

Bytes

Bits are almost always bundled into 8-bit collections called Bytes.

One byte can be represented by 2 hex digits

A      F  
1010    1111 - 1 byte (2 nibbles)

With 8 bits, we can represent  $2^8 = 256$  different 'things'  
for example we can represent the numbers 0 to 255

A single byte is also sufficient to represent all the characters in a western keyboard.

The American System for Coded Information Interchange (ASCII) character set is a standardized table in which each number from 0 to 255 is associated with either a symbol (that we call a printable character) or with an action (that we call an unprintable character)

## Ascii character Set.

Dec	Hx	Oct	Char		Dec	Hx	Oct	Html	Chr		Dec	Hx	Oct	Html	Chr		Dec	Hx	Oct	Html	Chr
0	0	000	NUL	(null)	32	20	040	&#32;	Space		64	40	100	&#64;	Ø		96	60	140	&#96;	~
1	1	001	SOH	(start of heading)	33	21	041	&#33;	!		65	41	101	&#65;	A		97	61	141	&#97;	a
2	2	002	STX	(start of text)	34	22	042	&#34;	"		66	42	102	&#66;	B		98	62	142	&#98;	b
3	3	003	ETX	(end of text)	35	23	043	&#35;	#		67	43	103	&#67;	C		99	63	143	&#99;	c
4	4	004	EOT	(end of transmission)	36	24	044	&#36;	\$		68	44	104	&#68;	D		100	64	144	&#100;	d
5	5	005	ENQ	(enquiry)	37	25	045	&#37;	%		69	45	105	&#69;	E		101	65	145	&#101;	e
6	6	006	ACK	(acknowledge)	38	26	046	&#38;	&		70	46	106	&#70;	F		102	66	146	&#102;	f
7	7	007	BEL	(bell)	39	27	047	&#39;	!		71	47	107	&#71;	G		103	67	147	&#103;	g
8	8	010	BS	(backspace)	40	28	050	&#40;	(		72	48	110	&#72;	H		104	68	150	&#104;	h
9	9	011	TAB	(horizontal tab)	41	29	051	&#41;	)		73	49	111	&#73;	I		105	69	151	&#105;	i
10	A	012	LF	(NL line feed, new line)	42	2A	052	&#42;	*		74	4A	112	&#74;	J		106	6A	152	&#106;	j
11	B	013	VT	(vertical tab)	43	2B	053	&#43;	+		75	4B	113	&#75;	K		107	6B	153	&#107;	k
12	C	014	FF	(NP form feed, new page)	44	2C	054	&#44;	,		76	4C	114	&#76;	L		108	6C	154	&#108;	l
13	D	015	CR	(carriage return)	45	2D	055	&#45;	-		77	4D	115	&#77;	M		109	6D	155	&#109;	m
14	E	016	SO	(shift out)	46	2E	056	&#46;	.		78	4E	116	&#78;	N		110	6E	156	&#110;	n
15	F	017	SI	(shift in)	47	2F	057	&#47;	/		79	4F	117	&#79;	O		111	6F	157	&#111;	o
16	10	020	DLE	(date link escape)	48	30	060	&#48;	0		80	50	120	&#80;	P		112	70	160	&#112;	p
17	11	021	DC1	(device control 1)	49	31	061	&#49;	1		81	51	121	&#81;	Q		113	71	161	&#113;	q
18	12	022	DC2	(device control 2)	50	32	062	&#50;	2		82	52	122	&#82;	R		114	72	162	&#114;	r
19	13	023	DC3	(device control 3)	51	33	063	&#51;	3		83	53	123	&#83;	S		115	73	163	&#115;	s
20	14	024	DC4	(device control 4)	52	34	064	&#52;	4		84	54	124	&#84;	T		116	74	164	&#116;	t
21	15	025	NAK	(negative acknowledgement)	53	35	065	&#53;	5		85	55	125	&#85;	U		117	75	165	&#117;	u
22	16	026	SYN	(synchronous idle)	54	36	066	&#54;	6		86	56	126	&#86;	V		118	76	166	&#118;	v
23	17	027	ETB	(end of trans. block)	55	37	067	&#55;	7		87	57	127	&#87;	W		119	77	167	&#119;	w
24	18	030	CAN	(cancel)	56	38	070	&#56;	8		88	58	130	&#88;	X		120	78	170	&#120;	x
25	19	031	EM	(end of medium)	57	39	071	&#57;	9		89	59	131	&#89;	Y		121	79	171	&#121;	y
26	1A	032	SUB	(substitute)	58	3A	072	&#58;	:		90	5A	132	&#90;	Z		122	7A	172	&#122;	z
27	1B	033	ESC	(escape)	59	3B	073	&#59;	;		91	5B	133	&#91;	[		123	7B	173	&#123;	{
28	1C	034	FS	(file separator)	60	3C	074	&#60;	<		92	5C	134	&#92;	\		124	7C	174	&#124;	
29	1D	035	GS	(group separator)	61	3D	075	&#61;	=		93	5D	135	&#93;	]		125	7D	175	&#125;	)
30	1E	036	RS	(record separator)	62	3E	076	&#62;	>		94	5E	136	&#94;	^		126	7E	176	&#126;	~
31	1F	037	US	(unit separator)	63	3F	077	&#63;	?		95	5F	137	&#95;	]		127	7F	177	&#127;	DEL

