

2a. Codes and number systems (continued)

How to get the binary representation of an integer:

special case of application of the inverse Horner scheme

⇒

repeated (integer) division by two.

Example:

What is the binary representation of 22 ?

$$22 \text{ div } 2 = 11, \text{ rest } 0$$

$$11 \text{ div } 2 = 5, \text{ rest } 1$$

$$5 \text{ div } 2 = 2, \text{ rest } 1$$

$$2 \text{ div } 2 = 1, \text{ rest } 0$$

$$1 \text{ div } 2 = 0, \text{ rest } 1$$

$$\Rightarrow 22_{10} = 10110_2$$

In case of fractions (numbers with decimal point):

- treat the fractional part separately

- apply *repeated multiplication by 2* (with integer results and fractional rest) to it

Example:

What is the binary representation of 0.375 ?

$$0.375 * 2 \rightarrow 0, \text{ rest } 0.75$$

$$0.75 * 2 \rightarrow 1, \text{ rest } 0.5$$

$$0.5 * 2 \rightarrow 1, \text{ rest } 0$$

$$\Rightarrow 0.375_{10} = 0.011_2$$

$$\Rightarrow 22.375_{10} = 10110.011_2$$

Integer representation in finite cells ("words" with fixed length):

Computer memory: organized in **finite cells**. Typically: Multiples of a byte.

How to store numbers in a 4-byte cell? Some encoding necessary. 2^{32} different values can be represented.

Example: $0 \dots 2^{32} - 1$ can be represented as binary numbers.

Example including negative numbers: $-2^{31} \dots 2^{31} - 1$ can be represented as two's complements numbers.

Two's complement: Most used representation for integers from range $-2^{n-1} \dots 2^{n-1} - 1$ (with n -bit cell).

Non-negative numbers: Are represented simply as binary numbers. Using n bits, the highest bit is always 0.

Negative numbers: (a) Represent their absolute value as binary number, (b) then invert all bits (including the infinite number of leading zeros, resulting in an infinite number of leading ones), and (c) add a 1. The last n bits are the two's complement of the value to be represented.

Example for the "Two's complement":

8-bit two's complement representation of -77

1. Represent $+77$ as a binary number: 1001101
2. Invert all bits, including the leading 0s: $\dots 1110110010$
3. Add 1: $\dots 1110110011$
4. Use only the lowest (= rightmost) 8 bits: 10110011

Notice:

For 16-bit cells, the result would be 1111111110110011.

decimal system	8-bit two's complement
-128	1000 0000
-127	1000 0001
-126	1000 0010
...	...
-2	1111 1110
-1	1111 1111
0	0000 0000
1	0000 0001
...	...
126	0111 1110
127	0111 1111

Properties of the two's complement:

Code represents numbers $-2^{n-1} \dots 2^{n-1} - 1$.

High bit represents **sign**.

Minimal value represented by 1000..., maximal by 0111....

-1 represented by 111....

Floating-point representations

Built analogously to the "scientific representation" of numbers in the form $m * 10^e$

- but using the binary system:

Represent numbers in the form

$$s * m * 2^e$$

with sign s (+1 or -1), non-negative mantissa m , and integer exponent e .

Representation is **normalized** if $1 \leq m < 2$.

Finite number of bits for sign, mantissa and exponent; often used: 32 bits (single precision), 64 bits (double precision), 80 bits (extended precision)

Typical layout of 32-bit floating point number:

Bit 31: represents s (1: negative; 0: positive)

Bits 30..23 (8 bits): represent e : Binary representation of $e + 127$, which allows the values $-126 \dots 127$. Value 0 is used in representation of number 0 and of unnormalized numbers. Value 255_{10} used to represent infinity and other exceptional values.

Bits 22..0 (23 bits): represent m , by binary representation of the integer part of $m * 2^{23}$, without the leading 1.

Example: representing $+26.625$ as a 32-bit normalized floating point number: $26.625_{10} = 11010.101_2$. Normalizing yields $1.1010'1010_2 * 2^4$. 32-bit floating point number ($s=0, e=131_{10}$):

0'10000011'101010100000000000000000

Digital representation of text

based on representation of *letters*

- depending on the alphabet: certain number of bits necessary
- for 26 letters: at least 5 bits necessary
($2^4 = 16 < 26$, $2^5 = 32 > 26$)
- but also encoding of digits, special signs, upper- and lower-case letters... desirable

traditional 7-bit code:

ASCII (= *American Standard Code for Information Interchange*)

ISO-646 norm

later extended to 8-bit code

examples: 00110000 = hex 30 = 48_{10} = digit 0

00110001 = hex 31 = 49_{10} = digit 1

...

