

# Rappresentazione dei dati

Andrea Passerini  
passerini@disi.unitn.it

Informatica

- Un *bit* (b) rappresenta una cifra binaria. E' l'unità minima di informazione.
- Un *Byte* (B) è costituito da 8 bit. Permette di codificare 256 entità di informazione distinte (e.g. caratteri dell'alfabeto, segni di interpunzione)
- Una *parola* (*word*) rappresenta l'insieme di Byte che possono essere trattati da un elaboratore in un'operazione. Il numero di Byte di una parola dipende dall'elaboratore, ad esempio 2, 3, 4, 8 corrispondenti a 16, 24, 32 o 64 bit.

# Multipli del Byte

- Kilobyte (KB) =  $2^{10}$  = 1024 Byte (e.g. un file di testo da 300KB)
- Megabyte (MB) = 1024KB  $\approx$  1 milione di Byte (e.g. un'immagine di 30MB)
- Gigabyte (GB) = 1024MB  $\approx$  1 miliardo di Byte (e.g. un hard disk da 80GB)
- Terabyte (TB) = 1024GB  $\approx$  1000 miliardi di Byte (e.g. un archivio da 20TB)

- Esprime la velocità di trasferimento dati in bit per secondo (bps).
- Multipli
  - Kilobit per secondo (Kbps) = 1000bps (e.g. un modem a 56Kbps)
  - Megabit per secondo (Mbps) = 1000Kbps (e.g. una ADSL a 4Mbps)

- La dimensione dei file viene generalmente espressa in multipli **binari** del **Byte** (e.g. 1MB =  $2^{10}$  KByte =  $2^{20}$  Byte = 1,048,576 Byte)
- Il bit rate viene generalmente espresso in multipli **decimali** del **bit** (e.g. 1Mbps = 1,000,000 bps)
- Quindi per scaricare un file da 4MB con un'ADSL a 4Mbps sono necessari (nella situazione ottimale):

$$\frac{4 \times 2^{20} \times 8}{4 \times 1000000} = 8.388608$$

ossia circa 8 secondi e mezzo.

- Questo non considerando la compressione e il tempo necessario a trasmettere l'intestazione dei pacchetti

# Nota sulla dimensione di HD

- La dimensione dei file viene generalmente espressa in multipli **binari** del **Byte** (e.g. 1MB =  $2^{10}$  KByte =  $2^{20}$  Byte = 1,048,576 Byte)
- La dimensione degli HD viene generalmente espressa in multipli **decimali** del **Byte** (e.g. 1GB =  $10^9$  Byte)
- Quindi un HD da 250GB contiene file per:

$$\frac{250 \times 10^9}{250 \times 2^{30}} \approx 233 \text{ GB}$$

# Hertz (Hz)

- Grandezza per misurare la rapidità (frequenza) dei dispositivi digitali
- Il nome deriva dal fisico Heinrich Rudolf Hertz
- 1 Hz corrisponde ad un ciclo o una oscillazione al secondo
- Multipli
  - Kilohertz (KHz) = 1000 Hz (si usa per misurare la frequenza di refresh dello schermo)
  - Megahertz (MHz) = 1000 KHz
  - Gigahertz (GHz) = 1000 MHz (si usa per misurare la frequenza di clock dei calcolatori)

- Il clock è un dispositivo che funziona come un metronomo sincronizzando tutte le operazioni dei dispositivi digitali
- Nota:
  - L'unità di misura Hz si può usare come indicatore relativo della velocità di elaborazione di un computer
  - ovvero si può usare per paragonare due processori con la stessa architettura (e.g. un Pentium 4 a 2.8GHz contro un Pentium 4 a 3.2GHz)
  - NON si possono però fare paragoni fra processori diversi (e.g. RISC Motorola contro Pentium)

# Dots Per Inch (DPI) ossia punti per pollice

- Grandezza per misurare la densità di punti o *definizione* o *risoluzione*.
- Un pollice quadrato corrisponde ad un area di  $2.54 \text{ cm}^2$
- Stampanti e schermi usano matrici di punti (*pixel*) per rappresentare immagini 2D.
- Maggiore è il numero di punti nell'unità di area maggiore è l'accuratezza con la quale si definisce un'immagine o un testo
- Esempio:
  - 300 dpi per le stampanti laser
  - 1200 dpi per gli scanner
- La risoluzione di uno schermo si indica riportando il numero di pixel visualizzati sul lato orizzontale e su quello verticale (e.g. 1280x800 o 860x640)

