

Data Representations

- We know that a computer's CPU and main memory can only store binary bits.
 - How are individual characters and character strings stored in a computer?
 - How are integers stored in a computer?
 - How are negative numbers represented?
 - How are floating point numbers, such as 2.75, represented in a computer?

Characters

- We need to distinguish between:
 - The concept of a character.
 - How a character looks on a screen or page.
- The concept of “upper case A” can have many appearances depending on the font used to display the character.
 - A (Arial) A (Times New Roman)
 - A (Comic Sans MS) *A (Baskerville MT)*
- The concept of “upper case A” is different from the concept of “numeral 4”.

Characters (2)

- Each character (concept of a character) is represented by a sequence of ones and zeros.
 - Rather than writing the individual ones and zeros we will group the bits into groups of 8 bits (called a byte) or groups of bytes (called words).
 - Each byte will be given in hexadecimal notation.
 - 0x5E is the bit sequence 01011110
- There are several standards for assigning a code to each character.

Characters (3)

- The 2 main standards are:
 - ASCII (American Standard Code for Information Interchange)
 - Unicode
 - An extension of ASCII to allow for different “languages”.
- IBM has its own “standard” called EBCDIC.

ASCII

USASCII code chart

Character	b7	b6	b5	b4	b3	b2	b1
0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	1
2	0	0	0	0	0	1	0
3	0	0	0	0	0	1	1
4	0	0	0	0	1	0	0
5	0	0	0	0	1	0	1
6	0	0	0	0	1	1	0
7	0	0	0	0	1	1	1
8	0	0	0	1	0	0	0
9	0	0	0	1	0	0	1
10	0	0	0	1	0	1	0
11	0	0	0	1	0	1	1
12	0	0	0	1	1	0	0
13	0	0	0	1	1	0	1
14	0	0	0	1	1	1	0
15	0	0	0	1	1	1	1
NUL	0	0	1	0	0	0	0
DLE	0	0	1	0	0	1	0
SP	0	0	1	0	1	0	0
@	0	0	1	0	1	0	1
P	0	0	1	0	1	1	0
Q	0	0	1	0	1	1	1
A	0	0	1	1	0	0	0
B	0	0	1	1	0	0	1
R	0	0	1	1	0	1	0
S	0	0	1	1	0	1	1
C	0	0	1	1	1	0	0
D	0	0	1	1	1	0	1
E	0	0	1	1	1	1	0
F	0	0	1	1	1	1	1
G	0	0	1	1	1	1	1
H	0	0	1	1	1	1	1
I	0	0	1	1	1	1	1
J	0	0	1	1	1	1	1
K	0	0	1	1	1	1	1
L	0	0	1	1	1	1	1
M	0	0	1	1	1	1	1
N	0	0	1	1	1	1	1
O	0	0	1	1	1	1	1
P	0	0	1	1	1	1	1
Q	0	0	1	1	1	1	1
R	0	0	1	1	1	1	1
S	0	0	1	1	1	1	1
T	0	0	1	1	1	1	1
U	0	0	1	1	1	1	1
V	0	0	1	1	1	1	1
W	0	0	1	1	1	1	1
X	0	0	1	1	1	1	1
Y	0	0	1	1	1	1	1
Z	0	0	1	1	1	1	1
[0	0	1	1	1	1	1
\	0	0	1	1	1	1	1
]	0	0	1	1	1	1	1
^	0	0	1	1	1	1	1
_	0	0	1	1	1	1	1
`	0	0	1	1	1	1	1
a	0	0	1	1	1	1	1
b	0	0	1	1	1	1	1
c	0	0	1	1	1	1	1
d	0	0	1	1	1	1	1
e	0	0	1	1	1	1	1
f	0	0	1	1	1	1	1
g	0	0	1	1	1	1	1
h	0	0	1	1	1	1	1
i	0	0	1	1	1	1	1
j	0	0	1	1	1	1	1
k	0	0	1	1	1	1	1
l	0	0	1	1	1	1	1
m	0	0	1	1	1	1	1
n	0	0	1	1	1	1	1
o	0	0	1	1	1	1	1
p	0	0	1	1	1	1	1
q	0	0	1	1	1	1	1
r	0	0	1	1	1	1	1
s	0	0	1	1	1	1	1
t	0	0	1	1	1	1	1
u	0	0	1	1	1	1	1
v	0	0	1	1	1	1	1
w	0	0	1	1	1	1	1
x	0	0	1	1	1	1	1
y	0	0	1	1	1	1	1
z	0	0	1	1	1	1	1
{	0	0	1	1	1	1	1
	0	0	1	1	1	1	1
}	0	0	1	1	1	1	1
~	0	0	1	1	1	1	1
DEL	1	1	1	1	1	1	1

