Standards for language encoding

Tomaž Erjavec Dept. of Knowledge Technologies Jožef Stefan Institute

ESSLLI 2011

A few words about me

- Tomaž Erjavec
 Department of Knowledge Technologies
 Jožef Stefan Institute
 Ljubljana
- <u>http://nl.ijs.si/et/, tomaz.erjavec@ijs.si</u>
- Areas of work: compilation and annotation of corpora and other language resources, encoding standards, digital libraries (text-critical editions)
- Web page for this course: <u>http://nl.ijs.si/et/teach/esslli11/</u> and, of course, Moodle

Overview of the course

- 1. Introduction, character sets
- 2. Structuring data: XML
- 3. Encoding for the humantities: TEI
- 4. Standards for LRs: ISO
- 5. Semantic Web: W3C standards

I. Introduction

What are standards?

- dictionary:
 - 1. an obligatory uniform regulation for measurement, quantity or quality
 - 2. that which specifies how something can or must be
- consensually accepted regulations, which are public and contain explicit definitions
- the main purpose is to harmonise industrial practice in various fields in order to enable interchange

History of standardisation

- XVIII century: in France each region (village) has its own units of measurement; also, different objects (say a field or forest) are measured differently
- how to define a uniform system of measurements: search for a single unit from which it would be possible to derive all other measures
- meter: one ten-millionth of the length of the meridian through Paris, from the North Pole to the equator
- the importance of standardisation grows with the industrial revolution: mechanical and electrical engineering, construction work...
- today, standards encompass even such "soft" fields as the organisation of business (ISO 9000)
- big business: companies that check compliance with standards

Standards and best practices

- National standard bodies: DIN, ANSI, SIST
- International standard bodies:
 - IEC: International Electrotechnical Commission
 - ISO: International organisation for standardisation
 - IETF: Internet Engineering Task Force
 - W3C:World Wide Web Consortium
 - Unicode consortium
- Initiatives:
 - MUFI: Medieval Unicode Font Initiative
 - TEI: the Text Encoding Initiative
- Best practices:
 - Penn Treebank PoS tags
 - TIGER annotation scheme

Language resources

Corpora

- monolingual and multilingual
- o general and domain specific
- raw text or annotated
- text or speech

Lexica

- o monolingual and multilingual
- words and lemmas
- o entities (names)
- o phrases (terms)

Annotations

- Morphological level: lemmas/stems, coarse (PoS) or fine (MSD) grained tags
- Syntax: syntactic trees or dependencies
- Semantics: word senses, semantic roles
- Named entities: names, dates, numeric expressions
- Terms, time & space expressions and relations
- Anaphora: anaphoric links
- Parallel corpora: sentence and word/phrase alignments
- Meta-data: information about the resource

Utility of LRs

A basis for:

- HLT development
 - Training: datasets for inducing language models
 - Testing: datasets for evaluating performance
- Empirically driven (applied) linguistics:
 - Corpus linguistics
 - Lexicography, Terminography
 - Language teaching

Why standards for encoding of digital data?

Traditionally, each developer made LRs to work with their particular software and for their particular needs Problems:

- Iongevity: advances in technology make programs soon obsolete and data bound to these programs becomes unreadable
- interchange: difficult to use data on other platforms or pass it between programs
- *exploitation:* difficult to re-use the data for other purposes
- intelligibility: no public and stable specifications of the format
- validation: we don't know whether certain data is written according to the specification or not

Standardisation of LRs

- LRs are expensive to produce so, a good idea if they are reusable and long-lasting
- LRs are becoming larger and with more complex annotations – no good for everyone to reinvent the wheel
- Familiarity with standards in this area helps to produce good resources and to be able to use the resources already produced:
 - freely available resources: Google, (MetaShare)
 - o LDC, ELRA
 - o corpora@uib.no

Levels of standardisation

- Characters: the basic building blocks
- XML: structuring the data and assigning annotations
- TEI: a large vocabulary of XML elements: "encoding text for scholarly purposes"
- ISO standards: some basic things like dates and languages; and recent attempts to standardise many different types of LRs
- Semantic Web: Meta-data and ontologies

II. Character sets

- Characters are the "atoms" of textual resources
- It still often happens that characters are garbled in processing, resulting in useless text
- Currently, Unicode is gaining ground but is still not the only character set in use
- Unicode is relatively complex

Character encoding

- Digital computers store data as (binary) numbers
- There is no a priori connection between these numbers and characters (of an alphabet)
- If there are no conventions for this mapping or if there are too many → chaos
- Standards and quasi standards: ASCII, ISO 8859, (Windows, Mac), Unicode

Basic concepts I.

character

- abstract concept (An "A" is something like a Platonic entity: it is the idea of an "A" and not the "A" itself)
- a character does not by itself have a mapping to a number or a specific graphical representation
- usually it is descriptivelly defined, e.g. "Greek letter lower-case alpha ", and the graphical representation is given only as a suggestion, "α"

Basic concepts II.