## Codifica binaria

- Qualsiasi informazione deve essere codificata in binario per poter essere trattata da un calcolatore.
- Vedremo come vengono codificati:
  - Numeri interi
  - Numeri “reali” (*real*)
  - Caratteri

## Numero di bit fissato

- Il calcolatore può fare operazioni su due numeri solo se sono codificati con lo stesso numero di bit.
- Si fissa il numero di bit  $k$  con cui si rappresenta un certo insieme di numeri (e.g. interi a 8 bit)
- Ogni numero di tale insieme deve essere rappresentato con  $k$  bit, eventualmente aggiungendo zeri a sinistra
- Esempio per interi positivi a  $k = 8$  bit

$$(54)_{10} = (00110110)_2$$

# Rappresentazione di interi

## Codifica valore assoluto con segno

- Si riserva il bit più significativo (quello più a sinistra) al segno (0=positivo, 1= negativo).
- Si codifica con i restanti bit il modulo del numero.
- Esempio con  $k = 8$  bit:

$$(+54)_{10} = (00110110)_2$$

$$(-54)_{10} = (10110110)_2$$

## Problema

- Problema: il *bit di segno* deve essere trattato in maniera diversa rispetto agli altri bit (complica i circuiti per somma e sottrazione).
- Si utilizza la codifica in *complemento a 2* che permette di realizzare sottrazioni con complementazioni e somme (non la vedremo).

- Si verifica quando la somma di due numeri è al di fuori del rango (*range*) di valori permessi nella rappresentazione scelta per i numeri.
- Ad esempio il rango dei numeri interi rappresentati in complemento a 2 con  $k$  bit è  $[-2^{k-1}, 2^{k-1} - 1]$
- Il risultato dell'operazione non è corretto nel caso si verifichi un overflow.
- Gli elaboratori elettronici verificano la condizione di overflow quando effettuano l'addizione binaria e la segnalano mettendo ad 1 uno speciale *bit di overflow*.

# Traslazione logica (*shift logico*)

- Consiste nello spostare a destra (*shift logico a destra*) o a sinistra (*shift logico a sinistra*) i bit di un numero binario.
- Nella traslazione a destra, il bit meno significativo viene perso, quello più significativo viene posto a 0.
- Un numero binario traslato a destra viene diviso per 2.
- Nella traslazione a sinistra, il bit più significativo viene perso, quello meno significativo viene posto a 0.
- Un numero binario traslato a sinistra viene moltiplicato per 2.

# Operazioni eseguite in termini di altre operazioni

- Con la codifica in complemento a 2, la sottrazione si realizza con una complementazione (operazione semplicissima) ed un'addizione.
- Si può realizzare un unico circuito che effettui somme e sottrazioni con notevoli risparmi di costi.
- Una moltiplicazione può essere realizzata tramite una sequenza di addizioni e di traslazioni a sinistra.
- Una divisione può essere realizzata tramite una sequenza di sottrazioni e di traslazioni a destra.
- Le operazioni più semplici vengono eseguite da appositi circuiti (a livello *hardware*).
- Operazioni più complesse sono eseguite in termini di esecuzione di altre operazioni più semplici sotto il controllo di programmi (a livello *software*).

# Rappresentazione di numeri “reali” (*real*)

## Rappresentazione in virgola fissa

- Si stabilisce un numero di bit  $k_1$  da assegnare alla parte intera ed un numero di bit  $k_2$  da assegnare alla parte frazionaria.
- Ad esempio per numeri a 32 bit se ne assegnano 16 alla parte intera e 16 alla parte frazionaria.
- Adatta solo a casi particolari in cui l'intervallo di valori da rappresentare è noto a priori.
- Inadatta nella maggior parte delle applicazioni scientifiche.

## Rappresentazione in virgola mobile (*floating point*)

- Si parte dalla rappresentazione scientifica in cui il numero è il prodotto di due parti: una parte frazionaria ed un fattore di scala, che è una potenza del 10.
- Ad esempio il numero 23.5 può essere rappresentato come:

$$23.5 \times 10^0$$

$$235 \times 10^{-1}$$

...

$$2.35 \times 10^1$$

$$0.235 \times 10^2$$

...

