

Università degli Studi di Milano «La Statale» Dipartimento di Informatica

# La rappresentazione di testi: caratteri e stringhe

(parte A)

1



# Rappresentiamo lettere, parole, testi: caratteri e stringhe

- Ora sappiamo come codificare numeri (naturali, interi, razionali)
- Anche le informazioni «alfabetiche» (o «testuali») vanno codificate!
- Idea base:

facciamo corrispondere ogni carattere (lettere) ad un numero

- (in modo arbitrario... purché standard)
   (cioè deciso una volta per tutte!)
- ▶ per es: 1 <==> A, 2 <==> B ....
- una parola, o un testo, è solo una sequenza di caratteri (una «stringa»)
  - memorizzata come sequenza di numeri (interi senza segno!)
  - (compreso un carattere «spazio» per staccare le parole e un carattere «accapo» per staccare le righe)
- Problemi apparentemente banali:
  - quali caratteri? quanti caratteri? quale corrispondenza caratteri-numeri?

Architettura degli elaboratori

- 2

Aritmetica binaria



#### La codifica dei caratteri

- Per il momento ci concentriamo sui caratteri base degli alfabeti occidentali, che, per motivi storici, sono stati i primi a essere codificati.
- Per gli usi più comuni, i caratteri usati sono relativamente pochi
  - ▶ 10 cifre numeriche 0,1,2,3,4,5,6,7,8,9
  - ▶ 26 minuscole: a,b,c, ...
  - ▶ 26 maiuscole: A,B,C, ...
  - ▶ 6 simboli di interpunzione ( . , ?!;:)
  - pochissimi operatori matematici di uso comune ( + < > )
  - apostrofi, accenti, virgolette (' " `)
  - un'altra manciatina di simboli uso comune (° @ #\$ | / \ \*)
  - In totale: meno di 100 caratteri:
     7 bit (128 combinazioni) ci bastano per rappresentarli tutti.
- Nota: per i numeri esiste una codifica "naturale" (ad es. lo zero è rappresentato da un insieme di bit nulli) per i caratteri possiamo usare ad esempio l'ordine alfabetico.

Architettura degli elaboratori

- 4 -

Aritmetica binaria

4



### La codifica dei caratteri: ASCII

- Con 7 bit, rappresento 128 caratteri diversi
- Basta fissare una corrispondenza fissa tra numeri appartenenti all'intervallo 0-127 e i caratteri.
- La cosa importante è che questa corrispondenza sia nota ed accettata da tutti (pena la scorretta interpretazione dei testi).
  - ▶ la codifica dei caratteri deve essere stabilita da uno standard.
- Lo standard più adottato è lo standard ASCII (persino prima dei computer! telescriventi, telegrafi...)

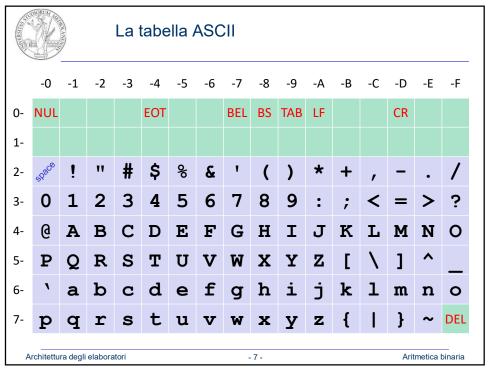
ASCII: American Standard Code for Information Interchange (si pronucia «àschii», ma in Ita a volte anche «àscii»)

Architettura degli elaboratori

- 5 -

Aritmetica binaria

6





# I caratteri ASCII «di controllo» cioè «non stampabili»

- Caratteri non visibili ma con un effetto sul testo, come:
- Sono una trentina (codici da 0x00 a 0x1F, e 0xFF)
- Originalmente, controllavano le funzionalità di dispositivi come telescriventi & telegrafi
- Molti di questi sono oggi senza senso per la rappresentazione dei testi, e sono inutilizzati.
  - es:
    - BEL «suona un campanello» :-D
      DEL e BS «cancella l'ultimo carattere» (in due modi diversi)
      e molti altri
  - ▶ purtroppo ce li portiamo dietro, solo per motivi storici ⊗
  - consumano spazio di rappresentazione (cioè lasciano meno codici liberi per rappresentare i caratteri utili)
- Alcuni hanno ancora degli usi, come «accapo» (

