

Exercises from the book

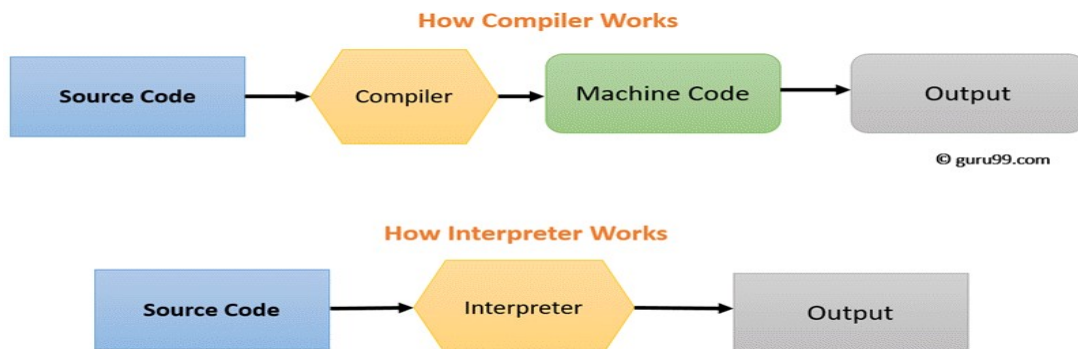
Essentials of computing systems

<https://doi.org/10.21814/uminho.ed.93>

Exerc. 1.1: Describe the following terms with your own words: **(a)** Compiler; **(b)** Interpreter; **(c)** Virtual machine.

a) A compiler is a computer program that takes as input the text of a program written in a high-level language like C, called the “source code” and outputs binary code, or machine code, that can be executed by the computer.

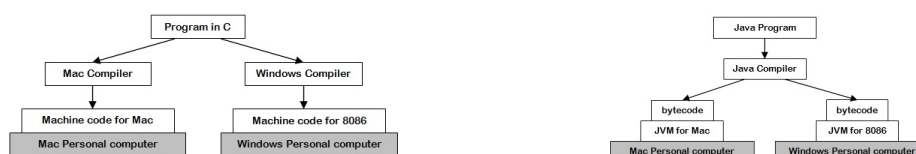
b) An interpreter executes the machine code as it reads and translates the source code (for instance Python). The compiler does the whole translation beforehand and does not execute the program. An interpreter is more flexible but is slower than a compiled program.



See <https://www.guru99.com/difference-compiler-vs-interpreter.html> for a more comprehensive comparison.

c) A Virtual Machine (VM) is a way of being able to run the same compiled program on different architectures and be more efficient than an interpreter. Different hardware has different compilations of source code into machine code. So, to run the same program on different architectures and if we do not want to use an interpreter, we can use a VM specific for that architecture. This is the approach taken by Java and C#. A Java or C# compiler translates the source code into “bytecode”. This is binary code BUT not yet machine code because machine code is specific to the architecture. The bytecode is run by a virtual machine. Each different architecture has its own version of the VM so the SAME bytecode can be translated to the specific machine code.

This may seem like the interpreter, but it is much faster due to the nature of the source code taking much more effort to translate to machine code than the bytecode. Languages that use VMs are faster than interpreted ones and can run the same “compiled program” on different architectures.



From <http://www.edu4java.com/en/concepts/compiler-interpreter-virtual-machine.html>

Exerc. 1.2: Write a small program in a given programming language. Compile it and try to calculate the ratio of source code statements to the machine language instructions generated by the compilation process. Add different types of statements to the high-level program, one at a time, and check how the machine language program is affected.

```
// https://godbolt.org/  
void dummy()  
{  
    // Uncomment the following lines, one by one,  
    // waiting for the compiler to update the output.  
  
    // int a = 2;  
    // int b = 5;  
    // int c = 0;  
    // c = a + b;  
  
    // change type of variable a from int to long  
}
```

Using <https://godbolt.org/> present several versions of a small C function (see below). The source code looks different. How about the machine translation?

```
int factorial(int n) {  
    int fact = 1;  
    for (int i = 1; i <= n; i++) {  
        fact *= i;  
    }  
    return fact;  
}
```

```
int factorial(int n) {  
    if (n <= 0) return 1;  
    else return n * factorial(n-1);  
}
```

Exerc. 1.3: On a **big-endian** computer, a 32-bit integer with value

00010010 00110100 01010110 01111000

is about to be stored in the memory at location 132,104. Indicate which memory cells are affected and which values are stored in each one.

132103	
132104	00010010
132105	00110100
132106	01010110
132107	01111000
132108	

Exerc. 1.4: Consider that part of the memory of a **little-endian** computer contains the values shown in the figure. Indicate the value of a 32-bit integer if it is read from the memory location 4365.

4362	0100 0011
4363	0111 0000
4364	0000 0011
4365	0001 0010
4366	1111 1111
4367	0000 0000
4368	0000 1111

00001111 00000000 111111 00010010

Exerc. 1.5: In a stored-program computer, both the instructions and the data of a program are located in the main memory while it is executed. What are the possible implications if a program accidentally modifies the value that is stored in a memory cell that is related to an instruction?

The program will not behave as expected. It can crash or produce incorrect output. Sometimes this is not an accident. Some viruses have its code encrypted. When loaded, the first instructions decode the main virus program, by modifying memory cells with the encrypted code.

Sometimes this happens because of external factors. See “Soft Error” (not to be confused with software error) in https://en.wikipedia.org/wiki/Soft_error

Exerc. 1.6: In a factory, the production process of a given product goes through four steps: preparation, assembly, testing, and packaging. Those steps take the following times, in seconds, to be executed: preparation (20), assembly (30), testing (35), and packaging (35). Calculate the time needed to produce 1000 replica of the product by:

(a) a single person.

$$120 * 1000 = 120000$$

(b) four persons working in a pipeline.

$$\text{Time to output the first item} = 120$$

$$\text{Output time per item after the first one: } 35$$

$$\text{Total processing time} = 120 + 999 \times 35 = 35085$$

Exerc. 2.1: To encode Roman numbers (from 1 to 899), the following binary encoding for the symbols has been proposed: I (01), V (100), X (00), L (101), C (110), D (111). Indicate whether this encoding is valid and, if so, what Roman number is represented by the binary pattern **111101000101**.

Answer: Yes, it is valid because each symbol has a distinct encoding and the symbols I,V,X,L,C and D are enough to represent the numbers 1 to 899 (DCCCXCIX).

111 101 00 01 01

111	101	00	01	01
D	L	X	I	I

DLXII = 562 (<https://www.calculatorsoup.com/calculators/conversions/roman-numeral-converter.php>)

Exerc. 2.2: Decode the following ASCII string: 1010101 0101110 0100000 1001101 1101001 1101110 1101000 1101111.

Answer:

ASCII table: <https://www.sciencebuddies.org/science-fair-projects/references/ascii-table>

1010101	0101110	0100000	1001101	1101001	1101110	1101000	1101111
U	.		M	i	n	h	o

U. Minho

Exerc. 2.3: A digital image has 128x128 pixels. Each pixel in the image stores information related to three channels (Red, Blue, Green). If each channel is capable of distinguishing 256 different tones, indicate the size in bytes of the image.

Answer:

Number of pixels: $128 \times 128 = 16384$

Each pixel needs 3 bytes (needs to represent 3 numbers ranging from 0 to 255 and to represent a number between 0 and 255 we need 8 bits = 1 byte)

Size of picture = $16384 \times 3 = 49152$ Bytes

Exerc. 2.4: An image occupies 192 kibibytes and has dimensions of 256x512 pixels. Each pixel is represented by three unsigned integer values, which indicate the intensity of each channel (Red-Green-Blue) in that pixel. Indicate, in binary and decimal, the maximum value that can be assigned to each of these integers, if they all have the same size.

Answer:

1 KiB (kibibyte) = 1024 bytes (book page 17)

Size of picture file: $192 \times 1024 = 196608$ bytes

Number of pixels in image: $256 \times 512 = 131072$

Number of bytes available for each pixel: $196608 / 131072 = 1,5$ bytes

One byte = 8 bits. 1,5 bytes = $8 + 4 = 12$ bits

So, for each value of the RGB, we have $12 / 3 = 4$ bits available.

With 4 bits we can go from 0000 to 1111 (binary) \Leftrightarrow 0 to 15 (decimal)

Exerc. 2.5: The CYMK* subtractive colour system is formed by Cyan, Magenta, Yellow and Black and works due to the absorption of light, as the colours that are seen come from the part of the light that is not absorbed. Each pixel is represented by four 6-bit patterns that indicate the intensity in each channel. Indicate how many different colours a pixel can have, assuming that the “00000-” (“000000” and “000001”) patterns cannot be used.

Answer: with 6 bits we have $2^6 = 64$ possible values for each of the components but we cannot use 2 of them (“000000” and “000001”) so we have 62 for each component. The color is obtained from the combination of the four values, so we have $62 \times 62 \times 62 \times 62 = 62^4 = 14776336$ possible colors.

Exerc. 2.6: Consider a digital image, with dimensions 1024x512, that codifies each pixel with six bits.

- (a) Calculate in KiB the size of this image.
(b) The image was proportionally resized and additionally only three bits are now used to represent the information of the pixels. What are the final dimensions of the image if the file occupies 48 KiB?