- character set
 - o a set of characters
 - o each character has an associated character code
- character code
 - 1-1 relation between the character from a character set and a number, e.g. A = 26, B = 27, ...
 - Note:

character codes are often written in hexadecimal: $0 \rightarrow 0, 1 \rightarrow 1, 2 \rightarrow 2, \dots 9 \rightarrow 9,$ $10 \rightarrow A, 11 \rightarrow B, \dots, 15 \rightarrow F,$ $16 \rightarrow 10, 17 \rightarrow 11, \dots,$ $254 \rightarrow FE, 255 \rightarrow FF, 266 \rightarrow 100$

Example: the ASCII character set

Dec	Char										
33	I.	49	1	65	Α	81	Q	97	a	113	q
34		50	2	66	В	82	R	98	b	114	r
35	#	51	3	67	С	83	S	99	с	115	s
36	\$	52	4	68	D	84	т	100	d	116	t
37	%	53	5	69	E	85	U	101	e	117	u
38	æ	54	6	70	F	86	v	102	f	118	v
39	•	55	7	71	G	87	w	103	g	119	w
40	(56	8	72	н	88	х	104	h	120	×
41)	57	9	73	1	89	Y	105	i	121	У
42	*	58	:	74	J	90	Z	106	j	122	z
43	+	59	;	75	ĸ	91	[107	k	123	{
44	,	60	<	76	L	92	X	108	ι	124	1
45	-	61	=	77	Μ	93	1	109	m	125	}
46		62	>	78	Ν	94	^	110	n	126	~
47	1	63	?	79	0	95	-	111	0	127	_
48	0	64	0	80	Р	96	•	112	Р		

e.g.

in the ASCII character set

the <u>character</u> lower case Latin a

has the <u>character</u> <u>code 97</u>

Basic concepts III.

glyph

- a graphical representation of a character
- one character can have more than one glyph
 e.g. the character "upper-case Latin A" ↔ glyphs A, A, A
- sometimes one glyph can be associated with more than one character, e.g. the glyph P ↔ characters "upper-case Latin P", "upper-case cyrillic R", "upper-case Greek Rho"

font

- a set of glyphs (for some character set):
 A, B, C, Č, D, ...
- sometimes a font does not cover the complete character set!

Some character sets

- ASCII oldest, contains only the letters of the English alphabet + punctuation, numbers
- Family of characters sets ISO 8879
- The Windows family of character sets ("code pages")
- Unicode

ASCII

- American Standard Code for Information Interchange (1950')
- 7-bit encoding: character codes 0-127
- 0-31 control characters + formatting characters:

Esc, Line Feed, tab, space,...

 32-126 – punctuation and special characters, numbers, upper- and lower-case letters:

!"#\$%&'()*+,-./0123456789:; <=>?@ABCDEFGHIJKLMNOP QRSTUVWXYZ[\]^_`abcdefgh ijklmnopqrstuvwxyz{|}~

The ISO 8859 family

- need for extra characters for national (European) alphabets:
 - o 80's
 - 8 bits, so twice as many chars as in ASCII
 - first half = ASCII, second half = new characters
- International Standards Organisation publishes character sets for particular groups of European languages: ISO 8859 (-1 .. -12)
- ISO 8859-1 (ISO Latin 1) Western European languages:

i¢£¤¥¦§[…]©^a≪¬®[−]°±²³´µ¶·,¹°»¼½ ¾;ÀÁÂÃÄÅÆÇÈÉÊËÌÍÎÏĐŇÒÓÔÕ ÖרÙÚÛÜÝÞßàáâãäåæçèéêëìíîï ðñòóôõö÷øùúûüýþÿ

