

# Zapis znakov in uporaba korpusov

Tomaž Erjavec  
Korpusno jezikoslovje  
UNG  
2008/2009  
9. 1. 2009

---

---

---

---

---

---

---

---

## Pregled predavanja

1. Zapis znakov v računalnikih
2. Primeri študij uporabe korpusov
3. Diskusija seminarских nalog

---

---

---

---

---

---

---

---

## Kodiranje znakov

- Digitalni računalniki shranjujejo podatke kot (binarne) številke
- Ne obstaja vnaprej dana povezava med temi številkami in znaki (abecede)
- Če ni konvencij za preslikavo ali jih je preveč → kaos
- Standardi in pol-standardi: ASCII, ISO 8859, (Windows, Mac), Unicode

---

---

---

---

---

---

---

---

## Osnovni pojmi I.

### znak (character)

- o abstrakten pojem (An „A“ is something like a Platonic entity: it is the idea of an „A“ and not the „A“ itself)
- o sam po sebi znak nima preslikave v številko ali določene grafične podobe
- o ponavadi je opisno definiran, npr. „grška črka mala alfa“, grafična podoba pa podana samo kot vodilo, „α“

---

---

---

---

---

---

---

---

## Osnovni pojmi II.

- **nabor znakov (character set)**
  - o množica znakov
  - o vsakemu znaku pripisana njegova koda
- **koda znaka (character code)**
  - o 1-1 relacija med znakom iz nekega nabora znakov in številko, npr. A = 26, B = 27, ...
  - o Pozor!  
Kode znakov se dostikrat zapisujejo šestnajstiškem sistemu:  
0 → 0, 1 → 1, 2 → 2, ... 9 → 9, 10 → A, 11 → B, ..., 15 → F, 16 → 10, 17 → 11, ..., 254 → FE, 255 → FF, 266 → 100

---

---

---

---

---

---

---

---

## Primer: nabor znakov ASCII

Below is the standard ASCII characters.

Dec	Char	Dec	Char	Dec	Char	Dec	Char	Dec	Char	Dec	Char
33	!	49	1	65	A	81	Q	97	a	113	q
34	"	50	2	66	B	82	R	98	b	114	r
35	#	51	3	67	C	83	S	99	c	115	s
36	\$	52	4	68	D	84	T	100	d	116	t
37	%	53	5	69	E	85	U	101	e	117	u
38	&	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	H	88	X	104	h	120	x
41	)	57	9	73	I	89	Y	105	i	121	y
42	*	58	:	74	J	90	Z	106	j	122	z
43	+	59	;	75	K	91	[	107	k	123	[
44	,	60	=	76	L	92	\	108	l	124	]
45	-	61	>	77	M	93	]	109	m	125	]
46	.	62	>	78	N	94	^	110	n	126	^
47	/	63	?	79	O	95	_	111	o	127	_
48	0	64	@	80	P	96	`	112	p		

npr.  
v naboru znakov ASCII  
ima znak mali latinični a  
kodo znaka 97

---

---

---

---

---

---

---

---

## Osnovni pojmi III.

- **pismenka (glyph)**
  - grafična predstavitev znaka
  - enemu znaku lahko ustreza več kot ena pismenka  
npr. znak "veliki latinični A" ↔ pismenke A, Ἀ, Α
  - redko tudi eni pismenki ustreza več znakov  
npr. pismenka P ↔ znaki "veliki latinični P", "veliki cirilični R", "veliki grški Ro")
- **font**
  - nabor pismenk (za nek nabor znakov):  
A, B, C, Ć, D, ...
  - včasih font ne pokriva celotnega nabora znakov!

---

---

---

---

---

---

---

---

## Nekateri nabori znakov

- ASCII - najstarejši, vsebuje samo črke ameriške abecede + ločila, številke
- Družina naborov znakov ISO 8879
- Družina naborov Windows
- Unicode

---

---

---

---

---

---

---

---

## ASCII

- American Standard Code for Information Interchange (1950')
- 7-bitni zapis znakov: kode znakov 0-127
- 0-31 - kontrolni znaki + znaki za formatiranje: Esc, Line Feed, tabulator, presledek, ...
- 32-126 – ločila in posebni znaki, številke, velike in male angleške črke :  
!"#\$%&'()\*+,-./0123456789:;  
<=>?@ABCDEFGHIJKLMN  
OPQRSTUVWXYZ[\]^\_`abcdefghijklmnopqrstuvwxyz{|}~

---

---

---

---

---

---

---

---

## Družina naborov znakov ISO 8859

- potreba po dodatnih znakih za nacionalne (evropske) pisave:
  - v 80's se pojavljajo novi nabori znakov
  - obsegajo ASCII kot podmnožico
- International Standard Organisation izda kodne nabore za posamezne skupine (evropskih) jezikov: družina standardov **ISO 8859**
- ISO 8859-1 (ISO Latin 1) – zahodnoevropski jeziki

¡ ¢ £ ¤ ¥ ¦ § ¨ © 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 ÿ

---

---

---

---

---

---

---

---

## Nabori znakov za ne-zahodno evropske jezike

- za slovenščino in ostale srednje- in vzhodno evropske (latinične) jezike - anarhija:
  - ISO 8859-2 (ISO Latin 2)
  - Windows CP1250 (grrr!)
  - lastni „standardi“: IBM, Apple, ...