Answer:

a)

1 KiB (kibibyte) = 1024 bytes

Size of picture file: $1024 \times 512 \times 6 \text{ bits} = 1024 \times 512 \times 6 / 8 \text{ bytes} = 512 \times 6 / 8 \text{ KiB} = 384 \text{ KiB}$

b)

$W \times H \times 3 / 8 / 1024 = 48 \text{ KiB}$

$W / H = 1024 / 512 \Leftrightarrow W / H = 2 \Leftrightarrow W = 2H$

$W \times H \times 3 / 8 / 1024 = 48 \Leftrightarrow 2H \times H \times 3 / 8 / 1024 = 48 \Leftrightarrow 6H^2 = 48 \times 8 \times 1024 \Leftrightarrow$

$H^2 = 65536 \Leftrightarrow H = 256$

$W \times 256 \times 3 / 8 / 1024 = 48 \Leftrightarrow 768W = 393216 \Leftrightarrow W = 512$

Final dimensions: 512x256

Exerc. 2.7: The SCB system for evaluating football players consists of three parameters: Strength, Courage and Braveness. Each parameter is represented by a binary pattern (of 7 bits each) that indicates the respective intensity. Indicate the number of different valid assessments with this system, if the “1111111” and “1111110” patterns represent evaluations that are still unknown or invalid, respectively.

Answer: For each of the three parameters (with 7 bits) we have 2^7 possibilities minus two (“1111111” and “1111110”). So, for each possibility we have $2^7 - 2 = 128 - 2 = 126$ possible assessments. Since we use three dimensions (Strength, Courage and Braveness) the total number of possible assessments will be $126 \times 126 \times 126 = 126^3 = 2000376$.

Exerc. 2.8: Calculate the size in kB of a sound file, if the recording lasts exactly 2 minutes and it is sampled using a sampling rate of 50 kHz and a sample resolution of 8 bits. What is the size in KiB?

Answer:

2 minutes = 120 seconds

50 kHz = 50 cycles/sec * 1000 = 50 000

Number of samples = 120 x 50000 = 6 000 000 samples

Each sample = 8 bits = 1 byte

Size of file = 6 000 000 x 1 byte = 6 000 000 bytes

6 000 000 bytes = 6 000 000 / 1000 Kb = 6000 Kb

6 000 000 bytes = 6 000 000 / 1024 KiB ≈ 5860 KiB (5859,375)

Exerc. 2.9: Calculate the sample resolution of a sound file with 4.5kB, if the recording lasts one minute with a sampling rate of 50 Hz.

Answer:

50 Hz = 50 cycles/sec

Number of samples taken in 1 minute (60 seconds) = 60 x 50 = 3000

Size of one sample = size of file 4500 bytes / 3000 samples = 1,5 bytes = 8+4 bits = 12 bits

Exerc. 2.10: Consider that a typical daily newspaper page contains 3700 Unicode characters including white spaces.

(a) How many bytes are needed to encode a 32-page edition of that daily newspaper, if it includes on average 50 photos (1.2MiB each one)?

(b) How many mebibytes are needed to store all the numbers of the newspaper published in a year?

(c) If a library contains 1250 different daily newspapers, which have on average 50 years of publication, how many pebibytes are stored there?

Answer:

UTF-8 is based on 8-bit code units. Each character is encoded as 1 to 4 bytes. Book page 19: "The majority of the common-use characters fit into the first 64k code patterns, which just require two bytes (16 bits)". We will assume that on average 1 Unicode character equals 2 bytes. So, One page = 3700 characters = 3700 x 2 = 7400 bytes.

a) Encode a 32 page edition:

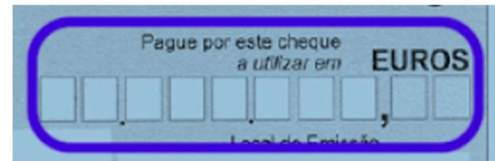
$32 \times 7400 + 50 \times 1,2 \times 1024^2 = 236800 + 62914560 = 63\,151\,360$ bytes

b) From a) one year will be $(365 \times 63151360) / 1024^2$ mebibytes = 23050246400 / 1048576 = 21982 MiB

c) One paper per year = 21982 MiB

50 years x 1250 titles x 21982 MiB = 1373875000 Mib x 1024^2 = 1440612352000000 bytes
 $1440612352000000 / 1024^5 = 1,28$ PiB

Exerc. 3.1: As the figure shows, bank cheques in Portugal have 10 boxes to indicate the amount to be paid. What are the minimum and the maximum values that can be written in a bank cheque?



Answer:

Minimum: 0,01

Maximum: 99999999,99

Exerc. 3.2: Represent the following decimal numbers as binary numbers:
(a) 131; (b) 511; (c) 888; (d) 4096.

Answer:

a) 131				b) 511			
$131_{10} = 10000011_2$				11111111_2			
Divide by the base 2 to get the digits from the remainder				Division by 2			
Division by 2	Quotient	Remainder (Digit)	Bit #	Division by 2	Quotient	Remainder (Digit)	Bit #
$(131)/2$	65	1	0	$(511)/2$	255	1	0
$(65)/2$	32	1	1	$(255)/2$	127	1	1
$(32)/2$	16	0	2	$(127)/2$	63	1	2
$(16)/2$	8	0	3	$(63)/2$	31	1	3
$(8)/2$	4	0	4	$(31)/2$	15	1	4
$(4)/2$	2	0	5	$(15)/2$	7	1	5
$(2)/2$	1	0	6	$(7)/2$	3	1	6
$(1)/2$	0	1	7	$(3)/2$	1	1	7
$= (10000011)_2$				$(1)/2$	0	1	8
				$= (11111111)_2$			

c) 888				d) 4096			
1101111000_2				100000000000_2			
Division by 2				Division by 2			
Division by 2	Quotient	Remainder (Digit)	Bit #	Division by 2	Quotient	Remainder (Digit)	Bit #
$(888)/2$	444	0	0	$(4096)/2$	2048	0	0
$(444)/2$	222	0	1	$(2048)/2$	1024	0	1
$(222)/2$	111	0	2	$(1024)/2$	512	0	2
$(111)/2$	55	1	3	$(512)/2$	256	0	3
$(55)/2$	27	1	4	$(256)/2$	128	0	4
$(27)/2$	13	1	5	$(128)/2$	64	0	5
$(13)/2$	6	1	6	$(64)/2$	32	0	6
$(6)/2$	3	0	7	$(32)/2$	16	0	7
$(3)/2$	1	1	8	$(16)/2$	8	0	8
$(1)/2$	0	1	9	$(8)/2$	4	0	9
$= (1101111000)_2$				$(4)/2$	2	0	10
				$(2)/2$	1	0	11
				$(1)/2$	0	1	12
				$= (100000000000)_2$			

Exerc. 3.3: What is the largest natural number that can be represented with (a) 5, (b) 10, (c) 18, and (d) 32 bits?

Answer:

a) $(11111)_2 = (1 \times 2^4) + (1 \times 2^3) + (1 \times 2^2) + (1 \times 2^1) + (1 \times 2^0) = (31)_{10}$

b) $1111111111_2 =$

$1 \times 2^9 + 1 \times 2^8 + 1 \times 2^7 + 1 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 1023_{10}$

c) $111111111111111111_2 = 1 \times 2^{17} + 1 \times 2^{16} + 1 \times 2^{15} + 1 \times 2^{14} + 1 \times 2^{13} + 1 \times 2^{12} + 1 \times 2^{11} + 1 \times 2^{10} + 1 \times 2^9 + 1 \times 2^8 + 1 \times 2^7 + 1 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 262143_{10}$

d) $11111111111111111111111111111111_2 = 1 \times 2^{31} + 1 \times 2^{30} + 1 \times 2^{29} + 1 \times 2^{28} + 1 \times 2^{27} + 1 \times 2^{26} + 1 \times 2^{25} + 1 \times 2^{24} + 1 \times 2^{23} + 1 \times 2^{22} + 1 \times 2^{21} + 1 \times 2^{20} + 1 \times 2^{19} + 1 \times 2^{18} + 1 \times 2^{17} + 1 \times 2^{16} + 1 \times 2^{15} + 1 \times 2^{14} + 1 \times 2^{13} + 1 \times 2^{12} + 1 \times 2^{11} + 1 \times 2^{10} + 1 \times 2^9 + 1 \times 2^8 + 1 \times 2^7 + 1 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 4294967295_{10}$

Exerc. 3.4: List all the digits and their binary representation in base 13.

Answer:

0	0000 ₂
1	0001 ₂
2	0010 ₂
3	0011 ₂
4	0100 ₂
5	0101 ₂
6	0110 ₂

7	0111 ₂
8	1000 ₂
9	1001 ₂
A	1010 ₂
B	1011 ₂
C	1100 ₂

Exerc. 3.5: Convert the following binary numbers to hexadecimal:

(a) 101111101101; (b) 1001110110;

(c) 11111111111; (d) 10100011110.

Answer:

HINT: Convert every 4 binary digits (from bit0) to hex digit.

(a) $101111101101 = 1011\ 1110\ 1101 = B\ E\ D = BED$

(b) $1001110110 = 10\ 0111\ 0110 = 2\ 7\ 6 = 276$

(c) $11111111111 = 111\ 1111\ 1111 = 7\ F\ F = 7FF$

(d) $10100011110 = 101\ 0001\ 1110 = 5\ 1\ E = 51E$

Exerc. 3.6: Convert the following hexadecimal numbers to binary:
(a) BEEF; **(b)** 1000.FF; **(c)** ABC.DEF; **(d)** DAC.34.