Extending ASCII

- While ASCII worked great for English and simple “typewriter” systems, it could not cope with accented characters very well.
 - â could be made by an “a” character followed by a backspace (BS) character and then a circumflex (“^”) character.
- Screen-based displays enabled simple graphics to be constructed from block graphics characters.

Extended ASCII Table

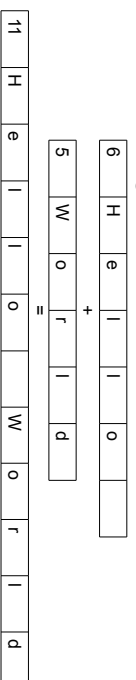
DEC	HEX	CHAR	DEC	HEX	CHAR	DEC	HEX	CHAR	DEC	HEX	CHAR
128	80	€	160	A0	À	192	C0	Ì	224	E0	Š
129	81	£	161	A1	Á	193	C1	Í	225	E1	Ṧ
130	82	¤	162	A2	Â	194	C2	Î	226	E2	Œ
131	83	¥	163	A3	Ã	195	C3	Ï	227	E3	Œ̇
132	84	¦	164	A4	Ä	196	C4	Ï̇	228	E4	Ɔ
133	85	§	165	A5	Å	197	C5	Ï̈	229	E5	Ɔ̇
134	86	¨	166	A6	Ä̇	198	C6	Ï̈̇	230	E6	Ɔ̈
135	87	©	167	A7	Å̇	199	C7	Ï̈̈	231	E7	Ɔ̈̇
136	88	ª	168	A8	Ä̈	200	C8	Ï̈̈̇	232	E8	Ɔ̈̈
137	89	«	169	A9	Å̈	201	C9	Ï̈̈̈	233	E9	Ɔ̈̈̇
138	8A	»	170	AA	Ä̈̇	202	CA	Ï̈̈̈̇	234	EA	Ɔ̈̈̈
139	8B	¼	171	AB	Å̈̇	203	CB	Ï̈̈̈̈	235	EB	Ɔ̈̈̈̇
140	8C	½	172	AC	Ä̈̈	204	CC	Ï̈̈̈̈̇	236	EC	Ɔ̈̈̈̈
141	8D	¾	173	AD	Å̈̈	205	CD	Ï̈̈̈̈̈	237	ED	Ɔ̈̈̈̈̇
142	8E	¿	174	AE	Ä̈̈̇	206	CE	Ï̈̈̈̈̈̇	238	EE	Ɔ̈̈̈̈̈
143	8F		175	AF	Å̈̈̇	207	CF	Ï̈̈̈̈̈̈	239	EF	Ɔ̈̈̈̈̈̇
144	90	À	176	B0	À̇	208	D0	Ï̈̈̈̈̈̇	240	F0	Ɔ̈̈̈̈̈̈
145	91	Á	177	B1	Á̇	209	D1	Ï̈̈̈̈̈̈	241	F1	Ɔ̈̈̈̈̈̈̇
146	92	Â	178	B2	Â̇	210	D2	Ï̈̈̈̈̈̈̇	242	F2	Ɔ̈̈̈̈̈̈̈
147	93	Ã	179	B3	Ã̇	211	D3	Ï̈̈̈̈̈̈̈̇	243	F3	Ɔ̈̈̈̈̈̈̈̇
148	94	Ä	180	B4	Ä̇	212	D4	Ï̈̈̈̈̈̈̈̈	244	F4	Ɔ̈̈̈̈̈̈̈̈
149	95	Å	181	B5	Å̇	213	D5	Ï̈̈̈̈̈̈̈̈̇	245	F5	Ɔ̈̈̈̈̈̈̈̈̇
150	96	Æ	182	B6	Æ̇	214	D6	Ï̈̈̈̈̈̈̈̈̈	246	F6	Ɔ̈̈̈̈̈̈̈̈̈
151	97	Ç	183	B7	Ç̇	215	D7	Ï̈̈̈̈̈̈̈̈̈̇	247	F7	Ɔ̈̈̈̈̈̈̈̈̈̇
152	98	È	184	B8	È̇	216	D8	Ï̈̈̈̈̈̈̈̈̈̈	248	F8	Ɔ̈̈̈̈̈̈̈̈̈̈
153	99	É	185	B9	É̇	217	D9	Ï̈̈̈̈̈̈̈̈̈̈̇	249	F9	Ɔ̈̈̈̈̈̈̈̈̈̈̇
154	9A	Ê	186	BA	Ê̇	218	DA	Ï̈̈̈̈̈̈̈̈̈̈̈	250	FA	Ɔ̈̈̈̈̈̈̈̈̈̈̈
155	9B	Ë	187	BB	Ë̇	219	DB	Ï̈̈̈̈̈̈̈̈̈̈̈̇	251	FB	Ɔ̈̈̈̈̈̈̈̈̈̈̈̇
156	9C	Ì	188	BC	Ì̇	220	DC	Ï̈̈̈̈̈̈̈̈̈̈̈̈	252	FC	Ɔ̈̈̈̈̈̈̈̈̈̈̈̈
157	9D	Í	189	BD	Í̇	221	DD	Ï̈̈̈̈̈̈̈̈̈̈̈̈̇	253	FD	Ɔ̈̈̈̈̈̈̈̈̈̈̈̈̇
158	9E	Î	190	BE	Î̇	222	DE	Ï̈̈̈̈̈̈̈̈̈̈̈̈̈	254	FE	Ɔ̈̈̈̈̈̈̈̈̈̈̈̈̈
159	9F	Ï	191	BF	Ï̇	223	DF	Ï̈̈̈̈̈̈̈̈̈̈̈̈̈̇	255	FF	Ɔ̈̈̈̈̈̈̈̈̈̈̈̈̈̇

