result.  
 $\Rightarrow N_a(a) \cdot N_n(b) = a \cdot b$  and  $N_n(a) \cdot b = N_{2n}(a \cdot b) = a \cdot N_n(b)$  (mit  $a, b \geq 0$ )

Note that we don't need to care about signed, if  
we do not overwrite the MSB!

Division  
should be avoided, slow operation compared to others  
32 bit -> 20 times as long as multiplication  
64 bit -> 80 times as long as multiplication  
- can be replaced with right shift for potences  
see /10 for decimal numbers.  
signed and unsigned division are completely different

- Unsigned: Iteratives Verfahren for  $i = n - 1$  bis  $i = 0$ :
  - Man überprüft ob  $b \cdot 2^i$  in die  $n - i$  obersten Bits von  $a$  passt
  - Wenn ja, dann setzt man im Ergebnis Bit  $i$  und zieht  $b \cdot 2^i$  von  $a$  ab

Inversion  
 $N(b+1) > 2^n - 1 - b > 2^n == 0$   
 $0 - 1 = -1 > 11111... > -1 - b = \overline{b}$   
Okay, the idea is that -1 is the number with all bits set to 1  
This means that no matter what you do, you can't have overflow.  
In fact this means that b can be inverted by subtracting it to -1.

this is -1 !!

	1	1	1	1	1
-1					
b					
$\overline{b}$					

Unsigned in java  
You can't declare unsigned integer etc, instead you just use the  
unsigned functions.  
compareUnsigned, divideUnsigned, remainderUnsigned  
for Integer, Short, Byte, Long

Java comparators  
==, !=, <, >, <=, >=  
• a != b ↔ !(a == b)  
• a > b ↔ !(a < b)  
• a <= b ↔ (a < b) | (a == b)  
• a > b ↔ !(a <= b) ↔ !((a < b) | (a == b))  
we only need == <

Beause processors are super fast in addition an subtraction,  
we can simplify equality checks by using these 2 operations.  
- example unsigned, check if a < b:  
==, !=, <, >, <=, >=  
• a != b ↔ !(a == b)  
• a > b ↔ !(a < b)  
• a <= b ↔ (a < b) | (a == b)  
• a > b ↔ !(a <= b) ↔ !((a < b) | (a == b))  
we only need == <

c is the carry bit, and it is only set to 1 if a < b!!  
- example signed:  
case 1, overflow(wrong prefix)

	(+a) - (-b) = (-d)	0111 - 1000 = (1)1111	7 - (-8) = -1
	(-a) - (+b) = (+d)	1000 - 0111 = 0001	(-8) - 7 = 1

$0 = \overline{s_a} s_b s_d \vee s_a \overline{s_b} \overline{s_d}$  S = MSB of variable.

The 0 stands for overflow, however it only checks  
for an overflow of the prefix. Aka it checks wether  
or not the prefix makes sense.  
when we subtract something negative then we  
expect a positive outcome, which we DIDN'T get  
above!!

case 2 correct prefix

1. result is positive ->  $s_d = 0 \rightarrow a > b$   
2. result is negative ->  $s_d = 1 \rightarrow a < b$

$a < b \rightarrow 0 \oplus s_d = 1$

- the check for all 0 or 1 is also fast -> AND...  
In java we only work with the signed interpretation by default  
use the before mentioned special functions for unsigned!

Character Encoding  
A character is encoded via a bijectional function  
 $E'(z_0..z_{n-1}) = f(E(z_0)..E(z_{n-1}))$   
The formula must obviously be known by both parties, as the  
text can't be decoded otherwise.

ASCII  
 $2^7 = 128 - > 7bit$   
- uses both printable characters and  
nonprintable control characters  
- introduction 1963, enacted by US president 1968



30 control characters 0-9 33 special symbols ! " # ... 26 upper case letters A-Z 26 lower case letters a-z

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOF	ENQ	ACK	NEL	BS	TAB	LF	VT	FF	CR	SD	SI
1	DEL	DC1	DC2	DC3	DC4	DC5	DC6	DC7	DC8	DC9	DC10	DC11	DC12	DC13	DC14	DC15
2	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	P	Q	R	S	T	U	V	W	X	Y	Z	[	\	]	^	_
6	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

00<sub>h</sub> \0 NULL: marker for End of text  
08<sub>h</sub> \b Backspace  
09<sub>h</sub> \t Tabspace  
0A<sub>h</sub> \n Line Feed, Newline  
0D<sub>h</sub> \r Carriage Return

Problems with ASCII  
- Just another standard of many, even if it's the first  
- Egnlish only, other languages can't to be encoded in ASCII  
-> hence no ä ö ü ...