Answer:

(a) BEEF = B E E F = 1011 1110 1110 1111 = 1011111011101111

(b) 1000.FF = 1 0 0 0 . F F = 0001 0000 0000 0000 . 1111 1111 = 1000000000000.11111111₂

(c) ABC.DEF = 1010 1011 1100 . 1101 1110 1111 = 101010111100.110111101111₂

(d) DAC.34 = 1101 1010 1100 . 0011 0100 = 110110101100.00110100₂

Exerc. 3.7: Convert the following decimal numbers to base 5:
(a) 77 **(b)** 131 **(c)** 511 **(d)** 1000.

Answer:

(a) 302₅

Division	Quotient	Remainder (Digit)	Digit #
77/5	15	2	0
15/5	3	0	1
3/5	0	3	2
= (302) ₅			

(b) 1011₅

Division	Quotient	Remainder (Digit)	Digit #
131/5	26	1	0
26/5	5	1	1
5/5	1	0	2
1/5	0	1	3
= (1011) ₅			

(c) 4021₅

511/5	102	1	0
102/5	20	2	1
20/5	4	0	2
4/5	0	4	3
= (4021) ₅			

(d) 13000₅

1000/5	200	0
200/5	40	0
40/5	8	0
8/5	1	3
1/5	0	1
= (13000) ₅		

Exerc. 3.8: Convert the following base-9 numbers to decimal: **(a)** 66; **(b)** 123; **(c)** 317; **(d)** 800.

Answer:

(a) $66_9 = 6 \times 9^1 + 6 \times 9^0 = 60_{10}$

(b) $123_9 = 1 \times 9^2 + 2 \times 9^1 + 3 \times 9^0 = 102_{10}$

(c) $317_9 = 3 \times 9^2 + 1 \times 9^1 + 7 \times 9^0 = 259_{10}$

(d) $800_9 = 8 \times 9^2 + 0 \times 9^1 + 0 \times 9^0 = 648_{10}$

Exerc. 3.9: Convert the following numbers from the given base to the indicated bases:

- (a) 66_{10} to bases 2, 7 and 9
 (b) $13F.4_{16}$ to bases 10 and 12
 (c) 1110010.1_2 to bases 3, 4 and 7
 (d) AB_{13} to bases 2, 6, and 8.

Answer:

a) 66_{10}

1000010 ₂			123 ₇			73 ₉		
66/2	33	0	Division	Quotient	Remainder (Digit)	66/9	7	3
33/2	16	1	66/7	9	3	7/9	0	7
16/2	8	0	9/7	1	2	= (73) ₉		
8/2	4	0	1/7	0	1			
4/2	2	0	= (123) ₇					
2/2	1	0						
1/2	0	1						
= (1000010) ₂								

3.9) b) $13F.4_{16}$

$$13F.4_{16} = 319.25_{10}$$

$$13F.4_{16} = 1 \times 16^2 + 3 \times 16^1 + 15 \times 16^0 + 4 \times 16^{-1} = 319.25_{10}$$

<https://www.rapidtables.com/convert/number/base-converter.html>

How to convert from any base to any base

- Convert from source base to decimal (base 10) by multiplying each digit with the base raised to the power of the digit number (starting from right digit number 0):

$$\text{decimal} = \sum (\text{digit} \times \text{base}^{\text{digit number}})$$
- Convert from decimal to destination base: divide the decimal with the base until the quotient is 0 and calculate the remainder each time. The destination base digits are the calculated remainders.

$$13F.4_{16} = 227.3_{12}$$

Base 16 to decimal calculation: See left.

Decimal to base 12 calculation:

Multiply the number with the destination base raised to the power of decimals of the result (number of decimal places -1) up to 6 digits resolution: $\text{floor}(319.25 \times 12^1) = 3831$

Divide by the base to get the digits from the remainders:

Division	Quotient	Remainder (Digit)	Digit #
3831/12	319	3	0
319/12	26	7	1
26/12	2	2	2
2/12	0	2	3
$= (227.3)_{12}$			

3.9) (c) 1110010.1₂ to bases 3, 4 and 7

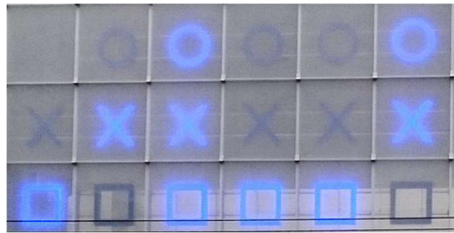
11020.111111 ₃	1302.2 ₄	222.333333 ₇																																																																											
<p>Base 2 to decimal calculation: $1110010.1_2 = 1 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 + 1 \times 2^{-1} = 114.5_{10}$</p> <p>Decimal to base 3 calculation: Multiply the number with the destination base raised to the power of decimals of the result up to 6 digits resolution: $\text{floor}(114.5 \times 3^6) = 83470$</p> <p>Divide by the base to get the digits from the remainders:</p>	<p>Base 2 to decimal calculation: $1110010.1_2 = 1 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 + 1 \times 2^{-1} = 114.5_{10}$</p> <p>Decimal to base 4 calculation: Multiply the number with the destination base raised to the power of decimals of the result up to 6 digits resolution: $\text{floor}(114.5 \times 4^1) = 458$</p> <p>Divide by the base to get the digits from the remainders:</p>	<p>Base 2 to decimal calculation: $1110010.1_2 = 1 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 + 1 \times 2^{-1} = 114.5_{10}$</p> <p>Decimal to base 7 calculation: Multiply the number with the destination base raised to the power of decimals of the result up to 6 digits resolution: $\text{floor}(114.5 \times 7^6) = 13470810$</p> <p>Divide by the base to get the digits from the remainders:</p>																																																																											
<table> <tr><td>83470/3</td><td>27823</td><td>1</td></tr> <tr><td>27823/3</td><td>9274</td><td>1</td></tr> <tr><td>9274/3</td><td>3091</td><td>1</td></tr> <tr><td>3091/3</td><td>1030</td><td>1</td></tr> <tr><td>1030/3</td><td>343</td><td>1</td></tr> <tr><td>343/3</td><td>114</td><td>1</td></tr> <tr><td>114/3</td><td>38</td><td>0</td></tr> <tr><td>38/3</td><td>12</td><td>2</td></tr> <tr><td>12/3</td><td>4</td><td>0</td></tr> <tr><td>4/3</td><td>1</td><td>1</td></tr> <tr><td>1/3</td><td>0</td><td>1</td></tr> </table>	83470/3	27823	1	27823/3	9274	1	9274/3	3091	1	3091/3	1030	1	1030/3	343	1	343/3	114	1	114/3	38	0	38/3	12	2	12/3	4	0	4/3	1	1	1/3	0	1	<table> <tr><td>458/4</td><td>114</td><td>2</td></tr> <tr><td>114/4</td><td>28</td><td>2</td></tr> <tr><td>28/4</td><td>7</td><td>0</td></tr> <tr><td>7/4</td><td>1</td><td>3</td></tr> <tr><td>1/4</td><td>0</td><td>1</td></tr> </table> <p>= (1302.2)₄</p>	458/4	114	2	114/4	28	2	28/4	7	0	7/4	1	3	1/4	0	1	<table> <tr><td>13470810/7</td><td>1924401</td><td>3</td></tr> <tr><td>1924401/7</td><td>274914</td><td>3</td></tr> <tr><td>274914/7</td><td>39273</td><td>3</td></tr> <tr><td>39273/7</td><td>5610</td><td>3</td></tr> <tr><td>5610/7</td><td>801</td><td>3</td></tr> <tr><td>801/7</td><td>114</td><td>3</td></tr> <tr><td>114/7</td><td>16</td><td>2</td></tr> <tr><td>16/7</td><td>2</td><td>2</td></tr> <tr><td>2/7</td><td>0</td><td>2</td></tr> </table>	13470810/7	1924401	3	1924401/7	274914	3	274914/7	39273	3	39273/7	5610	3	5610/7	801	3	801/7	114	3	114/7	16	2	16/7	2	2	2/7	0	2
83470/3	27823	1																																																																											
27823/3	9274	1																																																																											
9274/3	3091	1																																																																											
3091/3	1030	1																																																																											
1030/3	343	1																																																																											
343/3	114	1																																																																											
114/3	38	0																																																																											
38/3	12	2																																																																											
12/3	4	0																																																																											
4/3	1	1																																																																											
1/3	0	1																																																																											
458/4	114	2																																																																											
114/4	28	2																																																																											
28/4	7	0																																																																											
7/4	1	3																																																																											
1/4	0	1																																																																											
13470810/7	1924401	3																																																																											
1924401/7	274914	3																																																																											
274914/7	39273	3																																																																											
39273/7	5610	3																																																																											
5610/7	801	3																																																																											
801/7	114	3																																																																											
114/7	16	2																																																																											
16/7	2	2																																																																											
2/7	0	2																																																																											

3.9) (d) AB7₁₃ to bases 2, 6, and 8.