Unicode

- Even with the extended ASCII table, the table is very oriented towards Latin based alphabets.
- As computers became more prevalent across the world it became necessary to accommodate other alphabets.
 - Cyrillic
 - Kanji and Hiragana
 - Arabic
- Unicode uses 16-bit codes.
 - 256 tables of 256 characters each.
 - ASCII is bottom half of the first table.

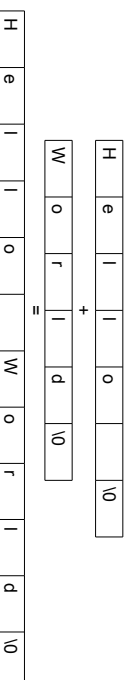
Strings

- Strings are sequences of characters.
- Fixed length strings
 - Older programming languages, e.g. Pascal, used fixed length strings.
 - The first byte gave the length of the string (number of characters in the string) and the remaining bytes were the characters in the string.



Strings (2)

- Variable length string
 - These are much more flexible than fixed length strings.
 - Variable length strings are terminated by a NULL character (0x00 or \0).
 - Concatenating strings (joining strings together is a little more complicated).



Integers

- A single byte can represent 256 different values, 0 up to 255.
- For larger numbers we need to use multiple bytes.
 - A 2 byte word can represent $2^{16}=65,536$ values, 0 up to 65535.
 - A 4 byte word can represent $2^{32}=4,294,967,296$ values.
 - No matter how many bytes you add there is always a limit.

2's Complement Numbers

- Lets us limit ourselves to bytes for the moment.
- There are 2 ways to implement negative numbers.
 - Use a sign bit.
 - If the most significant bit is 1 then the number is negative, and if it is 0 then it is positive.
 - Though simple it has a few problems.
 - There are two ways to represent zero.
 - $0x00 = 00000000 = +0$
 - $0x80 = 10000000 = -0$

2's Complement Numbers (2)

- Arithmetic manipulation is messy
 - Try adding 1 to -1 and you get:
 - $0x81 + 0x01 = 0x82 = -2$???????
 - You would need different hardware for adding two positive numbers and a positive number and a negative number.
- 2's complement
 - To convert a positive number into a negative number subtract it from 0x00.