Extensions for ASCII / Codepage  
ASCII can be extended since it only uses 7 bit of the available 8. This means that the 8th bit can be used for other characters like ä.  
range: 80<sub>h</sub> to FF<sub>h</sub>  
Extension announced with MSB = 1  
No extension MSB = 0  
Extended ASCII remains compatible to regular ASCII  
Many different Codepages defined in ISO.  
Codepages have the problem that the sender and receiver need to again know what codepage is being used, since otherwise the receiver might not be able to decode certain characters properly.

Unicode  
UTF-8 (Standard) UTF-16 UTF-32  
max 17 · 64K - 2K = 1M + 62K = 1,112,064 Code-Points  
work in progress

Codepoints(CP): Number of a character  
CodeUnit(CU): Unit for displaying encoded character  
UTF-32 - all CP displayed with one CU  
UTF-16 - all CP displayed with 1-2 CUs  
UTF-8 - all CP displayed with 1-4 CUs  
UTF-16 and 32 need to check for endianness since they are above 8 bit.  
UTF-32 is limited to 21 bit!!  
Rest often used for other things.

UTF-8  
UTF-8 it itself another ASCII extension  
UTF-8 -> 7 bit for characters  
max 4 Bytes chained together -> 4x UTF-8 block

Codepoint in [80, 7FFF]:  
•  $U_1 = B_1 = 110 P_{10} \dots P_6 (+ C0_h)$   
•  $U_0 = B_0 = 10 P_5 \dots P_0 (+ 80_h)$

Codepoint in [800, FFFF]:  
•  $U_2 = B_2 = 1110 P_{15} \dots P_{12} (+ E0_h)$   
•  $U_1 = B_1 = 10 P_{11} \dots P_8 (+ 80_h)$   
•  $U_0 = B_0 = 10 P_7 \dots P_0 (+ 80_h)$

next would have 21 significant-bits

UTF-16  
max 4 Bytes chained together -> 2x UTF-16 block  
CP in [0, FFFF] without [D800, DFFF]  
[D800, DFFF] is not valid  
Rest  $U_0 = P_{15} \dots P_0$

CP in [1'0000, 10'FFFF] 0 not used here!! therefore - 10000h !!!  
•  $Q = P - 1'0000_h$  (-> Q in [0, FFFF] and  $Q_{15} \dots Q_0 = P_{15} \dots P_0$ )  
•  $U_1 = 1101'10 Q_{19} \dots Q_{10} (+ D800_h)$ ,  $U_0 = 1101'11 Q_9 \dots Q_0 (+ DC00_h)$

encode -1'0000<sub>h</sub>  
1-16 with 0-15, hence -10000h

• The entire block is always interpreted as Big-Endian!!

- UTF-16BE: Big Endian
  - CP with 1 CU:  $U_0 = B_1 B_0$
  - CP with 2 CUs:  $U_1 U_0 = B_3 B_2 B_1 B_0$
- UTF-16LE: Little Endian
  - CP with 1 CU:  $U_0 = B_0 B_1$
  - CP with 2 CUs:  $U_1 U_0 = B_2 B_3 B_0 B_1$

$U_1 (D800 \leq U_1 \leq DBFF)$  is the high surrogate  
 $U_0 (DC00 \leq U_0 \leq DFFF)$  is the low surrogate  
codepoints above 10000<sub>h</sub> are supplementary codepoints  
codepoints c with  $p0000_h \leq c \leq pFFFF_h$  are on plane p  
p = 0 is the basic multilingual plane -> 1 UTF-16 block  
 $1 \leq p \leq 16$  are the supplementary planes

UTF-32  
first 11 bits used for other things  
Big-Endian by default -> convert for little endian  
21 significant bits  
max 4 bytes -> ONLY 1 UTF-32 block  
Most used internally in programming languages and or applications.