00111010 = hex 3A = 58_{10} = ':'

...

01000001 = hex 41 = 65_{10} = 'A'

01000010 = hex 42 = 66_{10} = 'B'

...

01100001 = hex 61 = 97_{10} = 'a'

...

ASCII Table:

Non-printable characters						Printable characters					
Dez	Okt	Hex	Char	Code	Remark	Dez	Okt	Hex	Char	Remark	
0	000	0x00	Ctrl-@	NUL	Null prompt	32	040	0x20		blank	
1	001	0x01	Ctrl-A	SOH	Start of heading	33	041	0x21	!	exclamation mark	
2	002	0x02	Ctrl-B	STX	Start of text	34	042	0x22	"	quotation mark	
3	003	0x03	Ctrl-C	ETX	End of Text	35	043	0x23	#		
4	004	0x04	Ctrl-D	EOT	End of transmission	36	044	0x24	\$	Dollar character	
5	005	0x05	Ctrl-E	ENQ	Enquiry	37	045	0x25	%		
6	006	0x06	Ctrl-F	ACK	Acknowledge	38	046	0x26	&		
7	007	0x07	Ctrl-G	BEL	Bell	39	047	0x27	'	apostroph	
8	010	0x08	Ctrl-H	BS	Backspace	40	050	0x28	(
9	011	0x09	Ctrl-I	HT	Horizontal tab	41	051	0x29)		
10	012	0x0A	Ctrl-J	LF	Line feed	42	052	0x2A	*	asterisk	
11	013	0x0B	Ctrl-K	VT	Vertical tab	43	053	0x2B	+	plus sign	
12	014	0x0C	Ctrl-L	FF	Form feed	44	054	0x2C	,	comma	
				NP	New page	45	055	0x2D	-	minus sign	
13	015	0x0D	Ctrl-M	CR	Carriage return	46	056	0x2E	.	dot	
14	016	0x0E	Ctrl-N	SO	Shift out	47	057	0x2F	/	slash	
15	017	0x0F	Ctrl-O	SI	Shift in	48	060	0x30	0		
16	020	0x10	Ctrl-P	DLE	Data link escape	49	061	0x31	1		
17	021	0x11	Ctrl-Q	DC1	X-ON	50	062	0x32	2		
18	022	0x12	Ctrl-R	DC2		51	063	0x33	3		
19	023	0x13	Ctrl-S	DC3	X-Off	52	064	0x34	4		
20	024	0x14	Ctrl-T	DC4		53	065	0x35	5		
21	025	0x15	Ctrl-U	NAK	No achnowledge	54	066	0x36	6		
22	026	0x16	Ctrl-V	SYN	Synchronous idle	55	067	0x37	7		
23	027	0x17	Ctrl-W	ETB	End transmission blocks	56	070	0x38	8		
24	030	0x18	Ctrl-X	CAN	Cancel	57	071	0x39	9		
25	031	0x19	Ctrl-Y	EM	End of medium	58	072	0x3A	:	colon	
26	032	0x1A	Ctrl-Z	SUB	Substitute	59	073	0x3B	;	semicolon	
27	033	0x1B	Ctrl-[ESC	Escape	60	074	0x3C	<	less than	
28	034	0x1C	Ctrl-\	FS	File separator	61	075	0x3D	=	euquality character	
29	035	0x1D	Ctrl-]	GS	Group separator	62	076	0x3E	>	greater than	
30	036	0x1E	Ctrl-^	RS	Record separator	63	077	0x3F	?	interrogation mark	
31	027	0x1F	Ctrl- <u> </u>	US	Unit separator	64	0100	0x40	@	at	
127	0177	0x7F		DEL	Delete or rubout	65	0101	0x41	A		
						66	0102	0x42	B		
						67	0103	0x43	C		
						68	0104	0x44	D		
						69	0105	0x45	E		
						70	0106	0x46	F		
						71	0107	0x47	G		
						72	0110	0x48	H		
						73	0111	0x49	I		
						74	0112	0x4A	J		
						75	0113	0x4B	K		
						76	0114	0x4C	L		
						77	0115	0x4D	M		
						78	0116	0x4E	N		