0000 0010 = 0x02	2
0000 0001 = 0x01	1
0000 0000 = 0x00	0
1111 1111 = 0xFF	-1
1111 1110 = 0xFE	-2
1111 1101 = 0xFD	-3

2's Complement Numbers (3)

- Now adding positive and negative numbers does not cause any complications.
 - $-1 + 2 = 0xFF + 0x02 = 0x01 = 1$
 - And we did not have to even involve a carry bit!
- Multiplication by 2 still works as a shift left operation.
 - $4 \times 2 = 0x04 \ll 1 = 0000\ 0100 \ll 1 = 0000\ 1000 = 0x08 = 8$
 - $-2 \times 2 = 0xFE \ll 1 = 1111\ 1110 \ll 1 = 1111\ 1100 = 0xFC = -4$

2's Complement Numbers (4)

- Division by 2 is still a shift right but with a slight twist.
 - $4 / 2 = 0x02 \gg 1 = 0000\ 0100 \gg 1 = 0000\ 0010 = 0x02 = 2$
 - $4 / 2 = 0xFC \gg 1 = 1111\ 1100 \gg 1 = 1111\ 1110 = 0xFE = -2$
- If dividing a positive number by 2 we need to shift in a 0 to the most significant bit.
- If dividing a negative number by 2 we need to shift in a 1 to the most significant bit.
- Most CPUs have shift instructions that shift the carry bit into the register being shifted. The carry bit can be set depending on whether the register contains a negative 2's complement number (as determined by a bit in the status register).

Floating Point Numbers

- If we only have binary numbers, how can we represent numbers such as:
 - 0.00000025
 - 1,000,000,000,000,000
 - -127.5
 - -0.075
- First let us think back to how we humans have learnt to deal with very big and very small numbers.

Scientific Notation

- Using scientific notation we can reduce very large numbers and very small numbers into manageable entities.

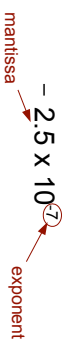
$$- 0.00000025 = 2.5 \times 10^{-7}$$

$$- 1,000,000,000,000,000 = 1.0 \times 10^{15}$$

$$- 127.5 = -1.275 \times 10^2$$

$$- 0.075 = -7.5 \times 10^{-2}$$

- Each part of the scientific notation number has a specific name.



The diagram shows the scientific notation $- 2.5 \times 10^{-7}$. A red arrow points from the label "mantissa" to the number "2.5". Another red arrow points from the label "exponent" to the number "-7".

Scientific Notation in Binary

- The standard for scientific notation in binary is the IEEE 754 standard.
- First what does the binary number 0.1001 represent?
 - $-(0 \times 2^0) + (1 \times 2^{-1}) + (0 \times 2^{-2}) + (0 \times 2^{-3}) + (1 \times 2^{-4})$
 - $- 0 + (1 \times 0.5) + (0 \times 0.25) + (0 \times 0.125)$
 - $+ (1 \times 0.0625)$
 - $- 0.5625$

Scientific Notation in Binary (2)

- Now let us try 1234.5625
- Convert the part to the left of the decimal point into binary and the part to the right of the decimal point into binary.
- 1234.5625 = 10011010010.1001
- Now normalise the binary number
 - 1.00110100101001) x 2¹⁰¹⁰
- A nice feature of normalised binary numbers is that the digit to the left of the binary point is always a 1!

mantissa

exponent

Scientific Notation in Binary (3)

- IEEE 754 covers 32-bit, 64-bit and 128-bit floating point numbers.
- A 32-bit IEEE 754 number is calculated and stored in memory as follows.
 - Convert integer part into binary.
 - Convert fractional part into binary
 - Normalise, drop leading 1 and add sufficient 0's to the start to get a 23-bit mantissa
 - Add 127 to the number of places shifted to get an 8-bit exponent (shift left is positive, shift right is negative).