11100110000 ₂	12304 ₆	3460 ₈																																																												
<p>Base 13 to decimal calculation: $AB7_{13} = 10 \times 13^2 + 11 \times 13^1 + 7 \times 13^0 = 1840_{10}$</p> <p>Decimal to base 2 calculation: Divide by the base to get the digits from the remainders:</p>	<p>Base 13 to decimal calculation: $AB7_{13} = 10 \times 13^2 + 11 \times 13^1 + 7 \times 13^0 = 1840_{10}$</p> <p>Decimal to base 6 calculation: Divide by the base to get the digits from the remainders:</p>	<p>Base 13 to decimal calculation: $AB7_{13} = 10 \times 13^2 + 11 \times 13^1 + 7 \times 13^0 = 1840_{10}$</p> <p>Decimal to base 8 calculation: Divide by the base to get the digits from the remainders:</p>																																																												
<table> <tr><td>1840/2</td><td>920</td><td>0</td></tr> <tr><td>920/2</td><td>460</td><td>0</td></tr> <tr><td>460/2</td><td>230</td><td>0</td></tr> <tr><td>230/2</td><td>115</td><td>0</td></tr> <tr><td>115/2</td><td>57</td><td>1</td></tr> <tr><td>57/2</td><td>28</td><td>1</td></tr> <tr><td>28/2</td><td>14</td><td>0</td></tr> <tr><td>14/2</td><td>7</td><td>0</td></tr> <tr><td>7/2</td><td>3</td><td>1</td></tr> <tr><td>3/2</td><td>1</td><td>1</td></tr> <tr><td>1/2</td><td>0</td><td>1</td></tr> </table>	1840/2	920	0	920/2	460	0	460/2	230	0	230/2	115	0	115/2	57	1	57/2	28	1	28/2	14	0	14/2	7	0	7/2	3	1	3/2	1	1	1/2	0	1	<table> <tr><td>1840/6</td><td>306</td><td>4</td></tr> <tr><td>306/6</td><td>51</td><td>0</td></tr> <tr><td>51/6</td><td>8</td><td>3</td></tr> <tr><td>8/6</td><td>1</td><td>2</td></tr> <tr><td>1/6</td><td>0</td><td>1</td></tr> </table>	1840/6	306	4	306/6	51	0	51/6	8	3	8/6	1	2	1/6	0	1	<table> <tr><td>1840/8</td><td>230</td><td>0</td></tr> <tr><td>230/8</td><td>28</td><td>6</td></tr> <tr><td>28/8</td><td>3</td><td>4</td></tr> <tr><td>3/8</td><td>0</td><td>3</td></tr> </table>	1840/8	230	0	230/8	28	6	28/8	3	4	3/8	0	3
1840/2	920	0																																																												
920/2	460	0																																																												
460/2	230	0																																																												
230/2	115	0																																																												
115/2	57	1																																																												
57/2	28	1																																																												
28/2	14	0																																																												
14/2	7	0																																																												
7/2	3	1																																																												
3/2	1	1																																																												
1/2	0	1																																																												
1840/6	306	4																																																												
306/6	51	0																																																												
51/6	8	3																																																												
8/6	1	2																																																												
1/6	0	1																																																												
1840/8	230	0																																																												
230/8	28	6																																																												
28/8	3	4																																																												
3/8	0	3																																																												

Exerc. 3.10: The St. Gallen train station, in Switzerland, is equipped with a binary electronic clock to indicate the time of day (hours, minutes, seconds) on three lines. At what time the photograph was taken?

Answer:



Hours	01001	9
Minutes	011001	25
Seconds	101110	46

Exerc. 3.11: A given computer is equipped with 1,073,741,824 bytes of memory. Why did this odd number was chosen?

Answer:

Let's see if 1073741824 is a power of 2. If $\log_2(n)$ is integer than n is a power of 2, else not.
 $\log_2(1073741824) = 30$ so $1073741824 = 2^{30}$

This number was chosen probably because this computer uses 30 bits to address 2^{30} memory cells. With 30 bits we can address each memory byte individually.

Exerc. 3.12: A 32-bit signed integer on a little-endian computer contains the numerical value of 3. If it is transmitted to a big-endian computer byte by byte and stored there, with byte 0 in byte 0, byte 1 in byte 1, and so on, what is its numerical value on the big-endian machine if read as a 32-bit unsigned integer?

Answer:

Little Endian Byte Order: **The least significant byte** (the "little end") of the data **is placed at the byte with the lowest address**. The rest of the data is placed in order in the next bytes in memory.

Big Endian Byte Order: **The most significant byte** (the "big end") of the **data is placed at the byte with the lowest address**. The rest of the data is placed in order in the next bytes in memory.

Number 3 as a 32 bit signed integer on the the little endian computer will be interpreted like:
 $3_{10} = 11_2 = 00000000\ 00000000\ 00000000\ 00000011_2$

How this number is stored CORRECTLY on both architectures:

Little Endian		Big Endian	
Memory position X	00000011	Memory position X	00000000
X + 1	00000000	X + 1	00000000
X + 2	00000000	X + 2	00000000
X + 3	00000000	X + 3	00000011

Transmitting from little endian to big endian and storing the bytes using the same order is the same as reversing the interpretation of the bits. Big endian stored number will be interpreted like this:

Little Endian		Big Endian without changing order	
Memory position X	00000011	Memory position X	00000011
X + 1	00000000	X + 1	00000000
X + 2	00000000	X + 2	00000000
X + 3	00000000	X + 3	00000000

Or, writing from **most significant byte to least**, according to memory position, we get:
 $00000011\ 00000000\ 00000000\ 00000000_2 = 50331648_{10} = 300\ 0000_{16}$
 Very different from 3.

Exerc. 3.13: As of 2018, Iceland had about 23,000 registered footballers (male and female). Calculate the minimum number of bits that allows representing this value with an integer encoded in signal and amplitude.

Answer:

1 bit for signal

How many bits for amplitude? We have to represent 23000 different athletes.

$\log_2(23000)$ bits for amplitude = 14.489 Since it must be an integer number, we need 15 bits for amplitude. To represent with signal and amplitude we need $1 + 15 = 16$ bits

Exerc. 3.14: Monaco had 37,831 inhabitants in 2013. Calculate the minimum number of bits that are required to encode this value as a signed integer. What is the answer if the value is encoded as an unsigned integer?

Answer:

To represent 37831 in binary we need $\text{Log}_2(37831) = 15.2 \Rightarrow 16$ bits

As a signed integer: 1 for sign + 16 for amplitude = 17

Unsigned: 0 for sign + 16 for amplitude = 16

Exerc. 3.15: A given company has 19 employees, who are paid every two weeks. It is necessary to register the number of half hours that each employee worked in each workday (Monday to Friday). For health reasons, the law does not permit an employee to work more than 12h in a day. Indicate the minimum number of bits needed to represent this information for two weeks.

Answer:

Assuming we just need to record the work done by each employee every day. It is implicit that we already have a data structure for each employee.

Maximum work time per day = 12h

Each employee can work from 0 to 24 half hours per day.

To record a number from 0 to 24 we need $\text{Log}_2(24) = 4.5 \Rightarrow 5$ bits

Two work weeks = 10 days. Total bits = $19 \times 10 \times 5 = \mathbf{950 \text{ bits}}$

Assuming we record the work done by each employee every day and need to record the employee id

To identify each employee, we need to record a number from 0 to 18 $\text{Log}_2(19) = 4.2 \Rightarrow 5$ bits

Maximum work time per day = 12h

Each employee can work from 0 to 24 half hours per day.

To record a number from 0 to 24 we need $\text{Log}_2(24) = 4.5 \Rightarrow 5$ bits

To record this information for 19 employees for 10 days (2 weeks) we need:

$19 \times 10 \times 10 = \mathbf{1900 \text{ bits}}$

Assuming we just want to record the total number of half hours done in 2 weeks per employee

To identify each employee, we need to record a number from 0 to 18 $\text{Log}_2(19) = 4.2 \Rightarrow 5$ bits

Maximum number of half hours worked in 2 weeks per employee: $24 \times 10 = 240$

$\text{Log}_2(240) = 7.9 \Rightarrow 8$ bits

To process the two-week salary, we need:

$19 \times (5 + 8) = \mathbf{247 \text{ bits}}$

Exerc. 3.16: An European institute aims to assign a code to each of its members. To this end, it was decided to use the format AA / HHHH-BB, with A being a capital letter, H being a hexadecimal digit and B being a base 2 digit. The two letters indicate which of the 51 affiliated countries the member belongs to (e.g., BE for Belgium, PO for Portugal, LX for Luxembourg). The binary digits are used to encode the type of membership of the member with the Association (00: junior member, 01: regular member, 10: senior member). Indicate, for a given country, the maximum number of members that this code allows to register.

Answer:

For each country we can have 0 up to $FFFF_{16}$ members in each one of the 3 classes (junior, regular and senior). In total we can address/identify $3 \times FFFF$ members. $\log_2(FFFF) = 20$.
Maximum number of registered members for one country: 3×2^{20}

Exerc. 3.17: In 2014, the Psy's Gangnam Style video reached 2,147,483,648 views on YouTube. However, the number presented was negative. Explain why this happened and suggest the simplest solution to overcome it. To solve this issue, YouTube made an internal change that now allows the counters to go up to 9,223,372,036,854,775,807.

Answer:

$$2147483648_{10} = 10000000\ 00000000\ 00000000\ 00000000_2$$

This is a 32 bit number. If youtube showed a negative value we know that the representation was a 32 bit signed integer. The new maximum (9223372036854775807) tells us that they changed to a 64 bit signed integer because $\log_2(9223372036854775807) = 63$.