Byte-Order-Mark (BOM) FEFF<sub>h</sub>  
If this is at the start of a file, then you can read the encoding from it:

UTF-16BE FE FF  
UTF-16LE FF FE  
UTF-8 EF BB BF  
UTF-32BE 00 00 FE FF  
UTF-32LE 00 00 FE FF

Encoding examples

CP	UTF-32BE	UTF-32LE	UTF-8	UTF-16BE	UTF-16LE
A	41	00 00 00 41	41	00 41	41 00
ä	E4	00 00 00 E4	C3 A4	00 E4	E4 00
α	3B1	00 00 03 B1	CE B1	03 B1	B1 03
ä	1EB7	00 01 E1 B7	E1 BA B7	1E B7	B7 1E
Ð	10330	00 01 03 30	F0 9C B0	D8 00 DF 30	00 D8 30 DF

Enian(ness)  
Big-endian and Little-endian  
Big-endian -> MSB leftmost bit, LSB rightmost bit  
Little-endian -> MSB rightmost bit, LSB leftmost bit

Big-Endian Little-Endian

Important since for example Intel processors have 1 Byte as the smallest addressable unit. This means you need to know the Endian of that byte, as you otherwise wouldn't be certain if you hit the LSB or the MSB first.

Endian Problem  
Little-endian is objectively better  
example: accessing the first 8 bits from a 16 bit number. It would make sense to only request the 8 bits, but this doesn't work for big Endian as it ends with the smallest number.

Big-Endian Little-Endian

For little endian, the access could have been smaller -> 8 bit. Thank fibonacci for fucking up his translation of numbers.

Endian Problem solution  
- shift to Mixed or Middle endian  
- Acceptance of both Endians  
These solutions are needed since it is unlikely that there will unity with endianness.

Encoding in Java  
Char receives a UTF-16 Code Unit  
Strings are an object with length and CodePoints  
This is different from C where you have to iterate to find the NULL symbol 0x00. Java uses more memory but it is faster to check length etc.

char c = 'i'; // char literal, in UTF-8: 0x31  
char d = '\u0031'; // Unicode codepoint literal  
char e = 1; // int literal converted to char 0x0001

charAt(i) codePointAt(i)  
chars codePoints  
length codePointCount  
charAt(i-1) codePointBefore(i)  
charAt(i+n) offsetByCodePoints(i, n)

Certain operations are not permitted on Strings directly such as + - delete, insert, replace etc as they would change the string, and a string is immutable!  
To mitigate this problem, use buffers and StringBuilder

ByteBuffer, CharBuffer, IntBuffer  
contains: array, capacity, limit and pos  
- array is the storage for the type -> char, int  
- capacity is the size of the entire buffer, it is immutable  
- limit is the amount of changeable elements inside the buffer  
- pos is the relative position for read and write operations

- byte get(): get array[pos], increment pos by 1
- get(byte[] d, int s, int c)  
copy array[pos]-array[pos+c-1] into d[s]-d[s+c-1]  
pos += c  
-> exception when pos+c >= limit
- get(byte[] d) = get(d, 0, d.length)
- put(byte b)  
put(byte[] d, int s, int c) same as get, inverted.  
put(byte[] d) = put(d, 0, d.length)

```
while (hasData)
{
    prepareData(...); // calls buffer.put(...);
    buffer.flip(); limit = pos -> pos = 0 ->
    //tell the buffer this is the end of the data.
    //We can now iterate through the data again from the
    //beginning. In other words, end of write, start read.
    processData(...); // calls buffer.get()
    buffer.clear(); clear the buffer for reuse.
}
```

InputStreams  
Sequential processing of bytes  
int read() -> 0-255 per byte  
int read(byte[] d, int s, int c) amount of read bytes  
int read(byte[] d)  
all of these function return -1 for exceptions  
FileInputStream  
reads data via functions from the operating system.  
should only be used in combination with read(byte[], ...)  
as it would be extremely slow otherwise.  
BufferedInputStream  
reads data form another Stream and buffers it.  
very fast read()  
if size is unknown then other streams should be used.  
var stream = new BufferedIS(new FileIS("..."));

InputStreamReader  
read form an InputStream  
has to convert bytes into chars  
It uses a decoder d as specified:  
ISR(IS s, CharsetDecoder d):  
ISR(IS s, Charset c): Wie ISR(s, c.newDecoder())  
ISR(IS s, String n): Wie ISR(s, Charset.forName(n))  
ISR(IS s): like ISR(s, Charset.defaultCharset())

Charsets in Java  
Charset.newDecoder() or Charset.newEncoder()  
create decoder or encoder.  
static SortedMap<String,Charset> Charset.availableCharsets()  
returns all available Charsets for this VM.  
FileReader (extends FileInputStream)  
- just calls super(FileInputStream) - only available with default en  
- for more encoders you have to call FileInputStream manually  
BufferedReader  
- reads data blockwise from a reader and buffers it  
check BufferedInputStream for similarity  
String readLine(): gives all chars until next line  
var reader = new BufferedReader(new FileReader("..."));