79	0117	0x4F	O	
80	0120	0x50	P	
81	0121	0x51	Q	
82	0122	0x52	R	
83	0123	0x53	S	
84	0124	0x54	T	
85	0125	0x55	U	
86	0126	0x56	V	
87	0127	0x57	W	
88	0130	0x58	X	
89	0131	0x59	Y	
90	0132	0x5A	Z	
91	0133	0x5B	[
92	0134	0x5C	\	backslash
93	0135	0x5D]	
94	0136	0x5E	^	caret
95	0137	0x5F	_	low line
96	0140	0x60	`	back quote
97	0141	0x61	a	
98	0142	0x62	b	
99	0143	0x63	c	
100	0144	0x64	d	
101	0145	0x65	e	
102	0146	0x66	f	
103	0147	0x67	g	
104	0150	0x68	h	
105	0151	0x69	i	
106	0152	0x6A	j	
107	0153	0x6B	k	
108	0154	0x6C	l	
109	0155	0x6D	m	
110	0156	0x6E	n	
111	0157	0x6F	o	
112	0160	0x70	p	
113	0161	0x71	q	
114	0162	0x72	r	
115	0163	0x73	s	
116	0164	0x74	t	
117	0165	0x75	u	
118	0166	0x76	v	
119	0167	0x77	w	
120	0170	0x78	x	
121	0171	0x79	y	
122	0172	0x7A	z	
123	0173	0x7B	{	
124	0174	0x7C		
125	0175	0x7D	}	
126	0176	0x7E	~	

ASCII not sufficient for alphabets of the non-Anglo-american world (not even for European alphabets with ä, ö, ü, ß, é, ø, ñ, å...)

Unicode:

2 byte (= 16 bit) code for multilingual text processing
- can represent 65536 characters

amongst them: 27786 Chinese-Japanese-Korean characters

11172 Hangul characters (Korean)

ancient Nordic runes

Tibetan characters

Cherokee characters ...

complete list see <http://www.unicode.org/charts/>

Unicode "Escape sequence" (to utilise it in the programming language Java):

e.g., `\u0041` = 'A' (0041 = hexadecimal representation)

Some characters occur more frequently in texts than others:

better use *variable-length* code

UTF-8: Universal Transformation Format

Characters encoded with variable number of bytes

⇒ for texts with many ASCII characters (like on many web pages) shorter as Unicode

Strings (or *words*): sequences of characters

encoded by sequences of the corresponding code words

Digital representation of pictures:

Gray levels: encode each gray level by a number from a fixed interval (e.g. 0, 1, ..., 255: 8-bit representation – 0 = black, 255 = white)

Colours:

several colour models possible

the most frequently used one:

RGB model

(red / green / blue: primary colours for additive colour composition)

Each colour from a certain range ("gamut") can be mixed from these primary colours

examples with 8-bit intensities:

black	(0, 0, 0)
white	(255, 255, 255)
medium gray	(127, 127, 127)
red	(255, 0, 0)
green	(0, 255, 0)
blue	(0, 0, 255)
light blue	(127, 127, 255)
yellow	(255, 255, 0)

Pictures:

typically represented as raster images – rectangular array (matrix) of *pixels*, each pixel represented by its 3 colour values.

Representation of text documents (book pages, web pages...):

Level of representation is important.

(1) Is there text on the page? – One bit.

(2) What is the text on the page? –
Representation of letter sequence (e.g., string of ASCII characters).

(3) What is the exact layout of the text on the page? –
"formatted text"
- use special characters for formatting, or
- represent the page by a rasterized black-and-white image.

Text documents with graphical elements:

- represent all as a single raster image, or
- use combined representation: several data files, one for the text, the other for the pictorial parts
→ HTML web pages are built like this

file <name>.html or <name>.htm contains text,
layout information and links to other pages
files <name>.gif or <name>.jpg or <name>.png
contain images

Messages and redundancy

Message: A finite sequence of letters, used to transfer some information via encoding/transfer/decoding

Signal: The physical representation of the message (examples: as voltage pattern or light pattern)

Redundancy: Part of a message which is not necessary for the transferred information (later explained more exactly)

Error correction by **redundant** codes: Natural languages allow to detect many errors.

Example in informatics: **Parity bits**. Even parity: 9 bits per byte. 9th bit makes number of one-bits even. Allows detection of single-bit errors. Computer memory sometimes uses 9 bits per byte for this purpose.

Other example: ISBN code (International Standard Book Number) – last character is a parity character

Entropy and quantification of information

Shannon's information theory:
Information as a measurable, statistical property of signals

How can we measure information and redundancy of characters in a message?

Assumption: N -character alphabet $\{ x_1, x_2, \dots, x_N \}$

Number of bits per character:

$$H_0 = \log_2 N$$

($\log_2 N = (\log N)/(\log 2)$)

Information content of a single character x_i : $\log_2 \frac{1}{p(x_i)}$

Entropy = average value of information content of all characters

$$= H = \sum_{k=1}^N p(x_k) * \log_2 \frac{1}{p(x_k)}$$

Binary encoding needs at least, on average, H bits per character.

Redundancy: $R = H_0 - H$.