If they used an unsigned 64 bit integer the maximum would be:

$$2^{64} = 18\ 446\ 744\ 073\ 709\ 551\ 616$$

Exerc. 3.18: Find the value of X, so that: **(a)** $23_x = 10101_2$; **(b)** $4X_7 = 35_9$.

Answer:

(a) $23_x = 10101_2 = 21_{10}$

$$2X^1 + 3X^0 = 21 \Leftrightarrow 2X + 3 = 21 \Leftrightarrow 2X = 18 \Leftrightarrow X = 9$$

Brute force:

From $23_x = 21_{10}$, we know that the last digit in base X is 3. To find X we divide 21 by X and the first remainder must be 3 (and the second 2). Let's try $X = 9$:

Division	Quotient	Remainder (Digit)	Digit #
21/9	2	3	0
2/9	0	2	1
= $(23)_9$			

(b) $4X_7 = 35_9 = 32_{10}$ Since the base is 7 we just need to convert 32_{10} to base 7: 44_7

$$X = 4 \text{ Or } 4x7^1 + Xx7^0 = 32 \Leftrightarrow x = 32 - 28 = 4$$

Exerc. 3.19: Add the following natural numbers:

(a) $110011_2 + 10101_2$; **(b)** $129B_{12} + 239_{12}$; **(c)** $CBA_{16} + 987_{16}$.

Answer:

Natural numbers = 1, 2, 3, 4, 5, ... so we know they are all positive.

(a) $110011_2 + 10101_2 = 1001000_2 = 72_{10}$

(b) $129B_{12} + 239_{12} = (2135 + 333)_{10} = 2468_{10} = 1518_{12}$

(c) $CBA_{16} + 987_{16} = (3258 + 2439)_{10} = 5697_{10} = 1641_{16}$

Exerc. 3.20: Indicate the ten's-complement of the following decimal numbers:
(a) 1236; (b) 90037; (c) 111122.

Answer:

The "Ten's-complement" is the number we add to make 10.
For a 4 digit number is the number we add to have 10000

(a) 1236 Ten's-complement = $10000 - 1236 = 8764$

(b) 90037 Ten's-complement = $100000 - 90037 = 9963$

(c) 111122 Ten's-complement = $1000000 - 111122 = 888878$

Exerc. 3.21: Indicate the two's-complement of the following binary numbers:
(a) 0011100₂; (b) 110011001₂; (c) 00000001₂; (d) 1110000001₂.

Answer:

How to calculate two's-complement for a binary number?

1. Find the one's complement by inverting 0s & 1s of a given binary number.
2. Add 1 to the one's complement to get the two's complement.

	Binary number	one's-complement	two's-complement
(a)	0011100	1100011	1100100
(b)	1100110011	0011001100	0011001101
(c)	00000001	11111110	11111111
(d)	11100000001	00011111110	00011111111

Exerc. 3.22: Write the 8-bit sign-magnitude, one's-complement, two's-complement representations for decimal numbers: (a) +18; (b) +121; (c) -33; (d) -100.

For the **8-bit sign-magnitude** we first convert the number to binary and then set the first bit to 0 if the number is positive or 1 if negative.

One's-complement: A negative number needs to be converted to its complement (by flipping 0s to 1s and vice versa), which should have a '1' in the MSB. **Positive numbers, which have a '0' in the MSB, are used as is, i.e., they are not converted to their complements.**

Two's complement is just one's complement incremented by 1. To find the two's complement of a binary number, one just needs to flip bits and add 1. Again, positive numbers, which have a '0' in the MSB, are used as is, i.e., they are not converted to their complements.

	8-bit sign-magnitude	one's-complement	two's-complement
(a) +18	00010010	00010010	00010010
(b) +121	01111001	01111001	01111001
(c) -33	10100001	-(00100001) = 11011110	11011110 + 1 = 11011111
(d) -100	11100100	-(01100100) = 10011011	10011011 + 1 = 10011100

Exerc. 3.23: Calculate the value of the 10-bit binary number $10110\ 00111_2$ in the following representations:

(a) sign-magnitude; (b) one's-complement; (c) two's-complement. (d) excess-511.

Answer:

(a) sign-magnitude – the first bit tells us the number is negative. The remaining bits are the magnitude. $011000111 = 199$ so, **Sign-magnitude $10110\ 00111_2 = -199$**

(b) one's-complement. If this number is in one's-complement form, we know it is negative (from the most significant bit.)

Revert from one-complement $1011000111 \rightarrow 0100111000 \rightarrow 312$ but we know it is negative so **$10110\ 00111_2 = -312$**

(c) two's-complement. To revert from two's complement we subtract 1 and then flip the bits:
 $1011000111 \rightarrow 1011000110 \rightarrow 0100111001_2 = 313_{10}$ but we know it is a negative number.
 So, $1011000111_2 = -313_{10}$

(d) excess-511. In an excess-b representation, an n-bit pattern, whose unsigned integer value is V ($0 \leq V < 2^n$) represents the signed integer V-b, where b is the bias (or offset) of the numeral system. The representable numeric values range from -b to 2^n-1-b . In this case b = 511.

V = unsigned value of the bit pattern. X is the represented value.

$X = V - b$

Our number $1011000111_2 = 711_{10}$ We know that this bit pattern represents V which in turn represents the number $X = V - b$

$V - b = X \Leftrightarrow 711 - 511 = X \Leftrightarrow \mathbf{X = 200}$. In excess-511, the number $10110\ 00111_2 = \mathbf{200}_{10}$

Exerc. 3.24: Represent the number -233 in the following 10-bit representations:

(a) sign-magnitude; (b) one's-complement; (c) two's-complement; (d) excess-511.

Answer:

(a) sign-magnitude $-233_{10} = 1011101001_2$
 $(1110\ 1001_2 = 233_{10})$

(b) one's-complement with 10 bits:

$-233_{10} = -0011101001_2 \rightarrow \mathbf{1100010110_2}$

(c) two's-complement

$-233_{10} \rightarrow$ one's complement $1100010110_2 \rightarrow$ two's complement (add 1) **1100010111_2**

(d) Represent -233 in excess-511 with 10-bit representation.

$X = V - b$

V = unsigned value of the bit pattern. X is the represented value and b is the excess value

$X = -233$

$b = 511$

$X = V - b \Leftrightarrow -233 = V - b \Leftrightarrow -233 = V - 511 \Leftrightarrow V = -233 + 511$

$-233 + 511 = 278_{10} = 0100010110_2$

Exerc. 3.25: Perform binary subtraction by taking the two's-complement of the subtrahend:
(a) $100110_2 - 111_2$; **(b)** $100110_2 - 10000_2$; **(c)** $1010101_2 - 11_2$; **(d)** $1000001_2 - 1000000_2$.

Answer:

1. In the first step, find the 2's complement of the subtrahend.
2. Add the complement number with the minuend.
3. If we get the carry by adding both numbers, we discard this carry and the result is positive else take 2's complement of the result which will be negative.

(a) $100110 - 111 \Leftrightarrow 100110 - 000111 \Leftrightarrow 100110 + 111001 = [1]011111 = 011111 = 1111$
 Two's complement of $000111 = 111000 + 1 = 111001$

(b) $100110 - 10000 \Leftrightarrow 100110 - 010000 \Leftrightarrow 100110 + 110000 = [1]010110 = 010110$
 Two's complement of $010000 = 101111 + 1 = 110000$

(c) $1010101 - 11 \Leftrightarrow 1010101 + 1111101 = [1]010010 = 010010$

(d) $1000001 - 1000000 \Leftrightarrow 1000001 + 1000000 = [1]000001 = 000001$
 (complement of $1000000 \rightarrow 0111111 + 1 \rightarrow 1000000$ then add and drop carry bit)

Exerc. 3.26: Add the following pairs of unsigned binary numbers, explicitly indicating the carries:

(a) $\begin{array}{r} 11010 \\ + 1010 \\ \hline \end{array}$ **(b)** $\begin{array}{r} 111010 \\ + 101010 \\ \hline \end{array}$ **(c)** $\begin{array}{r} 1001111010 \\ + 1011010 \\ \hline \end{array}$ **(d)** $\begin{array}{r} 1101011 \\ + 1011000 \\ \hline \end{array}$

Answer:

(a)	(b)	(c)	(d)
$\begin{array}{r} \textcolor{red}{111} \\ 11010 \\ + 1010 \\ \hline 100100 \end{array}$	$\begin{array}{r} \textcolor{red}{1111} \\ 111010 \\ + 101010 \\ \hline 1100100 \end{array}$	$\begin{array}{r} \textcolor{red}{11111} \\ 1001111010 \\ + 1011010 \\ \hline 1011010100 \end{array}$	$\begin{array}{r} \textcolor{red}{1111} \\ 1101011 \\ + 1011000 \\ \hline 11000011 \end{array}$

Exerc. 3.27: Add the following pairs of **8-bit two's-complement** numbers, explicitly indicating situations of overflow:

(a) 1001 1010 + 1000 1010 <hr style="width: 100%;"/>	(b) 0111 1010 + 0110 1010 <hr style="width: 100%;"/>	(c) 1101 1101 + 1110 1101 <hr style="width: 100%;"/>	(d) 0110 1011 + 0101 1000 <hr style="width: 100%;"/>
---	---	---	---

Answer:

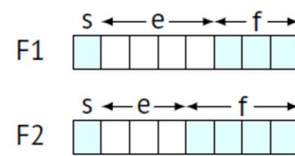
From page 35 of class manual: “an addition overflows whenever the signs of the addends are the same (both numbers are either positive or negative) and the sign of the sum is different from the addends’ sign.” In other words:

“In two’s complement, If the sum of two positive numbers yields a negative result, the sum has overflowed. If the sum of two negative numbers yields a positive result, the sum has overflowed. Otherwise, the sum has not overflowed.”