ISO 8859-2 C&E European languages

	NBSP	Ą	v	Ł	¤	Ľ	Ś	§		Š	Ş	Ť	Ź	SHY	Ž	Ż
A –	00A0	0104	02D8	0141	00A4	013D	015A	00A7	00A8	0160	015E	0164	0179	00AD	017D	017B
	160	161	162	<i>163</i>	164	165	166	16 7	<i>168</i>	<i>169</i>	<i>170</i>	171	172	173	174	175
	0	ą	c.	ł	-	ľ	ś	~		š	Ş	ť	ź	~	ž	ż
B -	00B0	0105	02DB	0142	00B4	013E	015B	02C7	00B8	0161	015F	0165	017A	02DD	017E	017C
	176	177	178	179	180	181	182	183	184	185	186	18 7	188	189	190	191
	Ŕ	Á	Â	Ă	Ä	Ĺ	Ć	Ç	Č	É	Ę	Ë	Ě	Í	Î	Ď
C –	0154	00C1	00C2	0102	00C4	0139	0106	00C7	010C	00C9	0118	00CB	011A	00CD	00CE	010E
	192	193	1 <i>9</i> 4	195	196	19 7	198	199	200	201	<i>202</i>	<i>203</i>	204	205	206	20 7
	Ð	Ń	Ň	Ó	Ô	Ő	Ö	x	Ř	Ů	Ú	Ű	Ü	Ý	Ţ	ß
D-	0110	0143	0147	00D 3	00D4	0150	00D6	00D7	0158	016E	OODA	0170	OODC	OODD	0162	OODF
	<i>208</i>	<i>209</i>	210	211	212	213	214	215	216	217	218	219	220	221	222	223
	ŕ	á	â	ă	ä	ĺ	ć	ç	č	é	ę	ë	ě	í	î	ď
E-	0155	00E1	00E2	0103	00E4	013A	0107	00E7	010D	00E9	0119	OOEB	011B	OOED	OOEE	010F
	224	225	226	227	228	229	<i>230</i>	231	232	233	234	235	<i>236</i>	237	238	239
	đ	ń	ň	ó	ô	ő	ö	÷	ř	ů	ú	ű	ü	ý	ţ	•
F-	0111	0144	0148	00F3	00F4	0151	00F6	00F7	0159	016F	OOFA	0171	00FC	OOFD	0163	02D9
	240	241	242	243	244	245	246	247	248	2 4 9	250	251	252	253	254	255
	-0	-1	-2	-3	-4	-5	-6	-7	-8	-9	- A	-В	- c	-D	-Е	-F

Confusion

- Microsoft also developed their own code pages:
 - Windows CP1252 v.s. ISO-8859-1
 - Windows CP1250 v.s. ISO-8859-2
- Also other "standards": IBM, Apple, …
- Problems with Web pages

8-bit character sets (ISO 8859, Windows)

- advantages to ASCII:
 - we can directly write the characters for national alphabets (slovenščina)
- disadvantages:
 - we cannot write multilingual texts in the same character set
 - o confusion due to competing character sets
 - no coverage for East Asian languages or more complex characters: punctuation, math operators, diacritics, historical characters, Klingon...
 - the file gives no indication which character set it uses:

© Global publishing ~ Ž Global publishing

The final solution

- Need a character set that would be universal, i.e. would contain all the world's characters
- Must be well documented and open
- Has to be consensually developed and maintained
- Still needs some room for "private" characters

Unicode

1991 – Unicode Consortium: <u>http://www.unicode.org/</u> Unicode Standard / ISO 10646 "Universal Character Set"

- The most recent major revision is Unicode 6.0. (2011): 109,000 characters, 93 scripts
- code charts for visual reference + reference data files
- encoding methodology, character properties, rules for normalization & decomposition, collation, rendering, bidirectional display: complex!
- As yet unrealised ambition: completely replace other character sets

Visual reference: Character code charts

Latin Extended-A

Position	Decimal	Name	Appearance
0x0100	256	LATIN CAPITAL LETTER A WITH MACRON	Ā
0x0101	257	LATIN SMALL LETTER A WITH MACRON	ā
0x0102	258	LATIN CAPITAL LETTER A WITH BREVE	Ă
0x0103	259	LATIN SMALL LETTER A WITH BREVE	ă
0x0104	260	LATIN CAPITAL LETTER A WITH OGONEK	Ą
0x0105	261	LATIN SMALL LETTER A WITH OGONEK	ą
0x0106	262	LATIN CAPITAL LETTER C WITH ACUTE	Ć
0x0107	263	LATIN SMALL LETTER C WITH ACUTE	ć
0x0108	264	LATIN CAPITAL LETTER C WITH CIRCUMFLEX	Ĉ
0x0109	265	LATIN SMALL LETTER C WITH CIRCUMFLEX	ĉ
0x010A	266	LATIN CAPITAL LETTER C WITH DOT ABOVE	Ċ
0x010B	267	LATIN SMALL LETTER C WITH DOT ABOVE	ċ
0x010C	268	LATIN CAPITAL LETTER C WITH CARON	Č
0x010D	269	LATIN SMALL LETTER C WITH CARON	č