Example: Four-letter alphabet $\{a, b, c, d\}$

Probabilities: $p_a = 0.5, p_b = 0.25, p_c = 0.125, p_d = 0.125$

Thus:

$H_0 = 2$ bits per character encodable

Entropy: $0.5 * 1 + 0.25 * 2 + 0.125 * 3 + 0.125 * 3 = 1.75$ bits per character encoded

Redundancy: 0.25 bits per character

Examples:

– $a \mapsto 00, b \mapsto 01, c \mapsto 10, d \mapsto 11$: on average 2 bits per character

– $a \mapsto 0, b \mapsto 10, c \mapsto 110, d \mapsto 111$: on average 1.75 bits per character (optimal, no redundancy)

3. Computer architecture

Principal design of computers

Standard scheme, realised in most contemporary computers (with small modifications):

"von Neumann architecture"

Seven principles of the von Neumann architecture:

1. The computer consists of four components:

(i) A MEM (memory) containing both programs and data,

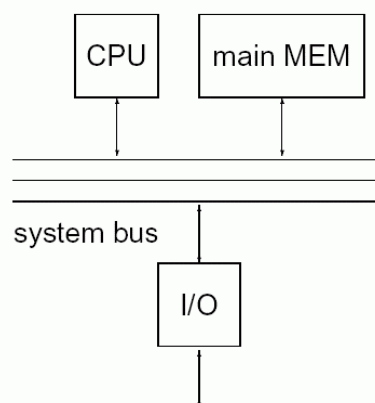
(ii) a CU (control unit) which interprets the program,

(iii) an ALU (arithmetical and logical unit) which performs the operations on the data values, and

(iv) an IO unit (input/output unit) for interacting with external components (like keyboard and mouse, screen and printer, disk memory) and other computers.

CU, ALU and a part of the MEM ("registers") are often put together in one component, the CPU (central processing unit). The other part of MEM is called "main memory".

Main structure of a computer



2. The structure of the computer is largely *independent of the programme*. ("General purpose computer")
3. Main memory locations are used *for both programme and data*; they can both be read and written by the machine.
4. The main memory is split into *cells of the same size*. Each of these cells can be *addressed* by its number.
5. A programme consists of a number of separated instructions, encoded as cell contents of the main memory. They are mostly executed *in the order in which they occur* in the main memory.
6. The order of execution of instructions can be changed by conditional or unconditional *jumps* to another programme address.
7. *Binary codes* are used to represent numbers and programme instructions.

Computer components

main MEM, I/O, CPU, and system bus for connection

main MEM, I/O: address (input), data (both input and output), selection (input), r/w (input)

CPU: ALU, data registers, CU, processor bus for connecting other components

system bus: address, data, control

Address: the binary number, designating, for example, a cell (often a byte) in the memory

System bus

Connections between different components of a computer, i.e. between CPU, MEM and IO unit.

Consists of a number of **electrical lines**

Three important types of lines:

- **Address bus:** Is used for addresses, during access to memory or I/O.

- **Data bus:** Is used for the data when it is **written** from the CPU to the MEM or to the IO unit, or when it is **read** by the CPU from the MEM or the IO unit.

- **Control bus:** for deciding if data is read or written and what kind of access takes place (e.g. MEM or IO).

Example of an interface of 8 MB byte-organised MEM:

All lines are binary, i.e. either 0 or 1.

1 selection line: To define if the memory is being accessed at the moment.

23 address lines: When selection line is active, to define which of the $2^{23} = 8$ million bytes is to be read or written.

1 read/write line: When selection line is active, to define if the byte at the address defined by the address lines is being read or written.

8 data lines: When memory is being read, the content of the addressed byte is put on these lines; when memory is being written to, the content of these lines is put into the byte location.

How are larger memories organized?

64 MB byte-wide memory: From 8 such chips.

- Computer must use 26 address lines and 8 data lines.
- Lower 23 address lines are used as address lines of memory chips.
- Upper 3 address lines are used to select which of the 8 memory chips is addressed.

Example

Let chips be numbered from 0..7.

Address lines are numbered A0..A25.

In order to address RAM chip 6, the following is done:

Since $6_{10} = 110_2$, the chip is to be accessed if $(A_{25}A_{24}A_{23}) = (110)$.

The latter can be expressed by the Boolean function $A_{25} \wedge A_{24} \wedge \neg A_{23}$.

this means "*A25 and A24 and not A23*".

Boolean: having only a truth value (true or false) as result.

The Boolean function is built with logic chips. Its output is used for the selection line of the RAM.

32 MB 4-byte-wide memory: From 4 such chips.

- Computer uses 23 address lines and 32 data lines.
- All address lines are used for all chips.
- Each chip is connected to another 8 lines of the 32 data lines.