Streams and Resource-Leaks  
Because most Streams need resources from the OS, we need to close these streams again. Otherwise leak!  
The try is for possible exceptions, obviously  
FileReader fr = new FileReader(path);  
BufferedReader br = new BufferedReader(fr);  
try {  
 return br.readLine();  
} finally {  
 br.close();  
 fr.close();  
}

or the simplified and more safe way (Try-With-Resources):  
(em.close can cause exception!)  
try (FileReader fr = new FileReader(path);  
 BufferedReader br = new BufferedReader(fr))  
{  
 return br.readLine();  
}

- types have to extend the closeable type.  
- resources are automatically closed in the opposite manner  
- try with resources guarantees safe closing.

Units in Programming languages Obviously there is no ty  
We therefore use comments to signify the km or kg for humans  
Keep in mind that this means we also have to implement the  
conversion ourselves, as the computer doesn't know you want  
to calculate meters + kilometers!!  
// in km/h int length1 = 10; // in cm  
int speed; int length2 = length1 \* 10; // in mm  
int tempC = 20; // in C  
// in CHF int tempF = tempC \* 9 / 5 + 32; // in F  
int amount; Doesn't work!! TRUNCATE!!

Fixed Point Numbers: Integers with Commas!  
int sizeK = 5; // in K  
int sizeMk = sizeK \* (1 << 10); // in mK  
adding K to K then converting is very easy:  
0 ... 0 1 0 1, 0 ... 0 0  
2<sup>31</sup> 2<sup>13</sup> 2<sup>12</sup> 2<sup>11</sup> 2<sup>10</sup> 2<sup>9</sup> 2<sup>1</sup> 2<sup>0</sup>  
2G 8K 4K 2K 1K 512 2 1  
(4K + 1K) · 2<sup>-10</sup> = 5K · 2<sup>-10</sup> = 5.0  
even adding a single bit is not that hard  
however note the ending 5.0000... oof  
0 ... 0 1 0 1, 0 ... 0 1  
2<sup>31</sup> 2<sup>13</sup> 2<sup>12</sup> 2<sup>11</sup> 2<sup>10</sup> 2<sup>9</sup> 2<sup>1</sup> 2<sup>0</sup>  
2G 8K 4K 2K 1K 512 2 1  
(4K + 1K + 1) · 2<sup>-10</sup> = (5K + 1) · 2<sup>-10</sup> = 5.0009765625