- Rappresentazione scientifica normalizzata: la parte frazionaria ha la cifra più a sinistra diversa da 0 e subito seguita dal punto decimale (e.g.  $2.35 \times 10^1$ ).
- Nella rappresentazione floating point si rappresenta il numero come una coppia: la *mantissa* corrispondente alla parte frazionaria, la *caratteristica* (o *esponente*) che corrisponde all'esponente del fattore di scala.

# Rappresentazione floating point binaria

- Nel caso binario mantissa ed esponente sono rappresentati in binario, ed il fattore di scala è una potenza del 2.
- Esempi:

Numero	Rapp. norm.	Mantissa	Esponente
+101010.0	$+1.010100 \times 10^{101}$	+1.010100	+101
+0.000110	$+1.100000 \times 10^{-100}$	+1.100000	-100
-110.1100	$-1.101100 \times 10^{10}$	-1.101100	+010

- Repertorio: insieme di caratteri considerati, definito mediante i nomi dei caratteri e magari una loro rappresentazione visiva.
- Numero di codice: tabella che associa ad ogni carattere un numero da un dato insieme di numeri naturali.
- Codifica: un metodo per associare a ciascun numero di codice una sequenza di bit.
- Nel caso più semplice ogni carattere ha un numero tra 0 e 127 e la codifica è semplicemente la codifica binaria del numero in 7 bit.

- Acronimo di *American Standard Code for Information Interchange*.
- 7 bit per carattere, si possono rappresentare  $2^7 = 128$  caratteri distinti. I codici sono tipicamente scritti in notazione esadecimale.
- I codici da 0 a 1F sono usati per *caratteri di controllo*
- I codici da 20 a 7E sono usati per *caratteri stampabili*
- Ordine alfabetico: cifre 0-9 prima dei caratteri alfabetici, maiuscole prima delle minuscole (si rifletterà nell'ordinamento lessicografico delle stringhe).