---

---

---

---

---

---

---

---

## ISO 8859-2 (zgornja polovica)

	NSBP	À	Á	Â	Ã	Ä	Å	Ā	Ą	Ć	Č	Š	Ś	Ŝ	Ț	Ž	ŠHY	Ž	Ž
A-	00A0	0104	0108	0141	00A4	013D	015A	00A7	00A8	0160	015E	0164	0179	00AD	017D	017E			
	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175			
	°	ª	¸	¸	¸	¸	¸	¸	¸	¸	¸	¸	¸	¸	¸	¸			
B-	0080	0105	0208	0142	0084	013E	015B	02C7	0088	0161	015F	0165	017A	0200	017F	017C			
	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191			
	Ħ	Á	Â	Ã	Ä	Å	Ā	Ą	Ć	Č	Š	Ś	Ŝ	Ț	Ž	Š	Ÿ	Ž	Ž
C-	0114	00C1	00C2	0102	00C4	0139	0106	00C7	010C	00C9	0118	00C8	011A	00C0	00CE	010E			
	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207			
	Đ	đ	Ħ	Ó	Ô	Õ	Ö	×	Ř	Ú	Û	Ü	Ý	Ť	Ř				
D-	0110	0143	0147	0203	0204	0150	0205	0207	0158	016E	020A	0170	020C	0200	0162	020F			
	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223			
	ƒ	á	â	ã	ä	å	ā	ą	ć	č	š	ś	ŝ	ț	ž	š	ÿ	ž	ž
E-	0155	00E1	00E2	0103	00E4	013A	0107	00E7	010D	00E9	0119	00E8	011B	00E0	00E4	010F			
	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239			
	đ	ñ	ñ	ó	ô	õ	ö	×	ř	ú	û	ü	ý	ť	ř				
F-	0111	0144	0148	0201	0204	0151	0206	0207	0159	016F	020B	0171	020E	0200	0163	020F			
	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255			
	-0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-A	-B	-C	-D	-E	-F			

---

---

---

---

---

---

---

---

## 8-bitni nabori (ISO 8859, Windows)

- prednosti:
  - lahko zapišemo znake posameznih nacionalnih abeced (slovenščina)
- slabosti:
  - v istem kodnem naboru ne moremo zapisati večjezičnih besedil
  - zmeda zaradi večih kodnih naborov, ki pokrivajo iste jezike
  - ni pokrītja npr. za vzhodno-azijske jezike ali bolj zahtevne znake: ločila, matematični simboli, naglasna znamenja, ...
  - datoteka ne vsebuje podatka v katerem kodnem naboru je vsebina:  
© Global publishing ~ Ž Global publishing

---

---

---

---

---

---

---

---

## Unikod I.

Unikod (Unicode oz. ISO 10646)

1991 – Unicode Consortium: <http://www.unicode.org/>

- definira univerzalni nabor znakov
- vsebuje 30 svetovnih abeced, ki pokrivajo več sto jezikov, definiranih približno 40.000 znakov
- ...arabščina, sanskrt, kitajščina, japonsščina, korješčina, ...
- tudi zgodovinske pisave, ločila, matematični simboli, naglasna znamenja, ...
- Unikod razdeli znake v „bloke“  
npr. Basic Latin, Latin-1 Supplement, Latin Extended-A, Latin Extended-B, IPA Extensions, Combining Diacritical Marks, Greek, Cyrillic, ...