<p>This can still be displayed, but:</p> $5 \cdot \frac{1}{1024} = 5.0009765625 < 5.001 < 5.001953125 = 5 \cdot \frac{2}{1024}$ <p>5.001 can't be displayed, only the one on the left and right. Inaccuracy due to integer limitations! only multiples of <math>2^{-10}</math> can be displayed</p>	<p>• <math>z = 0 \mid 11111111111111111111111111111111 = \text{Float.MAX\_VALUE}</math></p> $k = 2^8 - 2 \quad m = 2^{23} - 1$ $W(z) = 2^{28-2-(2^7-1)} \cdot (1 + (2^{23} - 1) \cdot 2^{-23})$ $= 2^{27-1} \cdot (1 + 1 - 2^{-23}) = 2^{127} \cdot 2 \cdot (1 - 2^{-24})$ $= 2^{128} - 2^{104}$ <p>• <math>z = 0 \mid 00000001 \mid 000000000000000000000000 = \text{Float.MIN\_NORMAL}</math></p> $W(z) = 2^{1-(2^7-1)} = 2^{-126}$	<p>Guard digits</p> <ul style="list-style-type: none"> <li>– just as hidden as the hidden bit</li> <li>– implemented in hardware not software</li> <li>– essentially a flag for rounding!</li> </ul> <p>• <math>x = 1.999'999'880'790'710'449'218'750 = 2 - 2^{-23}</math></p> <p>• <math>y = 0.000'000'059'604'644'775'390'625 = 2^{-24}</math></p> <p>• <math>x = 0 \mid 01111111 \mid 111111111111111111111111</math></p> <p>• <math>y = 0 \mid 01101111 \mid 000000000000000000000000</math></p> <p>• <math>x' = 0 \mid 01111111 \mid (1)111111111111111111111111</math></p> <p>• <math>y' = 0 \mid 01111111 \mid (0)000000000000000000000000 \quad a + 24, m \cdot 2^{-24}</math></p> <p>• <math>z' = 0 \mid 01111111 \mid (1)111111111111111111111111 = x' + y'</math></p> <p>• <math>z = 0 \mid 01111111 \mid 111111111111111111111111</math></p> <p>• <math>z = x</math></p> <p>• <b>result is wrong, should be rounded up!!</b></p> <p><math>x = 1.999'999'880'790'710'449'218'750 = 2 - 2^{-23}</math></p> <p><math>y = 0.000'000'059'604'644'775'390'625 = 2^{-24}</math></p> <p><math>x = 0 \mid 01111111 \mid 111111111111111111111111</math></p> <p><math>y = 0 \mid 01101111 \mid 000000000000000000000000</math></p> <p><math>x' = 0 \mid 01111111 \mid (1)111111111111111111111111(0) \quad \text{Guard-Digit}</math></p> <p><math>y' = 0 \mid 01111111 \mid (0)000000000000000000000000(1) \quad \text{Guard-Digit}</math></p> <p><math>z' = 0 \mid 01111111 \mid (1)111111111111111111111111 = x' + y'</math></p> <p><math>z'' = 0 \mid 01111111 \mid (10)000000000000000000000000(0) \quad \text{rounded}</math></p> <p><math>z = 0 \mid 10000000 \mid 000000000000000000000000</math></p>
<p>Biggest number in Klotz</p> $2^{-10} \cdot \sum_{k=0}^{31} 2^k = 2^{-10} \cdot (2^{32} - 1) = 2^{22} - 2^{-10} = 4M - \frac{1}{1K}$ <p>Note that <math>\sum_{k=0}^3 1</math> due to Integer being max 32 bit!</p> <p>Cut Numbers after comma</p> $2^0 \cdot \sum_{k=0}^{n-1} 2^k = 2^n - 1$ <p>Cut Numbers before comma</p> $2^{-n} \cdot \sum_{k=0}^{n-1} 2^k = 2^{-n} \cdot (2^n - 1) = 2^0 - 2^{-n} = 1 - 2^{-n}$ <p>Addition: <math>x \cdot 2^{-a} + y \cdot 2^{-a} = (x + y) \cdot 2^{-a}</math></p> <p>Subtraction: <math>x \cdot 2^{-a} - y \cdot 2^{-a} = (x - y) \cdot 2^{-a}</math></p> <p>2 complement: <math>2^n \cdot 2^{-a} - x \cdot 2^{-a} = (2^n - x) \cdot 2^{-a}</math> (Zweierkomplement)</p> <p>Shiftoptions: <math>(x \cdot 2^{-a}) \cdot 2^b = (x \cdot 2^b) \cdot 2^{-a}</math></p> <ul style="list-style-type: none"> <li>• For multiplication, only the comma position changes</li> </ul> $(x \cdot 2^{-a}) \cdot (y \cdot 2^{-a}) = (x \cdot y) \cdot (2^{-a} \cdot 2^{-a}) = (x \cdot y) \cdot 2^{-2a}$ <ul style="list-style-type: none"> <li>• for <math>a = n \rightarrow x \cdot y &lt; 2^{2n}</math>:</li> </ul> $(x \cdot 2^{-n}) \cdot (y \cdot 2^{-n}) = (x \cdot y) \cdot 2^{-2n} < 2^{2n} \cdot 2^{-2n} = 1$	<p><b>Rounding Modes</b> mathematical operations are done before rounding!