Unicode definitions for IPA

内	🚈 Adobe Acrobat Professional - [U0250[1].pdf]											
艿	<u>File E</u> dit <u>V</u> iew <u>D</u> ocument	Comme	nts Tools <u>A</u> dvanced <u>W</u> indow <u>H</u> elp									
1 7	Treate PDF 🗸 🦉 Comment & Markup 🗸 Send for Review 🗸 🔒 Secure 🗸 🆉 Sign 🗸 📑 Forms 🔹 🥮 Note Tool 🕂 Te											
2)] 🗈 Select 📷 🛛 🔍	•										
	0250	e	LATIN SMALL LETTER TURNED A									
arks			 low central unrounded vowel 									
, щ	0251	α	LATIN SMALL LETTER ALPHA									
B			= LATIN SMALL LETTER SCRIPT A									
8			 low back unrounded vowel 									
Page			\rightarrow 03B1 α greek small letter alpha									
0	0252	υ	LATIN SMALL LETTER TURNED ALPHA									
nature			 low back rounded vowel 									
Sig	0253	6	LATIN SMALL LETTER B WITH HOOK									
			• implosive bilabial stop									
			• Pan-Nigerian alphabet									
			\rightarrow 0181 B latin capital letter b with hook									
	0254	2	LATIN SMALL LETTER OPEN O									
	0201	0	• typographically a turned c									
			• lower mid back rounded yours									
			\rightarrow 0186 J latin capital letter open o									

Reference data

- Unicode Character Database
- e.g. <u>http://www.unicode.org/Public/UNIDATA/NamesList.txt</u>
- CSV files, XML, PDFs
- Various software uses this information: Java classes, Perl modules, e.g.
 - m/[[:upper:][:punct:]]/;
 matches any upper case letter or punctuation symbol
 - charnames::viacode(4532)
 returns: LATIN CAPITAL LETTER U WITH DOUBLE
 GRAVE

Unicode and diacritics

- many letters are available with diacritics as individual characters: áâãäååääåå
- but diacritics also exist as combining characters (combining diacritical marks)
- eg.: a + ^ + _ = a_.
- although problems with display of complex combinations, e.g. a + ^+ ° = a²
- solved by specialised fonts



Private Use Area

- Not all characters are (or could be) included in Unicode
- Unicode allows for addition of new characters, but only after an extended process
- For extra characters, a Private Use Area (PUA) is designated
- Fonts are free to use PUA, but with the understanding that these characters are not portable
- Example use: <u>Freising Manuscripts</u>