Source : [www.LookupTables.com](http://www.LookupTables.com)

## Extended ASCII Codes

Source: [www.LookupTables.com](http://www.LookupTables.com)

This table is effectively a way to associate a character with a number. Printable characters mean something to us, as humans, but apart from that, there is nothing special about them.

Note the distinction between the character '5' and the number 5.

The 8 bit representation for the character '5' is:

00110101 (i.e., 35<sub>hex</sub>)

Whereas the number 5 is :

00000101 (i.e., 05<sub>hex</sub>)

## Representing Characters in a Program.

Consider the following C program fragment in which we declare 2 Variables: ch, of type character, and i, of type integer.

The Program also Contains a number of Constants :

'5' , 53 , 0x35 , 065.

```
char ch;
int i;

int main(){
    ch = '5';
    i = 53;
    printf("ch as an int = %d, in hex = %x, in octal = %o and as a char = %c\n", ch, ch, ch, ch);
    printf("i as an int = %d, in hex = %x, in octal = %o and as a char = %c\n", i, i, i, i);
    printf("-----\n");

    ch = 0x35;
    i = 0x35;
    printf("ch as an int = %d, in hex = %x, in octal = %o and as a char = %c\n", ch, ch, ch, ch);
    printf("i as an int = %d, in hex = %x, in octal = %o and as a char = %c\n", i, i, i, i);
    printf("-----\n");

    ch = 065;
    i = 065;
    printf("ch as an int = %d, in hex = %x, in octal = %o and as a char = %c\n", ch, ch, ch, ch);
    printf("i as an int = %d, in hex = %x, in octal = %o and as a char = %c\n", i, i, i, i);
    printf("-----\n");
```

'5' is a character Constant

53 is a decimal number Constant

0x35 is a hex number Constant

065 is an Octal number Constant

All of these Constants have the  
Same bit Sequence : 00110101

In other words, they all have the same "Count value" - the distance from 0 to each number, in its respective base, is the same.

Our program assigns these constants, in turn, to both the character variable and to the Integer Variable and it Prints the values of those variables in decimal, hex, octal and character representations.

The program output for this fragment is:

ch = '5'	ch as an int = 53, in hex = 35, in octal = 65 and as a char = 5
i = 5	i as an int = 53, in hex = 35, in octal = 65 and as a char = 5
ch = 53	ch as an int = 53, in hex = 35, in octal = 65 and as a char = 5
i = 53	i as an int = 53, in hex = 35, in octal = 65 and as a char = 5
ch = 0x35	ch as an int = 53, in hex = 35, in octal = 65 and as a char = 5
i = 0x35	i as an int = 53, in hex = 35, in octal = 65 and as a char = 5
ch = 0b110101	ch as an int = 53, in hex = 35, in octal = 65 and as a char = 5
i = 0b110101	i as an int = 53, in hex = 35, in octal = 65 and as a char = 5

Regardless of the representation we use, each variable is given the same bit sequence and, when printing, we can write that sequence out in any one of the desired representations.

Does this mean that integer variables and character variables are the same?

?

At least at the bit representation level, they appear to be.

There are, however, some important differences

- (1) Character variables are only 1-byte in size  
 Integer variables are usually many bytes in size  
 -In our example program, they are 4-bytes in size.
- (2) Strongly-typed programming languages are strict in distinguishing between characters and integers and will not allow the programmer to combine them in expression.