</p> <p>Round to nearest:</p> <p>submode 1: <b>Ties to Even:</b> round to an even number, LSB = 0</p> <p>0.1 -&gt; 0, 1.1 -&gt; 2.0, 3.2 -&gt; 2.0</p> <ul style="list-style-type: none"> <li>– cancels itself out over time</li> <li>– better when numbers with 0.1 are added</li> <li>– default variant</li> </ul> <p>submode 2: <b>Ties away from Zero</b></p> <p>Round to next number 1.1 -&gt; 2.0</p> <p>Round toward Zero: truncate -&gt; 1.5325 -&gt; 1.0</p> <p>Round toward +Infinity: Ceiling 1.1 -&gt; 2.0</p> <p>Round toward -Infinity: Floor, 1.9 -&gt; 1.0</p> <p><b>Rounding errors</b></p> <p>• <math>z = 0 \mid 01111111 \mid 000000000000000000000001</math></p> $W(z) = (-1)^0 \cdot 2^{27-1-(2^7-1)} \cdot (1 + 1 \cdot 2^{-23}) = 1 + 2^{-23}$ <p>• Die Zahl <math>z' = 0 \mid 01111111 \mid 000000000000000000000001</math> this has too many bits, we can't display this!! the precision for float is 23 not 24!</p> $W(z') = (-1)^0 \cdot 2^{27-1-(2^7-1)} \cdot (1 + 1 \cdot 2^{-24}) = 1 + 2^{-24}$ <p>⇒ The solution in programming languages is to simply round the number. However this is unprecise... for our number here it would round to <math>z</math>, aka <math>1 + 2^{-23}</math></p> <p>This error can now be quantified: first we display a real number <math>x</math> as a floating point number</p> <p><b><math>W(x) \rightarrow x \in \mathbb{R}</math></b></p> <p>The absolut error <b>E</b> is:</p> <p><b><math>E(x) =  W(x) - x </math></b></p> <p>The relative error <b>e</b> is:</p> <p><b><math>e(x) = \left  \frac{E(x)}{x} \right  = \left  \frac{R(x) - x}{x} \right </math></b></p> <p>For the number from before, an example:</p> $E(z') = W(z) - W(z') = 1 + 2^{-23} - (1 + 2^{-24}) = 2^{-23} - 2^{-24} = 2^{-24}$ $e = \frac{W(z) - W(z')}{W(z')} = \frac{2^{-24}}{1 + 2^{-24}} = \frac{1}{2^{24} + 1} < \frac{1}{2^{24}}$ <p><math>z' = 0 \mid 01111111 \mid 000000000000000000000001</math></p> <p>In this case the number would be rounded down, to 1</p>	<p>compared to before, not a single digit is the same... the guard digit did the trick by flagging <math>z</math> to be rounded after addition.</p> <p><b>Special cases for floating point numbers <math>\pm 0 \rightarrow a = 0</math> ..</b> subnormal number <math>\rightarrow a = 0 \dots 0</math> &amp; <math>z \text{ NOT } 0 \dots 0</math></p> <p><math>\pm \infty \rightarrow a = 1 \dots 1</math> &amp; <math>z = 0 \dots 0</math></p> <p>NaN <math>\rightarrow a = 1 \dots 1</math> &amp; <math>z \text{ NOT } 0 \dots 0</math></p> <p>The reason we want things such as infinity is this:</p> <pre>float x2 = x * x; if (x2 &lt; max_float) {     float y2 = y * y;     if (y2 &lt; max_float) {         float r = x2 + y2;         if (r &lt; max_float) {             r = sqrt(r);         }     } }</pre> <p>We either check for overflows the entire time! Or we just continue calculation with infinity And then at THE END we check if it's infinite</p> <pre>float r = sqrt (x*x + y*y); if (r == Float.POSITIVE_INFINITY) { ... }</pre> <p>Calculation with <math>\infty</math> exchange <math>\infty</math> with <math>x</math> and define: <math>\lim_{x \rightarrow \infty}</math> if the limit exists, the result is the limit if the limit doesn't exist, result is NaN</p> <p><math>y \pm \infty = \pm \infty \gg \lim_{x \rightarrow \infty} y + x = 0</math></p> <p><math>\pm x / 0 = \pm \infty</math> for <math>x \neq 0</math></p> <p><math>y / \pm \infty = 0</math> for <math>y \neq 0 \gg \lim_{x \rightarrow \infty} y / x = 0</math></p> <p><math>\sqrt{\infty} = \infty \gg \lim_{x \rightarrow \infty} \sqrt{x} = \infty</math></p>