(a)	(b)	(c)	(d)
<div style="border: 1px solid red; padding: 5px; margin-bottom: 10px;"> 10011010 + 10001010 <hr style="width: 100%;"/> 100100100 </div> <p>With 8 bits we have Overflow: adding 2 negatives we have a positive</p>	<div style="border: 1px solid red; padding: 5px; margin-bottom: 10px;"> 01111010 + 01101010 <hr style="width: 100%;"/> 11100100 </div> <p>Overflow: adding 2 positive returns a negative</p>	<div style="border: 1px solid red; padding: 5px; margin-bottom: 10px;"> 11011101 + 11101101 <hr style="width: 100%;"/> 111001010 </div> <p>Result: 11001010 Dropped last carry</p>	<div style="border: 1px solid red; padding: 5px; margin-bottom: 10px;"> 01101011 + 01011000 <hr style="width: 100%;"/> 11000011 </div> <p>Overflow: adding 2 positive returns a negative</p>

Exerc. 3.28

Exerc. 3.28: Consider two floating-point formats F1 and F2, with 8 bits, based on all the principles presented in Section 3.6, namely normal numbers, subnormal numbers, special values, etc.



- Indicate the mathematical expressions that can be used to calculate the normal numbers in both formats.
- For each format, indicate the bit patterns and the respective decimal value for i) the smallest positive subnormal number, ii) the largest subnormal number, iii) the smallest positive normal number, iv) one, and the v) largest normal number.
- Calculate the decimal values of the following bit patterns for the F1 format: i) 10110011, ii) 01111010, iii) 10010001, iv) 00000011, v) 11000001.
- Represent in the F1 format, the following values: i) -111.01_3 , ii) 128_{10} , iii) 111.01_{10} , iv) $-18C_{16}$, v) 0.005_8 .
- Convert the following numbers represented in the F1 format into the F2 format: i) 00110011, ii) 11101001, iii) 00010000, iv) 11001110, v) 10000010. Overflow must be represented by $\pm\infty$, underflow by ± 0 and the roundings must be made to the closest value.

Answers:

The exponent is encoded in an excess format. The bias value is a number near the middle of the range of possible values that is selected to represent zero. The bias typically equals $2^{k-1}-1$, where k is the number of bits in the exponent. The actual exponent is found by subtracting the bias from the stored exponent.

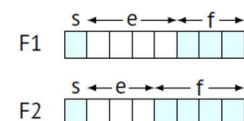
The mantissa is always normalized (1.xyz...) so it always starts with a 1. Because of that we can "ignore" that 1 in the representation. This bit is also the only one that is to the left of the binary point. So, the only part of the mantissa that needs to be represented in the bit pattern is its fractional part.

(a) Indicate the mathematical expressions that can be used to calculate the normal numbers in both formats.

Equation 3.6 page 37:

$$V = (-1)^s \times (1 + f) \times 2^{e-b}$$

The bias typically equals $2^{k-1}-1$, where k is the number of bits in the exponent.



General formula for **normal numbers**:

$$V = -1^s \times (1.0 + 0.M) \times 2^{e-bias}$$

$$V = -1^s \times (1.0 \times 2^0 + bit \times 2^{-1} + bit \times 2^{-2} + bit \times 2^{-3} + \dots) \times 2^{e-bias}$$

General formula for **subnormal numbers**:

$$V = -1^s \times (Mantissa) \times 2^{e-bias}$$

$$V = -1^s \times (bit \times 2^0 + bit \times 2^{-1} + bit \times 2^{-2} + bit \times 2^{-3} + \dots) \times 2^{e-bias}$$

For subnormal numbers, one does not assume the leading 1, hence it is directly used as it is represented. This means the range of the mantissa for subnormals is $[0,1]$, unlike normal numbers where it's $[1,2]$ due to the implicit leading 1.

$$F1: b = 2^{k-1}-1 = 2^3-1 = 7$$

$$V = (-1)^s \times (1 + f) \times 2^{e-7}$$

$$F2: b = 2^{k-1}-1 = 2^2-1 = 3$$

$$V = (-1)^s \times (1 + f) \times 2^{e-3}$$

3.28 (b) For each format, indicate the bit patterns and the respective decimal value for...

Book page 38:

"The all-zeros exponent is reserved to represent subnormal numbers and zero. A subnormal number (or denormalised number) is a non-zero number with magnitude smaller than the smallest positive normal number. Its exponent value is fixed to be 1-bias and the mantissa M is restricted by the condition $0 \leq M < 1$ (there is no leading 1). So, for the all-zeros e exponent, the value V of a subnormal number is given by the following equation: $V = (-1)^S \times f \times 2^{e-b}$

A floating-point number may be recognized as subnormal whenever its exponent is the least value possible. For the mantissa the interpretation is that if the exponent is non-minimal, there is an implicit leading 1, and if the exponent is minimal, there isn't, and the number is subnormal.

when converting a number to binary excess format, offset is added to the original number and when retrieving original number, it's subtracted.

3.28 (b) i) the smallest positive subnormal number

0 for the signal, 0000 for the exponent (smallest possible) and the smallest possible mantissa with 3 digits and not being zero: 001

0	0	0	0	0	0	0	1
---	---	---	---	---	---	---	---

$$V = (-1)^0 \times 2^{-2} \times 2^{0-7} = 2^{-9} = 1 / 512$$

NOTE: In the mantissa we write 2^{-2} instead of 2^{-3} because this is a subnormal number.

For F2

0	0	0	0	0	0	0	1
---	---	---	---	---	---	---	---

$$V = (-1)^0 \times 2^{-3} \times 2^{0-3} = 2^{-6} = 1 / 64$$

3.28 (b) ii) the largest subnormal number

0 for the signal, 0000 for the exponent (smallest possible) and the largest possible mantissa with 3 digits: 111

0	0	0	0	0	1	1	1
---	---	---	---	---	---	---	---

$$V = (-1)^0 \times (2^0 + 2^{-1} + 2^{-2}) \times 2^{0-7} = (4 \times 2^{-2} + 2 \times 2^{-2} + 1 \times 2^{-2}) \times 2^{-7} = (4 \times 2^{-9} + 2 \times 2^{-9} + 1 \times 2^{-9}) = 7 \times 2^{-9}$$

$$V = 7 / 512$$

0	0	0	0	1	1	1	1
---	---	---	---	---	---	---	---

$$V = (-1)^0 \times (2^0 + 2^{-1} + 2^{-2} + 2^{-3}) \times 2^{0-3} = (8 \times 2^{-3} + 4 \times 2^{-3} + 2 \times 2^{-3} + 1 \times 2^{-3}) \times 2^{-3} = (8 \times 2^{-6} + 4 \times 2^{-6} + 2 \times 2^{-6} + 1 \times 2^{-6}) = 15 \times 2^{-6}$$

$$V = 15 / 64$$

3.28 (b) iii) the smallest positive normal number

0 for the signal, 0001 for the exponent (smallest possible and normal) and the smallest possible mantissa with 3 digits: 000

0	0	0	0	1	0	0	0
---	---	---	---	---	---	---	---

For normal number: $-1^S \times (1.M) \times 2^{e-bias}$

$$\text{Exponent} = e - bias = 1 - 7 = -6$$

$$V = (-1)^0 \times (1 + 0) \times 2^{1-7} = 2^{-6}$$

$$V = 1 / 64$$

0	0	0	1	0	0	0	0
---	---	---	---	---	---	---	---

$$\text{Exponent} = e - bias = 1 - 3 = -2$$

$$V = (-1)^0 \times (1 + 0) \times 2^{1-3} = 2^{-2}$$

$$V = 1 / 4$$

3.28 (b) iv) One

We know that the mantissa (for normal numbers) starts always with 1 (not represented). So, for the number 1, the mantissa bits will all be zero.

Now the exponent. We know the offset is 7 (see (a)). We want to store the number 0 in the exponent. So, we add $0 + 7 = 7 = 0111_2$. If we stored 0 directly it would be a subnormal number!

0	0	1	1	1	0	0	0
---	---	---	---	---	---	---	---

For normal number: $-1^S \times (1.0 + 0.M) \times 2^{e-bias}$

$$V = (-1)^0 \times (1 + 0) \times 2^{7-7} = 1 \times 1 \times 1$$

$$V = 1$$

For F2 we just need to represent the exponent. The bias is 3 (see (a)) So exponent = $0 + 3 = 011_2$

0	0	1	1	0	0	0	0
---	---	---	---	---	---	---	---

$$V = (-1)^0 \times (1 + 0) \times 2^{3-3} = 1 \times 1 \times 1$$

$$V = 1$$

3.28 (b) v) the largest normal number.

