



# *Computação em Larga Escala*

*General Problems – Algorithmic analysis 1*

António Rui Borges

# *Summary*

- *Text processing in Portuguese*
  - *Character encodings*
  - *Rules for text processing*

## *Text processing in Portuguese - 1*

*Character encoding* is essential to store and process textual information people routinely use to communicate among themselves. Many such codes were introduced in the computer world to express written contents as time went by.

The most popular ones are

- ASCII (*American Standard Code for Information Interchange*)  
it is a 7 bit code able to represent 95 graphical symbols and 33 control signals, which was used for many years to encode english language texts
- ISO/IEC 8859 (first published in 1987)  
it is a 8 bit extension of ASCII, which provides 193 graphical symbols of the Latin Script, covering most western european languages and standard romanizations of east asian languages
- Unicode (first published in late 1980s)  
it is a computing industry standard for consistent character encoding of world languages; it contains presently a repertoire of 137,439 graphical symbols covering 146 modern and historic languages; UTF-8, its most common implementation, uses 1 byte for the first 128 code points (ASCII characters) and up to 4 bytes for other characters.

## *Text processing in Portuguese - 2*

*Character encoding* in UTF-8 supposes a character representation in one up to four bytes. It encompasses ASCII encoding as the one byte character representation class. All other classes are multibyte and follow the rules described below.

UTF-8 encoding format			
Byte 0	Byte 1	Byte 2	Byte3
0xxxxxxx			
110xxxxx	10xxxxxx		
1110xxxx	10xxxxxx	10xxxxxx	
11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

## *Text processing in Portuguese - 3*

Portuguese special characters encodings					
Upper case	UTF-8	ISO/IEC 8859	Lower case	UTF-8	ISO/IEC 8859
á	C3A1	E1	Á	C381	C1
à	C3A0	E0	À	C380	C0
â	C3A2	E2	Â	C382	C2
ã	C3A3	E3	Ã	C383	C3
é	C3A9	E9	É	C389	C9
è	C3A8	E8	È	C388	C8
ê	C3AA	EA	Ê	C38A	CA
í	C3AD	ED	Í	C38D	CD
ì	C3AC	EC	Ì	C38C	CC
ó	C3B3	F3	Ó	C393	D3
ò	C3B2	F2	Ò	C392	D2
ô	C3B4	F4	Ô	C394	D4
õ	C3B5	F5	Õ	C395	D5
ú	C3BA	FA	Ú	C39A	DA
ù	C3B9	F9	Ù	C399	D9
ç	C3A7	E7	Ç	C387	C7

## *Text processing in Portuguese - 4*

- the uppercase and lowercase alphabets should be treated as the same on detecting the occurrence of each letter
- in the same way, á – à – â – ã should be treated as instances of the letter a, é – è – ê should be treated as instances of the letter e, í – ì should be treated as instances of the letter i, ó – ô – õ should be treated as instances of the letter o, ú – û should be treated as instances of the letter u and ç should be treated as an instance of the letter c
- a *word* is defined as any sequence of characters, consisting of alphanumeric or underscore characters delimited by white spaces and separation or punctuation symbols

## *Text processing in Portuguese - 5*

- a *white space* is a *space* character (0x20), a *tab* character (0x9), a *newline* character (0xA) or a *carriage return* character (0xD)
- a *separation symbol* is a *hyphen* (-), a *double quotation mark* (" 0x22 - “ 0xE2809C - ” 0xE2809D), a *bracket* ( [ ] ) or a *parentheses* ( ( ) )
- a *punctuation symbol* is a *full point* (.), a *comma* (,), a *colon* (:), a *semicolon* (;), a *question mark* (?), an *exclamation point* (!), a *dash* (— 0xE28093) or an *ellipsis* (... 0xE280A6)
- the *apostrophe* (' 0x27) and *single quotation marks* (‘ 0xE28098 - ’ 0xE28099) are considered here to merge two words into a single one.