<p><b>Floating Point Numbers</b></p> <p>prefix MSB s: 1 bit Significand z: BitString of number Exponent a: float -&gt; 8, double -&gt; 11 Encoded Exponent k: explained below bias b: <math>-(2^{e-1} - 1)</math> precision p: float -&gt; 24, double -&gt; 53</p> <p><math>W((a, z)) = z \cdot 2^a</math></p>	<p><b>Machine Epsilon: <math>\epsilon</math></b> The biggest possible error</p> <p><b><math>\epsilon = 2^{-24} = 2^{-2}</math></b></p> <p>The smallest number <math>\epsilon</math>, so that <math>1 + \epsilon &gt; 1</math></p> <p>Note: the worst error is therefore always the flip of a highest bit</p> <p>In this case the flip of <math>2^{-23}</math> to 1 after <math>2^{-24}</math> was rounded.</p> <p>In Programming languages the epsilon is defined as: the distance from 1 to the next smallest number <math>2^{-23}</math></p> <p>Therefore the biggest error is often a slight bit bigger than we have defined</p> <p>⇒ <math>\epsilon_c = 2^{-p+1} = 2\epsilon</math></p>	<p>Calculation with <b>NaN</b> <b>ALL OPERATION WITH NaN RESULT IN NaN!</b> Other operations might also result in NaN:</p> <p>ZB: <math>-\infty + \infty</math> ZB: <math>\pm \infty \cdot 0</math> ZB: <math>0/0</math> or <math>\pm \infty / \pm \infty</math> ZB: <math>\sqrt{x}</math> for <math>x &lt; 0</math></p> <p>There are also some implementations for exceptions inside floating</p> <p>qNaN: propagates exception sNaN: throws exception</p>
<p><b>Different forms of floating point numbers:</b></p> <ul style="list-style-type: none"> <li>• Single (Java float): 24 bit precision, 8 bit exponent</li> <li>• Double (Java double): 52 bit precision, 11 bit exponent</li> <li>• Quadruple (2008): 113 bit precision, 15 bit exponent</li> <li>• Octuple (2008): not used</li> <li>• (Single-Extended 43 Bit)</li> <li>• (Double-Extended 79 Bit)</li> <li>• (Encodings of base 10 32, 64, 128 Bit, 2008)</li> </ul>	<p><b>Normalized Numbers</b></p> <p><math>0.75 = 0.110 \cdot 2^0 = \boxed{1.100 \cdot 2^{-1}} = 0.111 \cdot 2^1 = 11.00 \cdot 2^{-2} = (0.)0011 \cdot 2^2 = \dots</math></p> <p>The normalized number always has the 1 before the comma!!</p> <p>1.something -&gt; normalized (the one is called the hidden bit) it is always set, a zero would not be detected!</p> <p>The hidden bit is NOT an actual bit!! It is used inside the formula for calculation!</p>	<p>Calculation with <b>NaN</b> <b>ALL OPERATION WITH NaN RESULT IN NaN!</b> Other operations might also result in NaN:</p> <p>ZB: <math>-\infty + \infty</math> ZB: <math>\pm \infty \cdot 0</math> ZB: <math>0/0</math> or <math>\pm \infty / \pm \infty</math> ZB: <math>\sqrt{x}</math> for <math>x &lt; 0</math></p> <p>There are also some implementations for exceptions inside floating</p> <p>qNaN: propagates exception sNaN: throws exception</p>
<p><b>Normalized Numbers</b></p> <p><math>0.75 = 0.110 \cdot 2^0 = \boxed{1.100 \cdot 2^{-1}} = 0.111 \cdot 2^1 = 11.00 \cdot 2^{-2} = (0.)0011 \cdot 2^2 = \dots</math></p> <p>The normalized number always has the 1 before the comma!!</p> <p>1.something -&gt; normalized (the one is called the hidden bit) it is always set, a zero would not be detected!</p> <p>The hidden bit is NOT an actual bit!! It is used inside the formula for calculation!</p>	<p><b>Float and Double in Java</b> double default when using . -&gt; 0.5 optionally defined with d -&gt; 0.5d float -&gt; 0.5f is necessary!</p> <p><b>Addition of floating point numbers</b></p> <ol style="list-style-type: none"> <li>1. check prefix, if different do subtraction</li> <li>2. set prefix for result -&gt; 0 or 1</li> <li>3. if exponents are different, shift smaller number to right</li> </ol> <p>First shift a single 1 -&gt; hidden bit ! Shift the rest with 0s</p> <ol style="list-style-type: none"> <li>4. do addition just like with integers</li> <li>5. (if carry == 1) normalize result</li> </ol> <p>increment exponent by one shift significand to the right by 1</p> <p><math>x = 1.5, y = 0.75 \quad (\text{bit} = \text{hidden bit})</math></p> <p><math>x = 0 \mid 01111111 \mid 