To illustrate the size difference between Integers (32 bits) and characters, (8 bits) Consider the following program fragments:

```
ch = 241
ch = 242
ch = 243
ch = 244
ch = 245
ch = 246
ch = 247
ch = 248
ch = 249
ch = 250
ch = 251
ch = 252
ch = 253
ch = 254
ch = 255
ch = 0
ch = 1
ch = 2
ch = 3
ch = 4
ch = 5
ch = 6
ch = 7
ch = 8

#include<stdio.h>
unsigned char ch;

int main(){
    printf("The values that ch can contain are: \n");
    for(ch = 0; ;ch++)
        printf("ch = %d\n", ch);
}
```

```
#include<stdio.h>
#define LARGEINT 4294967000

unsigned int i;

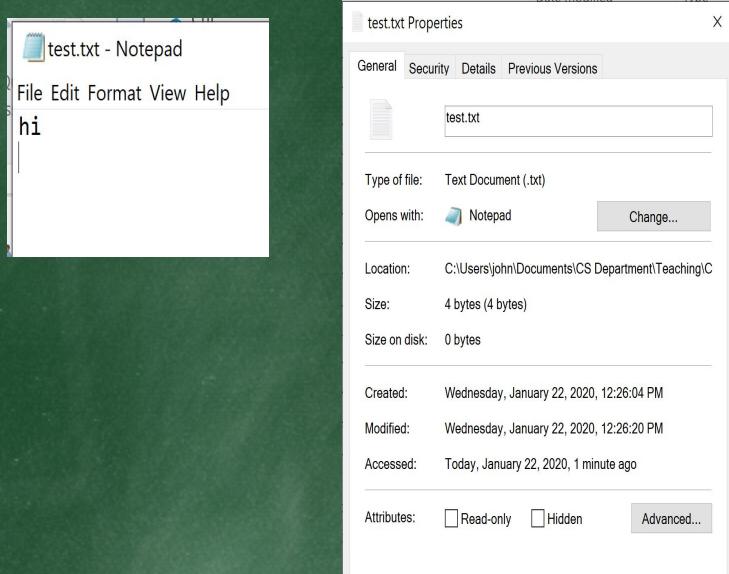
int main(){
    printf("The values that i can contain are: \n");
    for(i = LARGEINT; ;i++)
        printf("i = %u\n", i);
}
```

```
i = 4294967284
i = 4294967285
i = 4294967286
i = 4294967287
i = 4294967288
i = 4294967289
i = 4294967290
i = 4294967291
i = 4294967292
i = 4294967293
i = 4294967294
i = 4294967295
i = 0
i = 1
i = 2
i = 3
i = 4
i = 5
i = 6
i = 7
i = 8
i = 9
i = 10
i = 11
```

Bit Sequence overflows  
 after  $2^8 - 1$

Bit Sequence  
 overflows after  
 $2^{32} - 1$ .

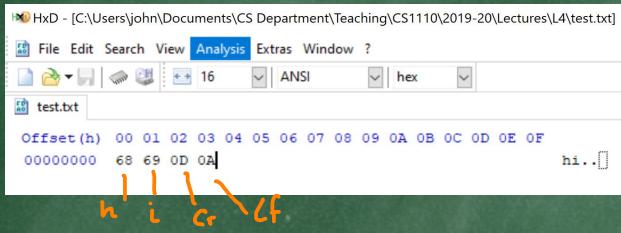
# How are characters manipulated by Editors?



The 4 bytes in test.txt are:

'h' 'i' '\r' '\n'

In Hex, the contents of test.txt is: 68 69 0D 0A



HxD is a Hex Editor

Compare the size of the files generated by notepad and word for the same 'file contents'.

Representing characters by numbers is another level of abstraction.

## Lots of Bytes

Multiples of bytes						V · T · E
Decimal			Binary			
Value	Metric		Value	IEC	JEDEC	
1000	kB	kilobyte	1024	KiB	<b>kibibyte</b>	KB kilobyte
$1000^2$	MB	megabyte	$1024^2$	MiB	<b>mebibyte</b>	MB megabyte
$1000^3$	GB	gigabyte	$1024^3$	GiB	<b>gibibyte</b>	GB gigabyte
$1000^4$	TB	terabyte	$1024^4$	TiB	<b>tebibyte</b>	—
$1000^5$	PB	petabyte	$1024^5$	PiB	<b>pebibyte</b>	—
$1000^6$	EB	exabyte	$1024^6$	EiB	<b>exbibyte</b>	—
$1000^7$	ZB	zettabyte	$1024^7$	ZiB	<b>zebibyte</b>	—
$1000^8$	YB	yottabyte	$1024^8$	YiB	<b>yobibyte</b>	—

Orders of magnitude of data

In this course, we will continue to use the older notation. Thus:

$$\text{Kilo} - \text{K} - 2^{10}$$

$$\text{Mega} - \text{M} - 2^{20}$$

$$\text{Giga} - \text{G} - 2^{30}$$

$$\text{Tera} - \text{T} - 2^{40}$$

$$\text{Peta} - \text{P} - 2^{50}$$

$$\text{Exa} - \text{E} - 2^{60}$$

$$\text{Zetta} - \text{Z} - 2^{70}$$

$$\text{Yotta} - \text{Y} - 2^{80}$$

1 yottabyte is close to the estimated number of atoms in a human body.