We will use the largest possible value for the mantissa (all 1s) and the largest for the exponent. Exponents of all zeros or all ones are reserved for subnormal numbers or special values. Largest exponent with 4 bits = 1110 (we cannot use 1111 and 0 on any other position would be smaller)

0	1	1	1	0	1	1	1
---	---	---	---	---	---	---	---

For normal number: $-1^S \times (1.0 + 0.M) \times 2^{e-bias}$

Exponent = $1110_2 = 14$; F1 bias = 7

$$V = (-1)^0 \times (1 + 2^{-1} + 2^{-2} + 2^{-3}) \times 2^{14-7} = (1 + 1/2 + 1/4 + 1/8) \times 2^7 = 1.875 \times 128 = 240$$

$$V = 240$$

0	1	1	0	1	1	1	1
---	---	---	---	---	---	---	---

Exponent = $110_2 = 6$; F2 bias = 3

$$V = (-1)^0 \times (1 + 2^{-1} + 2^{-2} + 2^{-3} + 2^{-4}) \times 2^{6-3} = (1 + 1/2 + 1/4 + 1/8 + 1/16) \times 2^3 = 1.9375 \times 8 = 15.5$$

3.28 (c) Calculate the decimal values of the following bit patterns for the F1 format: Exponent 4 bits, Mantissa 3 bits

3.28 (c) i) 10110011

1	0	1	1	0	0	1	1
---	---	---	---	---	---	---	---

$$V = (-1)^1 \times (1 + 0 \times 2^{-1} + 2^{-2} + 2^{-3}) \times 2^{6-7} = -1 \times (1.375) \times \frac{1}{2} = -0.687500$$

3.28 (c) ii) 01111010

0	1	1	1	1	0	1	0
---	---	---	---	---	---	---	---

NaN

All exponent bits are 1 and mantissa not zero. Table 3.4, page 39.

3.28 (c) iii) 10010001

1	0	0	1	0	0	0	1
---	---	---	---	---	---	---	---

$$V = (-1)^1 \times (1 + 2^{-3}) \times 2^{2-7} = -1 \times (1.125) \times 2^{-5} = -0.035156$$

3.28 (c) iv) 00000011

0	0	0	0	0	0	1	1
---	---	---	---	---	---	---	---

$$\text{SUBNORMAL } V = (-1)^0 \times (2^{-1} + 2^{-2}) \times 2^{0-7} = 0.75 \times 2^{-7} = 0.005859$$

3.28 (c) v) 11000001

1	1	0	0	0	0	0	1
---	---	---	---	---	---	---	---

$$V = (-1)^1 \times (1 + 2^{-3}) \times 2^{8-7} = -1 \times (1.125) \times 2 = -2,25$$

(d) Represent in the F1 format, the following values: i) -111.01_3 , ii) 128_{10} , iii) 111.01_{10} , iv) $-18C_{16}$, v) 0.005_8 .

3.28 (d) F1 format: Exponent 4 bits, Mantissa 3 bits

Answers:

3.28 (d) i) -111.01_3

Base 3 to decimal:

$$(-111.01)_3 = -[(1 \times 3^2) + (1 \times 3^1) + (1 \times 3^0) + (0 \times 3^{-1}) + (1 \times 3^{-2})] = -13.11111_{10}$$

To convert fraction to binary, start with the fraction in question and multiply it by 2 keeping notice of the resulting integer and fractional part. Continue multiplying by 2 until you get a resulting fractional part equal to zero. Then just write out the integer parts from the results of each multiplication.

$$0.11111 \times 2 = 0 + 0.22222$$

$$0.22222 \times 2 = 0 + 0.44444$$

$$0.44444 \times 2 = 0 + 0.88888$$

$$0.88888 \times 2 = 1 + 0.77776$$

$$0.77776 \times 2 = 1 + 0.55552$$

$$0.55552 \times 2 = 1 + 0.111$$

$$0.111 \times 2 = 0 + 0.222$$

$$0.222 \times 2 = 0 + 0.444$$

$$0.444 \times 2 = 0 + 0.888$$

$$0.888 \times 2 = 1 + 0.776$$

$$0.776 \times 2 = 1 + 0.55$$

....

Until we reach a form where the fraction is zero, something like $0.abc \times 2 = 1 + 0$

The 1s and 0s on the right side of the equal sign are collected to form our binary fraction. In our example we have: **0.00011100011** ...

So:

$$111.01_3 = -1101.00011100011010101_2$$

We must encode this binary number into F1 format: Exponent 4 bits, Mantissa 3 bits

Normalize:

$$-1101.00011100011010101_2 = -1.10100011100011010101_2 \times 2^3$$

Signal: $-1 \Rightarrow$ signal bit = 1

Our exponent = $3 + \text{bias} = 3 + 7 = 10 = 1010_2$

Our mantissa = **10100011100011010101**₂

1	1	0	1	0	1	0	1
---	---	---	---	---	---	---	---

3.28 (d) ii) 128_{10}

$$128_{10} = 10000000_2 = 1.000 \times 2^7$$

Signal bit = 0 (positive number)

Our exponent = $7 + 7 = 14 = 1110$

Our mantissa = 0000

0	1	1	1	0	0	0	0
---	---	---	---	---	---	---	---

3.28 (d) iii) 111.01_{10}

$$111.01_{10} = 1101111.0000001010001111011_2 = 1.101111000000101 \times 2^6$$

Signal bit = 0 (positive number)

$$\text{Exponent} = 6 + 7 = 13 = 1101$$

Mantissa = 1.101

0	1	1	0	1	1	0	1
---	---	---	---	---	---	---	---

3.28 (d) iv) $-18C_{16}$

F1 format: 4 bits exponent, 3 mantissa, bias 7

$$-18C_{16} = -0001\ 1000\ 1100 = -110001100 = -1.10001100 \times 2^8$$

Signal bit = 1 (negative number)

F1 exponent = $8 + 7 = 15 = 1111_2$ **OVERFLOW** since max in excess-7 = 14

Because the value is less than the smaller negative integer (biggest in absolute value) we should represent -infinity. By definition $s=1$; $e=\text{all } 1$ and $m=\text{all } 0$ (table 3.4 from the book)

1	1	1	1	1	0	0	0
---	---	---	---	---	---	---	---

Value = -infinity

3.28 (d) v) 0.005_8

$$0.005_8 = 0.005_8 = 0 \times 8^0 + 0 \times 8^{-1} + 0 \times 8^{-2} + 5 \times 8^{-3} = 0.009765625_{10}$$

Convert decimal 0.009765625 to binary:

$$0.009765625 \times 2 = 0 + 0.01953125$$

$$0.01953125 \times 2 = 0 + 0.0390625$$

$$0.0390625 \times 2 = 0 + 0.078125$$

$$0.078125 \times 2 = 0 + 0.15625$$

$$0.15625 \times 2 = 0 + 0.3125$$

$$0.3125 \times 2 = 0 + 0.625$$

$$0.625 \times 2 = 1 + 0.25$$

$$0.25 \times 2 = 0 + 0.5$$

$$0.5 \times 2 = 1 + 0$$

$$0.009765625_{10} = 0.000000101_2 = 1.01 \times 2^{-7} \text{ (normalized)}$$

Signal bit = 0 (positive)

Exponent = $-7 + 7 = 0000 \Rightarrow \text{SUBNORMAL}$

Since this is a subnormal number, we will use all mantissa bits

Mantissa = 101

0	0	0	0	0	1	0	1
---	---	---	---	---	---	---	---

Converting 00000101 to decimal, considering signal, 4 exponent digits and 3 mantissa digits we will get: 0.009766 (the original was 0.009765625)

3.28 (e) Convert the following numbers represented in the F1 format into the F2 format:

i) 00110011, ii) 11101001, iii) 00010000, iv) 11001110, v) 10000010.

Answers:

F1 format: Exponent 4 bits, Mantissa 3 bits

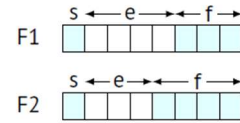
F2 format: Exponent 3 bits, Mantissa 4 bits

$$F1: b = 2^{k-1} - 1 = 2^3 - 1 = 7$$

$$V = (-1)^s \times (1 + f) \times 2^{e-7}$$

$$F2: b = 2^{k-1} - 1 = 2^2 - 1 = 3$$

$$V = (-1)^s \times (1 + f) \times 2^{e-3}$$



3.28 (e) i) 0 0110 011

Exponent bits = 0110 Exponent value (decimal) = 6

$$V_{F1} = (-1)^0 \times (1 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3}) \times 2^{6-7} = 1.375 \times 1/2 = 0.6875 \text{ (decimal)}$$

Convert to binary. To avoid the decimal separator, multiply the decimal number with the base raised to the power of decimals (4 in this example) in result: $0.6875 \times 2^4 = 11$ Faster than the method we used in exercise 3.28 (d)

$11_{10} = 01011_2$ To get the original value (with decimal point) we reverse the process:

$$01011_2 \times 2^{-4} = 0.1011_2 \text{ Normalizing: } = 0.1011_2 = 1.011 \times 2^{-1}$$

F2 Exponent value = $-1 + F2_bias = -1 + 3 = 2 = 10_2 = 010_2$ (F2 demands 3 bits for exponent)

Mantissa = 4-0110 (we need 4 bits for F2 mantissa)

$$F1 \ 0 \ 0110 \ 011 = F2 \ 0 \ 010 \ 0110$$

3.28 (e) ii) 1 1101 001

Exponent = 1101 = 13 decimal

$$V_{F1} = (-1)^1 \times (1 + 0 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3}) \times 2^{13-7} = -1 \times (1.125) \times 64 = -72$$

$$72 = 1001000 = 1.001000 \times 2^6$$

$$F2 \text{ exponent} = 6 + F2_bias = 6 + 3 = 9 = 1001_2$$

We cannot represent the exponent 1001_2 with 3 bits.