100000000000000000000000</math></p> <p><math>y = 0 \mid 01111110 \mid 100000000000000000000000</math></p> <p><math>x' = 0 \mid 01111111 \mid (1)100000000000000000000000</math></p> <p><math>y' = 0 \mid 01111111 \mid (0)110000000000000000000000 \quad a + 1, z \cdot 2^{-1}</math></p> <p><math>z' = 0 \mid 01111111 \mid (10)010000000000000000000000 = x' + y'</math></p> <p><math>z'' = 0 \mid 10000000 \mid (1)001000000000000000000000 \quad a + 1, z \cdot 2^{-1}</math></p> <p><math>z = 0 \mid 10000000 \mid 001000000000000000000000</math></p> <p><math>z = 2.25</math></p>	<p><b>Subnormal Numbers</b> No hidden bit! Used when the exponent is too small</p> <p><b><math>V(x) = (-1)^s \cdot 2^{-(b-1)} \cdot z \cdot 2^{p-1}</math></b></p> <p>Often used as a gradual underflow for semi results during a calculation</p> <p>• <math>y = s \mid 00000000 \mid 111111111111111111111111</math></p> $V(y) = \pm 2^{-126} \cdot (2^{23} - 1) \cdot 2^{-23}$ $= \pm 2^{-126} \cdot (1 - 2^{-23})$ $= \pm 2^{-126} - 2^{-149}$ <p>• <math>y = 0 \mid 00000000 \mid 000000000000000000000001 = \text{Float.MIN\_VALUE}</math></p> $V(y) = 2^{-126} \cdot 2^{-23} = 2^{-149}$
<p><b>Encoding of Exponents</b></p> <p>The bitpatterns 0.0 and 1...1 are not valid exponents</p> <p>The exponent is encoded as <math>k = a - b, b = -(2^{e-1} - 1)</math></p> <p>-&gt; <math>k \in [1, 2^e - 2]</math></p> <p>-&gt; <math>a \in [-(2^{e-1} - 2), 2^{e-1} - 1] = [-(b - 1), b]</math></p> <p>- for float : <math>b = 2^{8-1} - 1 = 127 \wedge a \in [-126, 127]</math></p> <p>- for double : <math>b = 2^{11-1} - 1 = 1023 \wedge a \in [-1022, 1023]</math></p> <p><b>Value of a floating point number</b></p> <p><b><math>W(x) = (-1)^s \cdot 2^{k+b} \cdot (1 + (z \cdot 2^{-(p-1)}))</math></b></p> <p>Example: Note <math>z</math> is the variable here, not the significand</p> <p><math>z = 1 \mid 10000000 \mid 011000000000000000000000</math></p> $s = 1 \quad k = 2^7 \quad m = 2^{21} + 2^{20}$ $W(z) = (-1)^1 \cdot 2^{27+(-2^7-1)} \cdot (1 + (2^{21} + 2^{20}) \cdot 2^{-23})$ $= -2^1 \cdot (1 + (2^{-2} + 2^{-3}))$ $= -(2 + 2^{-1} + 2^{-2})$ $= -2.75$ <p><b>Special Values:</b></p> <p>• <math>z = 0 \mid 01111111 \mid 000000000000000000000000</math></p> $s = 0 \quad k = 2^7 - 1 \quad m = 0$ $W(z) = (-1)^0 \cdot 2^{27-1-(2^7-1)} \cdot (1 + 0 \cdot 2^{-23}) = 2^0 = 1$ <p>• <math>z = s \mid \boxed{kkkk \ kkkk} \mid 000000000000000000000000</math> <math>m = 0</math> m = significand</p> $W(z) = (-1)^s \cdot 2^{k+b} \cdot (1 + 0 \cdot 2^{-23})$ $= \pm 2^{k+b} = \pm 2^a$	<p><b>Integers as Floating Point Numbers</b></p> <p><b><math>G(x) = 2^{k-b} \cdot 1</math></b></p> <p>Regular numbers obviously don't need comma shifting!!</p> <p>16million + 1 is the smallest number that can't be displayed</p> <ul style="list-style-type: none"> <li>- a=23 only integers can be displayed</li> <li>- a=24 only even numbers can be displayed</li> <li>- a=25 only numbers divisible by 4 can be displayed</li> <li>- a=26 only numbers divisible by 8 can be displayed</li> <li>- a ± 23 only numbers divisible by <math>2^a-23</math> can be displayed</li> </ul> <p>check page 2 for floating point comparison!</p>	<p><b>Integers as Floating Point Numbers</b></p> <p><b><math>G(x) = 2^{k-b} \cdot 1</math></b></p> <p>Regular numbers obviously don't need comma shifting!!</p> <p>16million + 1 is the smallest number that can't be displayed</p> <ul style="list-style-type: none"> <li>- a=23 only integers can be displayed</li> <li>- a=24 only even numbers can be displayed</li> <li>- a=25 only numbers divisible by 4 can be displayed</li> <li>- a=26 only numbers divisible by 8 can be displayed</li> <li>- a ± 23 only numbers divisible by <math>2^a-23</math> can be displayed</li> </ul> <p>check page 2 for floating point comparison!</p>