Unicode planes & BMP

v·d·e Unicode planes															Uni	code pl	anes and code point (character) ranges [hid						
Basic																	Supplementary						
0000-FFFF 10000-1FFFF										10	000-1	1FFFI	F		20000-2FFFF		30000-DFFFF	E0000-EFFFF	F0000-10FFFF				
Plane 0:								Plane 1:								Plane 2:	Plane 2: Planes			Planes 15-16:			
Basic Multilingual Plane									Supplementary Multilingual Plane							Plane	Supplementary Ideographic	Plane	Unassigned	Supplementary Special-purpose Plane	Private Use Area		
				BM	Р								SM	Р			SIP		-	SSP	S PUA A/B		
00 10 20 30 40 50 60 70 80 90 80 80 60	01 11 21 31 41 51 61 71 81 91 81 81	02 12 22 32 42 52 62 72 82 82 92 82 82 82 82	03 13 23 33 43 53 63 73 83 83 93 83 83 83	04 14 24 34 44 54 54 64 74 84 94 84 84	0 <mark>5</mark> 15 35 45 55 65 75 85 95 85 85	06 16 26 36 46 56 76 86 96 96 86 96	07 17 27 37 47 57 67 77 87 87 97 87 87	08 18 28 38 48 58 68 98 98 98 88 98 88	 09 19 29 39 49 59 69 79 89 99 A9 B9 C 	0A 1A 2A 3A 4A 5A 6A 7A 8A 9A 9A 8A 6A	08 18 38 38 58 68 78 88 98 88 98 88	0C 1C 2C 3C 4C 5C 6C 7C 8C 9C 8C 6C	0D 1C 2D 3D 4D 5D 6D 7D 8D 9D 4D 6D	0E 1E 2E 3E 5E 6E 8E 9E 8E 8E	 OF 2F 2F 4F 6F 8F 9F AF 6F 6F 		Latin scripts and symbols Linguistic scripts Other European scripts African scripts Middle Eastern and Southwest Asian scripts Central Asian scripts South Asian scripts Southeast Asian scripts East Asian scripts Unified CJK Han American scripts Symbols Diacritics UTF-16 surrogates and private use	Bas 000 (65 E00	sic multiling)0 - FFFF ,535 chara)0–F8FF -	gual plane acters) PUA			
D0 E0	D1 E1	D2 E2	D3 E3	D4 E4	D5 E5	D6 E6	D7 E7	D8 E8	D9 E9	DA EA	DB EB	DC EC	DD ED	DE EE	DF EF		Miscellaneous characters Unallocated code points						
FO	F1	F2	FЗ	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FΕ	FF								

Encoding Unicode

- ISO 8859, Windows charaters sets: 8bit, therefore limited to 256 characters
- Trivial mapping: one char = one byte
- But Unicode codepoints can be huge
- Necessary to use several bytes to encode one character
- All Unicode codepoints (numbers) fit into 4 bytes
- But it is in general very wasteful to use 4 bytes for one char
- Is there any better way to do it? Yes, several..

Unicode Transformation Format

UTF defines how to map codepoints to bytes (bits), which are then stored or transmitted

- UTF-32
 - 1 character = always 4 bytes
- UTF-16
 - 1 character = 2 bytes in basic multilingual plane
- UTF-8
 - if char in ASCII, then in 1 byte (compatibility!)
 - o otherwise 1-6 bytes for 1 char
 - cunning system, where not all byte sequences are valid (so, won't mix with 8 bit encodings)

Back to ASCII

ASCII is sometimes still the safest:

- o problems with input and display of chars
- o data transfer (e-mail)

Recoding to ASCII:

- e-mail: MIME standard
- HTML and XML character entities: š = Š = š

Defining the character set of a document

HTML:

<HTML> <HEAD> <TITLE>Recept za ribano kašo</TITLE> <META http-equiv="Content-Type" content="text/html; charset=ISO-8859-2"> </HEAD> <BODY>

XML:

. . .

. . .

<?xml version="1.0" encoding="utf-8"?> <recept> <naslov>Recept za ribano kašo</naslov>

- Some valid character sets:
 - o utf-8, iso-8859-X, us-ascii

i18n, l10n

- Internationalisation and localisation: enabling programs to work with different languages (and cultures)
- E.g. language of program messages and help; keyboard layout; date and number format
- <u>CLDR</u> Unicode Common Locale Data Repository
- For language resources:
 - o collating sequence: *a,b,c,č,d*, not *a,b,c,d,...č*
- Unix: system variables regulate which locale is selected, e.g. LC_COLLATE = sl_SI.UTF-8

Case study: Cleaning Gigafida

- 1,000,000,000 word tokens
- Unicode + XML + TEI
- Character profile: 1200
- Forbidden chars: 500
- Excel
- Character normalisation

Fixing chars in Gigafida

- Hyphens: \$s=~s/[\x{0336}\x{0096}\x{2010}]/-/g;
- Spaces: \$s=~s/[\x{00A0}\x{2002}\x{2008}\x{2009}\x{202F}]//g;
- Digraphs: \$s=~s/ffi/ffi/g; \$s=~s/ffl/ffl/g; \$s=~s/ff/ff/g; \$s=~s/fl/fl/g;
- Non-spacing diacritics:
 \$s=~s/ñ/ñ/g;
 \$s=~s/č/č/g;
 \$s=~s/š/š/g;
- Entities: \$s=~s/&/&/g; \$s=~s/ / /g; \$s=~s/©/©/g;





- This is the phonetic spelling of the Slovene word "čmrlj" (bumbelbee)
- 1. Write the characters in Word
- 2. Find which Unicode characters these are <u>http://www.unicode.org/charts/</u>