# Tabella dei codici ASCII

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b> (null)	32	20	040	<b>#32;</b> <b>Space</b>		64	40	100	<b>#64;</b> <b>@</b>		96	60	140	<b>#96;</b> <b>`</b>	
1	1	001	<b>SOH</b> (start of heading)	33	21	041	<b>#33;</b> <b>!</b>		65	41	101	<b>#65;</b> <b>A</b>		97	61	141	<b>#97;</b> <b>a</b>	
2	2	002	<b>STX</b> (start of text)	34	22	042	<b>#34;</b> <b>"</b>		66	42	102	<b>#66;</b> <b>B</b>		98	62	142	<b>#98;</b> <b>b</b>	
3	3	003	<b>ETX</b> (end of text)	35	23	043	<b>#35;</b> <b>#</b>		67	43	103	<b>#67;</b> <b>C</b>		99	63	143	<b>#99;</b> <b>c</b>	
4	4	004	<b>EOT</b> (end of transmission)	36	24	044	<b>#36;</b> <b>\$</b>		68	44	104	<b>#68;</b> <b>D</b>		100	64	144	<b>#100;</b> <b>d</b>	
5	5	005	<b>ENQ</b> (enquiry)	37	25	045	<b>#37;</b> <b>%</b>		69	45	105	<b>#69;</b> <b>E</b>		101	65	145	<b>#101;</b> <b>e</b>	
6	6	006	<b>ACK</b> (acknowledge)	38	26	046	<b>#38;</b> <b>&amp;</b>		70	46	106	<b>#70;</b> <b>F</b>		102	66	146	<b>#102;</b> <b>f</b>	
7	7	007	<b>BEL</b> (bell)	39	27	047	<b>#39;</b> <b>'</b>		71	47	107	<b>#71;</b> <b>G</b>		103	67	147	<b>#103;</b> <b>g</b>	
8	8	010	<b>BS</b> (backspace)	40	28	050	<b>#40;</b> <b>(</b>		72	48	110	<b>#72;</b> <b>H</b>		104	68	150	<b>#104;</b> <b>h</b>	
9	9	011	<b>TAB</b> (horizontal tab)	41	29	051	<b>#41;</b> <b>)</b>		73	49	111	<b>#73;</b> <b>I</b>		105	69	151	<b>#105;</b> <b>i</b>	
10	A	012	<b>LF</b> (NL line feed, new line)	42	2A	052	<b>#42;</b> <b>*</b>		74	4A	112	<b>#74;</b> <b>J</b>		106	6A	152	<b>#106;</b> <b>j</b>	
11	B	013	<b>VT</b> (vertical tab)	43	2B	053	<b>#43;</b> <b>+</b>		75	4B	113	<b>#75;</b> <b>K</b>		107	6B	153	<b>#107;</b> <b>k</b>	
12	C	014	<b>FF</b> (NP form feed, new page)	44	2C	054	<b>#44;</b> <b>,</b>		76	4C	114	<b>#76;</b> <b>L</b>		108	6C	154	<b>#108;</b> <b>l</b>	
13	D	015	<b>CR</b> (carriage return)	45	2D	055	<b>#45;</b> <b>-</b>		77	4D	115	<b>#77;</b> <b>M</b>		109	6D	155	<b>#109;</b> <b>m</b>	
14	E	016	<b>SO</b> (shift out)	46	2E	056	<b>#46;</b> <b>.</b>		78	4E	116	<b>#78;</b> <b>N</b>		110	6E	156	<b>#110;</b> <b>n</b>	
15	F	017	<b>SI</b> (shift in)	47	2F	057	<b>#47;</b> <b>/</b>		79	4F	117	<b>#79;</b> <b>O</b>		111	6F	157	<b>#111;</b> <b>o</b>	
16	10	020	<b>DLE</b> (data link escape)	48	30	060	<b>#48;</b> <b>0</b>		80	50	120	<b>#80;</b> <b>P</b>		112	70	160	<b>#112;</b> <b>p</b>	
17	11	021	<b>DC1</b> (device control 1)	49	31	061	<b>#49;</b> <b>1</b>		81	51	121	<b>#81;</b> <b>Q</b>		113	71	161	<b>#113;</b> <b>q</b>	
18	12	022	<b>DC2</b> (device control 2)	50	32	062	<b>#50;</b> <b>2</b>		82	52	122	<b>#82;</b> <b>R</b>		114	72	162	<b>#114;</b> <b>r</b>	
19	13	023	<b>DC3</b> (device control 3)	51	33	063	<b>#51;</b> <b>3</b>		83	53	123	<b>#83;</b> <b>S</b>		115	73	163	<b>#115;</b> <b>s</b>	
20	14	024	<b>DC4</b> (device control 4)	52	34	064	<b>#52;</b> <b>4</b>		84	54	124	<b>#84;</b> <b>T</b>		116	74	164	<b>#116;</b> <b>t</b>	
21	15	025	<b>NAK</b> (negative acknowledge)	53	35	065	<b>#53;</b> <b>5</b>		85	55	125	<b>#85;</b> <b>U</b>		117	75	165	<b>#117;</b> <b>u</b>	
22	16	026	<b>SYN</b> (synchronous idle)	54	36	066	<b>#54;</b> <b>6</b>		86	56	126	<b>#86;</b> <b>V</b>		118	76	166	<b>#118;</b> <b>v</b>	
23	17	027	<b>ETB</b> (end of trans. block)	55	37	067	<b>#55;</b> <b>7</b>		87	57	127	<b>#87;</b> <b>W</b>		119	77	167	<b>#119;</b> <b>w</b>	
24	18	030	<b>CAN</b> (cancel)	56	38	070	<b>#56;</b> <b>8</b>		88	58	130	<b>#88;</b> <b>X</b>		120	78	170	<b>#120;</b> <b>x</b>	
25	19	031	<b>EM</b> (end of medium)	57	39	071	<b>#57;</b> <b>9</b>		89	59	131	<b>#89;</b> <b>Y</b>		121	79	171	<b>#121;</b> <b>y</b>	
26	1A	032	<b>SUB</b> (substitute)	58	3A	072	<b>#58;</b> <b>:</b>		90	5A	132	<b>#90;</b> <b>Z</b>		122	7A	172	<b>#122;</b> <b>z</b>	
27	1B	033	<b>ESC</b> (escape)	59	3B	073	<b>#59;</b> <b>;</b>		91	5B	133	<b>#91;</b> <b>[</b>		123	7B	173	<b>#123;</b> <b>{</b>	
28	1C	034	<b>FS</b> (file separator)	60	3C	074	<b>#60;</b> <b>&lt;</b>		92	5C	134	<b>#92;</b> <b>\</b>		124	7C	174	<b>#124;</b> <b> </b>	
29	1D	035	<b>GS</b> (group separator)	61	3D	075	<b>#61;</b> <b>=</b>		93	5D	135	<b>#93;</b> <b>]</b>		125	7D	175	<b>#125;</b> <b>}</b>	
30	1E	036	<b>RS</b> (record separator)	62	3E	076	<b>#62;</b> <b>&gt;</b>		94	5E	136	<b>#94;</b> <b>^</b>		126	7E	176	<b>#126;</b> <b>~</b>	
31	1F	037	<b>US</b> (unit separator)	63	3F	077	<b>#63;</b> <b>?</b>		95	5F	137	<b>#95;</b> <b>_</b>		127	7F	177	<b>#127;</b> <b>DEL</b>	