Since the number is negative, we store -infinity. All exponent bits as 1 and the mantissa as zero:

11110000 (table 3.4 from the book)

3.28 (e) iii) 0 0010 000

Exponent = 0010 = 2 decimal

$$V_{F1} = (-1)^0 \times (1 + 0 \times 2^{-1} + 0 \times 2^{-2} + 0 \times 2^{-3}) \times 2^{2-7} = 2^{-5} = 0.03125$$

Convert to binary. To avoid the decimal separator, multiply the decimal number with the base raised to the power of decimals (5 in this example) $0.03125 \times 2^5 = 1$

$$1_{10} = 1_2 \text{ Reverse the multiplication we did: } 1_2 \times 2^{-5} = 0.00001_2$$

$$\text{Normalizing} = 0.00001_2 = 1_2 \times 2^{-5} = 1.0_2 \times 2^{-5}$$

F2 exponent = exponent + F2_bias = -5 + 3 = -2 In excess representation we can't have negative numbers. Lets try to **represent it as a subnormal number: $V = -1^s \times f \times 2^{1-\text{bias}}$**

By definition, for subnormal numbers, $e=1-\text{bias} \Leftrightarrow e = 1-3 = -2$ (for F2 format)

Move the decimal point so our exponent = -2 \Rightarrow Our number = $1.0_2 \times 2^{-5} = 0.010_2 \times 2^{-2}$

Answer: 0 000 0010 (exponent = 000, mantissa = 0010 in subnormal we take all bits)

3.28 (e) iv) 1 1001 110

$$V_{F1} = (-1)^1 \times (1 + 1 \times 2^{-1} + 1 \times 2^{-2} + 0 \times 2^{-3}) \times 2^{9-7} = -1 \times 1.75 \times 2 = -3.5$$

$$7_{10} = 111_2 \text{ Normalizing } 1.11 \times 2^2$$

$$\text{F2 exponent} = 2 + \text{F2_bias} = 2+3 = 5 = 101_2$$

Mantissa 1.11

F2 format: 1 101 1100

3.28 (e) v) 1 0000 010

Exponent bits = 0000 (subnormal number). Book page 38: *A subnormal number (or denormalised number) is a non-zero number with magnitude smaller than the smallest positive normal number.*

Its exponent value is fixed to be 1-bias [...] Value = $(-1)^s \times f \times 2^{1-\text{bias}}$

Subnormal number \Rightarrow exponent = 000

F1 bias = 7

$$V_{F1} (\text{subnormal}) = (-1)^1 \times (0 \times 2^0 + 1 \times 2^{-1} + 0 \times 2^{-2}) \times 2^{1-7} = - (0.5) \times 2^{-6} = -0.5/64 = -1/128 = -1 \times 2^{-7}$$

$$V_{F1} (\text{subnormal}) = -1 \times 2^{-7} = 0.0078125_{10} = -0.0000001_2 =$$

To store a subnormal number in F2 (bias = 3):

Signal = 1 (negative number)

$$\text{F2 exponent} = 000 \text{ (3 bits all zero)} \Leftrightarrow e+3=0 \Leftrightarrow e = -3$$

$$\text{To make } e = -3 \text{ we rewrite the value: } V_{F2} = -0.0000001_2 = -0.0000001_2 \times 2^0 = -0.0001_2 \times 2^{-3}$$

Mantissa = 0000 4 most significant bits. Since this is a subnormal number we do not assume the form $(1+f)$ for the mantissa. The value for subnormal numbers is Value = $(-1)^s \times f \times 2^{1-\text{bias}}$

F2 format: 1 000 0000

Exerc. 3.29: Write a C program that calculates the decimal value for a bit pattern that represents a floating-point number. The inputs, provided through the command line, are: the bit pattern (sequence of non-separated 0s and 1s), the number of bits of the exponent e , and the number of bits of the mantissa f . If no pattern is provided, the program lists a pair (bit pattern, decimal value) for all possible 2^{1+e+f} bit patterns.

```

/*
   Exercise 3.29 of the book https://doi.org/10.21814/uminho.ed.93
*/
#include <stdio.h>
#include <stdint.h>
#include <stdlib.h>
#include <math.h>
#include <string.h>

void convert2decimal(char bitPattern[], int nBitsExponent, int nBitsMantissa)
{
    if ((1 + nBitsExponent + nBitsMantissa) != strlen(bitPattern))
    {
        printf("Invalid bit pattern length: %d != %d",
               (int)strlen(bitPattern), 1 + nBitsExponent + nBitsMantissa);
        return;
    }

    int signBit = (int)(bitPattern[0] - '0');

    char* exponent = (char *)malloc(nBitsExponent + 1);
    memcpy(exponent, bitPattern + 1, nBitsExponent);
    exponent[nBitsExponent] = 0;

    char* mantissa = (char *)malloc(nBitsMantissa + 1);
    memcpy(mantissa, bitPattern + 1 + nBitsExponent, nBitsMantissa);
    mantissa[nBitsMantissa] = 0;

    int exponentValue = 0;
    for (int i = 0; i < nBitsExponent; i++)
        exponentValue += (int)(exponent[i] - '0') * pow(2, nBitsExponent - 1 - i);

    int count_ones = 0; int count_zeros = 0;
    for (int i = 0; i < nBitsExponent; i++)
    {
        if (exponent[i] == '1')
            count_ones++; // exponent all ones represents infinity or NaN
        else // table 3.4 https://doi.org/10.21814/uminho.ed.93
            count_zeros++; // exponent all zeros represent subnormal numbers and zero
    }

    int mantissaNonZeroBitCount = 0;
    for (int i = 0; i < nBitsMantissa; i++) mantissaNonZeroBitCount += (int)(mantissa[i] - '0');

    if(count_ones == nBitsExponent)
    {
        printf("%s\t", bitPattern);
        if (mantissaNonZeroBitCount == 0)
            printf("%sInfinity\n", (signBit == 1) ? "-" : "+");
        else
            printf("NaN\n");
        return;
    }

    // In normal numbers the first bit of the mantissa is  $1 \times 2^0 = 1$ 
    double mantissaValue = (count_zeros != nBitsExponent) ? 1 : 0;
    for (int i = 0; i < nBitsMantissa; i++)
    {
        if(count_zeros != nBitsExponent) // normal number
            mantissaValue += (int)(mantissa[i] - '0') * pow(2, -1 - i);
        else
            mantissaValue += (int)(mantissa[i] - '0') * pow(2, - i);
    }

    int bias = (int)pow(2, nBitsExponent - 1) - 1;
    double valueBase10 = pow(2, exponentValue - bias) * mantissaValue;
    if (signBit == 1) valueBase10 = -valueBase10;
}

```

```

printf("%s\t%f\t%s\n", bitPattern, valueBase10,
(count_zeros == nBitsExponent && mantissaNonZeroBitCount != 0) ? "(Subnormal number)" : "");
}

void printSeries(int nBitsExponent, int nBitsMantissa)
{
    int nBits = 1 + nBitsExponent + nBitsMantissa; // sign + exponent bits + mantissa bits
    char bitPattern[nBits + 1]; // +1 for null terminator
    bitPattern[nBits] = 0; // null terminator at end of pattern

    unsigned long long maxValue = (unsigned long long)pow(2, nBits-1) - 1;
    for (int sign = 0; sign <= 1; sign++)
        for (unsigned long long iterator = 0; iterator <= maxValue; iterator++)
        {
            for (int i = 0; i < nBits; i++)
            {
                bitPattern[nBits-i-1] = (iterator & (1 << i)) ? '1' : '0';
            }
            bitPattern[0] = (char)(sign + '0');
            convert2decimal(bitPattern, nBitsExponent, nBitsMantissa);
        }
}

int main(void)
{
    char bitPattern[256];
    int nBitsExponent; int nBitsMantissa; int go = 1;
    char c;
    while (go)
    {
        printf("Enter number of bits for exponent: ");
        scanf("%d", &nBitsExponent);

        printf("Enter number of bits for mantissa: ");
        scanf("%d", &nBitsMantissa);

        printf("Enter bit pattern. Enter * to list all possible values: ");
        scanf("%s", bitPattern);

        if (bitPattern[0] == '*')
            printSeries(nBitsExponent, nBitsMantissa);
        else
            convert2decimal(bitPattern, nBitsExponent, nBitsMantissa);

        printf("\nAnother? (y/n): ");
        scanf(" %c", &c);
        go = (c == 'y');
    }
    return 0;
}

```

Some useful links:

- <https://ncalculators.com/digital-computation/1s-2s-complement-calculator.htm>
- <https://www.rapidtables.com/convert/number/decimal-to-binary.html>
- <https://www.rapidtables.com/convert/number/binary-to-decimal.html>
- <https://www.rapidtables.com/convert/number/binary-to-hex.html>
- <https://www.rapidtables.com/convert/number/hex-to-binary.html>
- <https://www.rapidtables.com/convert/number/base-converter.html>
- <https://projects.klickagent.ch/prozessorsimulation/?converter=true>
- <https://www.omnicalculator.com/math/binary-subtraction> (with steps)
- <https://www.h-schmidt.net/FloatConverter/IEEE754.html>
- <https://trekhele.dev/blog/2021/binary-floating-point/> (only for normal numbers)
- <https://indepth.dev/posts/1019/the-simple-math-behind-decimal-binary-conversion-algorithms>