Architettura degli elaboratori

- 8 -

Aritmetica binaria

8



# I caratteri ASCII "di controllo" ancora usati (nella rappresentazione di testi)

- TAB: tabulazione (spazio orizzontale, per allineare);
- I due caratteri che segnalano la fine di una riga di testo:
  - CR carriage return (codice 13): originalmente: «riporta la testina scrivente all'inizio della riga»
  - ▶ LF line feed (codice 10): originalmente: «scorri il foglio della pagina giù di una riga»
  - Oggi (purtroppo) per «andare accapo»...
    - alcuni standard (Linux, MacOS) usano il 2°
    - altri (Windows) usano ancora entrambi i caratteri (in sequenza)
    - (altri ancora, in disuso (es. Commodore) usavano solo il 1°)
- NUL, il carattere nullo (codice 0):
  - spesso usato ad es per segnalare la fine di una stringa in memoria
  - ▶ stringa = una sequenza di caratteri in memoria
- EOT end of text (fine del testo)
  - a volte usato per segnalare la file di un file di testo;



### ASCII - una nota

- Per decenni, è stato l'unico standard diffuso per rappresentare testi.
- L'informatica ne porta ancora (e ne porterà sempre) le tracce. Infatti:
  - perché per la mail si usa il simbolo @?
  - perché per gli hashtag si usa #?
  - perché nei path dei file, o negli URL dei siti web, si usano / e \ ?
  - perché in HTML nei tag si aprono con < e chiudono con > ?
  - perché in molti URL appare la «tilde» ~ (per indicare homepage) ?
  - perché per dire «citazione», es nei forum o nelle mail, si usa > blah ?
  - perché il prompt dei comandi è denotata da >, ad es: C:\>?
  - perché il «pipe command» di Unix ha quella sintassi, con | ?
  - perché si sono così diffusi gli smiley, come :-) quando (oggi) ci sono gli emoticons, come © ?
  - perché il dash si trova scritto spesso come un meno ?
  - perché spesso si trova scritto perche ' al posto di perché ?
  - » perché il celebre sito slashdot si chiama proprio slashdot?

Architettura degli elaboratori

- 10 -

Aritmetica binaria

10



#### ASCII – una nota

- Per decenni, è stato l'unico standard diffuso per rappresentare testi.
- L'informatica ne porta ancora (e ne porterà sempre) le tracce.
- Perché, in quasi tutti i linguaggi di programmazione (dal C al Go)...
  - si usano parentesi (), [], {}, <> (con la più varia semantica)?
  - si usano spesso i simboli \$ , & e % (con la più varia semantica)?
  - ▶ la divisione si denota con / (e non, mettiamo, con ÷)?
  - ▶ la moltiplicazione si denota con \* (e non, mettiamo, con × o )?
  - il minore-o-uguale si denota con <= (e non con ≤)?</p>
  - diverso-da si denota con != oppure <> (e non con ≠)?
  - le virgolette sono aperte (e chiuse!) con " o ' (e mai con, p.e., « e »)?
  - ▶ per gli operatori logici si usa per es & | ! (invece che, per es, ∧ ∨ ¬)?
  - Ie frecce si scrivono con sequenze: -> o => (e non → o ⇒)?
  - ▶ per la radice quadrata si usano lettere: sqrt, e non √?
  - i quantificatori si scrivono ad es. come forEach e non ∀?
  - il valore di +infinito si scrive come +inf e non +∞
  - alla x alla seconda si scrive x<sup>2</sup> e non x<sup>2</sup>?
  - ▶ Etc

Architettura degli elaboratori

- 11 -

Aritmetica binaria



## Il codice ASCII non ha :\_(

- Lettere accentate (o nessun altro modificatore di lettera)
- Simboli matematici, eccetto i pochissimi basilari (+, -, <, >)
- Lettere greche
  - ▶ neanche quelle di uso comune, come pi-greco o epsilon
- Anzi... nessun alfabeto che non sia quello latino di base!
  - ▶ no: cirillico, giapponese, coreano, etc
- Simboli per valute che non siano il dollaro
- Interpunzioni rovesciate (per es. usate in spagna)
- Emoji ⊗
- etc...

Insomma la codifica ASCII è (a malapena) sufficiente solo per

- i testi più basici
- testi espressi in linguaggi formali (e non per caso: sono linguaggi progettati attorno a queste limitazione)

Architettura degli elaboratori

- 12 -

Aritmetica binaria

12



## Esercizi di riepilogo

#### In ASCII:

- Quale sequenza di 7 bit rappresenta la lettera «erre minuscola»?
- Quale sequenza di 7 bit rappresenta la lettera «erre maiuscola»?
- Quale sequenza di 7 bit rappresenta la lettera «effe minuscola»?
- Quale sequenza di 7 bit rappresenta la lettera «effe maiuscola»?
- Noti un pattern, nelle sequenze di 7 bit che rappresentano le lettere minuscole e maiuscole?
   (Fai altri esempi se necessario, e verifica di avere ragione)
- Scrivi in esadecimale il BYTE che rappresenta il punto esclamativo (il MSB – cioè il Most Significant Bit, che non viene utilizzato, vale 0)

Architettura degli elaboratori

- 13 -

Aritmetica binaria