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

# Limitazioni del codice ASCII

- I caratteri internazionali di numerose lingue europee (quali ad esempio è,ù,ç,å,æ,ü,ø) non sono rappresentabili
- nessuno dell'elevatissimo numero di simboli delle lingue asiatiche è rappresentabile
- La standardizzazione è importante: nella trasmissione e memorizzazione elettronica i caratteri sono rappresentati da insiemi di Byte ed è importante che il "trasmettitore" ed il "ricevente" adottino le stesse convenzioni!
- In assenza di opportune convenzioni, testi generati su un dato sistema possono risultare corrotti se visualizzati in un sistema diverso (capita facilmente con la posta elettronica).

# ISO Latin-1 (ISO 8859-1)

- Il repertorio contiene il repertorio ASCII come sottoinsieme e i codici per questi caratteri sono identici a quelli ASCII
- Il codice usa 8 bit: 256 caratteri distinti
- Contiene vari simboli usati dai linguaggi dell'Europa occidentale (Italiano, Francese, Spagnolo, Tedesco, Danese, etc.):

¡ ¢ £ ¤ ¥ ¦ § ¨ © ª « ¬ ® ¯ ° ± ² ³ ´ µ ¶ · ¸  
¹ º » ¼ ½ ¾ ¿ À Á Â Ã Ä Å Æ Ç È É Ê Ë Ì Í Î Ï Ð  
Ñ Ò Ó Ô Õ Ö × Ø Ù Ú Û Ü Ý Þ ß à á â ã ä å æ ç è  
é ê ë ì í î ï ð ñ ò ó ô õ ö ÷ ø ù ú û ü ý þ ÿ

- La famiglia contiene ben 15 alfabeti standard, tra cui ad esempio:

8859-2 (Latin-2) lingue dell'Europa centrale e orientale

8859-5 (Latin/Cyrillic) lingue Slave

8859-7 (Latin/Greek) Greco moderno

8859-9 (Latin-5) Turco

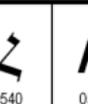
8859-10 (Latin-6) lingue nordiche (Islandese)

8859-15 (Latin-9) Latin-1 con l'Euro

- ISO 10646 è uno standard internazionale che definisce lo UCS (Universal Character Set)
- Il repertorio è molto ampio e contiene decine di migliaia di caratteri già definiti, con spazio per espansioni future
- Estende ISO Latin-1 nello stesso senso in cui ISO Latin-1 estende ASCII
- UNICODE è uno standard definito dal consorzio UNICODE (fondato nel '92) e definisce un repertorio e una codifica pienamente compatibili con ISO 10646
- UNICODE implementa solo una parte di ISO 10646

- Il vantaggio di UNICODE è di definire un insieme di caratteri adatti a trattare tutti i linguaggi, mentre la famiglia ISO 8859 definisce sottoinsiemi di caratteri adatti a trattare solo alcuni linguaggi alla volta
- Mentre con altri standard lo stesso carattere può avere codifiche diverse (a seconda dell'alfabeto usato), in UNICODE ogni carattere ha una codifica unica
- Ha una grande diffusione industriale (Apple, HP, IBM, Microsoft, Oracle, Sun, etc.)
- E' supportato da vari sistemi operativi ed internet browser più recenti.

- Il codice usa 16 bit (fino a 64K caratteri)
- I primi 256 codici coincidono con quelli di ISO Latin-1
- Esempio di codifica per caratteri armeni:

	053	054	055	056	057	058
0		Հ 0540	Ր 0550		Տ 0570	Ր 0580
1	Ա 0531	Զ 0541	Յ 0551	Մ 0561	Ճ 0571	Գ 0581
2	Բ 0532	Ղ 0542	Ի 0552	Բ 0562	Ղ 0572	Լ 0582
3	Գ 0533	Ճ 0543	Փ 0553	Գ 0563	Ճ 0573	Փ 0583
4	Դ 0534	Մ 0544	Ք 0554	Գ 0564	Մ 0574	Ք 0584