---

---

---

---

---

---

---

---

## Slovenske črke

### 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	Č
0x010B	267	LATIN SMALL LETTER C WITH DOT ABOVE	č
0x010C	268	LATIN CAPITAL LETTER C WITH CARON	Ć
0x010D	269	LATIN SMALL LETTER C WITH CARON	ć

---

---

---

---

---

---

---

---



## Nazaj v ASCII

ASCII je včasih še vedno edini varen zapis:

- o problemi pri vnosu ali izpisu znakov
- o problemi pri prenosu podatkov (elektronska pošta)

Prekodiranje v ASCII:

- elektronska pošta - standard MIME
- SGML (HTML) in XML - entitete za znake, s kodnimi mesti iz Unikoda  
npr. `&#353;`; = `&#x160;`; = `š`

---

---

---

---

---

---

---

---

## Določanje nabora znakov

### 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"?>
<receipt>
  <naslov>Recept za ribano kašo</naslov>
...
```

- Nekateri dovoljeni nabori znakov:
  - o `utf-8`, `iso-8859-2`, `us-ascii`

---

---

---

---

---

---

---

---

## Vaje iz zapisa znakov: Word

Pri tej vaji uporabimo Word, da spoznamo razlike med kodnimi nabori.

Besedilo:

*Mačka, miška in žolna so šli na izlet v Črnomelj, nato prav počasi v Šujico, na koncu so pa pristali v Žužemberku, kjer so srečali čmrjca. »Kako dolgo smo hodili?« je mu je potožila miška. Mačka pa mu je rekla »Cuj čmrjca, koliko samoglasnikov je pravzaprav v tvojem imenu?«*

- Odpri novo datoteko v Wordu, in gornje besedilo prilepi vanjo. Shrani jo kot besedilo (.txt), v kodnem naboru ISO-8859-2 (»Central European ISO«). Za katere znake javi Word, da bodo nepravilno shranjeni? Zakaj?
- Datoteko zapremo, nato .txt ponovno odpre v urejevalniku Word; katere kodne napore (od tih, ki smo jih omenili v predavanju) ponudi Word? Kaj se zgodi, če besedilo odpremo v prizeti kodni tabeli za Windows? Datoteko zapremo, ne da bi jo shranili.
- Datoteko odpremo in ponovno shranimo kot .txt, tokrat v kodnem naboru UTF-8 (»UTF«) in jo zapremo.
- Besedilo ponovno odpremo. Kako izgleda besedilo, če ga odpremo v kodnem naboru ISO-8859-2? Zapremo, ne da bi shranili.
- Besedilo pravilno odpremo (torej v UTF-8), nato pa shranimo v UTF-16 (»Unicode«), datoteko spet zapremo, in ponovno odpremo. Kaj se zgodi, če jo odpremo s katerim od 8-bitnih kodnih naborov? Zakaj?

---

---

---

---

---

---

---

---

## Vaje iz zapisa znakov: Unikod

Učimo se poiskati želeni znak v urejevalniku Word (vstavi simbol) in na spletnih straneh Unikoda, na <http://www.unicode.org/charts/>.

Fonetično želimo napisati »čmrļj«, ta ko da je med »mr« in »rļ« znak za polglasnik, namesto »č« pa znak iz mednarodne fonetične abecede IPA. Ker je prvi polglasnik naglašen, mu dodajte še ostrivec, torej:

ʈm̥éɾɯj

1. Najdi potrebne znake med »Vstavi znak« v urejevalniku Word. Pozor: v okenu levo zgoraj izberi font, ki podpira te fonetične znake, verjetno Arial Unicode MS. Kako so ponujeni znaki znaki urejeni? Kako se znaka imenujeta?
2. Znaka nato poišči še preko spletne strani Unicode. Kateri sta njuni kodni mesti?
3. V katerem kodnem naboru lahko shranimo to datoteko?

---

---

---

---

---

---

---

---

---

---

## Študije iz korpusnega jezikoslovja

- Knjiga  
Vojko Gorjanc: *Uvod v korpusno jezikoslovje*.  
Domžale: Izolit, 2005.  
163 str.
- prva knjiga o korpusnem jezikoslovju pri nas
- Predstavitve korpusov in korpusnega jezikoslovja
- Četrty del:  
*Korpusni opisi slovenskega jezika*  
(uporaba korpusa Fida)
- Dva primera:
  - spreminjanje jezika - terminologija
  - zajem semantičnih povezav med besedami

---

---

---

---

---

---

---

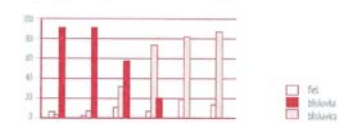
---

---

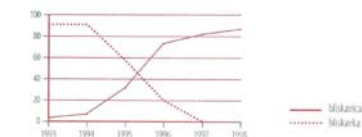
---

## Spreminjanje jezika

Razmerje med palminosovj/fel – bliskavka – bliskavica (%)



Razmerje med različnimi bliskavka in bliskavica (%)




---

---

---

---

---

---

---

---

---

---



## Semantične relacije

- enak, nasproten pomen:
  - sinonimi (neodvisnost ↔ samostojnost)
  - antonimi (lahek ↔ težek)
- nad, podpomen:
  - hipernimi (ptič → vrabec)
  - hiponimi (vrabec → ptič)
- del in celota:
  - holonimi (avto → vrata)
  - meronimi (vrata → avto)
- pomembno za izdelavo tezavrov in aplikacije v jezikovnih tehnologijah

---

---

---

---

---

---

---

---

---

---

## Uporaba vzorcev za odkrivanje semantičnih relacij

- Iskalni pogoj je niz, ki v besedilu vzpostavlja semantično relacijo, ki nas zanima

Pomenski označevalci glede na uspešnost zajetja pod-/nadpomenskosti

označevalec	abs. pog./pog. označ. podpomenskosti		uspešnost	
	A	B	A	B
je vrsta	90/19	13/4	0,21	0,31
je * vrsta	156/25	20/0	0,15	0
prištevamo med	12/9	5/5	0,75	1
sodi med	182/45	79/39	0,24	0,49
spada med	168/34	65/30	0,20	0,46
sodi v družino	13/4	0/0	0,30	0
uvrščamo med	60/22	7/3	0,36	0,43

---

---

---

---

---

---

---

---

---

---

## Sopomenke: imenovan tudi

Del konkludirnega niza iskalnega pogoja imenovan tudi/imenujemo tudi

opisan neposreden način odstiti delček oksid,	imenovan tudi	svetljali plus, znački katerega postane človek omamljen
Vitamin B1,	imenovan tudi	tansin, je sestava najbolj znan med listni vitamini
Vitamin B2,	imenovan tudi	riboflavin, je prazgavov deličen najstari požarnosti
Štepijo dostopa do kode	imenujemo tudi	človek procedur
numeričazbirne manjše, ta samostanski lučič,	imenovan tudi	ženski prvi legree, je v petletnosti
prinosnost moškega žuljenja, To skupno zveznik	imenujemo tudi	ehinoderme,
Je kdaj stisnjati, da žemlje	imenujemo tudi	moda planet?
Zato spletno strani	imenujemo tudi	HTML dokument, V osnovi je HTML dokument
Veljavno odelnost	imenujemo tudi	naznodelna sestnost, to je izraz, s katerim
karte med 1 : 10 000 in 1 : 5 000	imenujemo tudi	detajne geotiske karte, karte še večjih merilih
Odločanje hitri elektroni	imenujemo tudi	sevanje Zarkon II, ves poganje
Srčni v notnem agregatnem stanju	imenujemo tudi	trdišae, tudi je jih nas zanima, kako se
sta bogata predvsem z žlezovimi spojinami, jih	imenujemo tudi	žlezovata rta, ženski.

---

---

---

---

---

---

---

---

---

---

## “znan kot”

Del konordančnega niza iskalnega pojopa znan kot

<p>uživamo v vsakem večerju klasi mestjenski žarči, je proci svetlobo odda jernih galaksij. Začni pesni, ki je dases kot jih je imel slasti člen lambda (odl) esopkuje pred približno 17 milijardami let. Ti dogadjek je sočasnjen z zvezd, ki tvorijo anemzem, v obsevatarijo Maurt Palomar, ki je pamnje poslal antiki astronom je: 30 Al Rata v, v Evropi je delo Hertzsprung prikazal grafično. Graf, sedaj antihonke, kot epe, štali tip strepskeku. A rad n'ni in tuž povečane. Pajav je</p>	<p><b>znan kot</b> senjenduminščerica <b>znan kot</b> nedeč gremik, je abko iz spremembe barvne <b>znan kot</b> kuzmoišta korenanta, ki ga je <b>znan kot</b> "Vilki bum" (tuži "veliki pov") <b>znan kot</b> Juhova krata. Dve od teh zvezd imata lučiv <b>znan kot</b> Halov teleskop. Za postavitev tega teleskopa <b>znan kot</b> Almagestus. Dejav je v zadnji četrsti <b>znan kot</b> Hertzsprung-Russlowv diagram, prikazuje zvezo <b>znan kot</b> vrstnosna bitenja, za katero je možno, da je <b>znan kot</b> fara mozgana in je pogosto opazovan</p>
--	--

---

---

---

---

---

---

---

---

---

---

---

---

## Samostalnik (Samostalnik)

Del prečiščenega konordančnega niza iskalnega pojopa Sam (San) v podkapitulu naravnostne vede (Cobis)

<p>Enocelčni plazmodiji razgajajo rlela krma proces strjevanja ledaj, ko vmešamo v smolno sredstvo za vrste lesa, papir, kovine, stekla, acetata v vodi in kisk. Votlik se natira na negativni negativni električni (katodi), kisk pa na pozitivni lastnosti dimnih zavese spremenijo na optičnih pojavih zvišamo, odje začne zgovorati in pri tem nastane ogljen pišč bogov in kraljev, ki se v čino osrednjega dela ali telesa nevrons, več kažih, vejastih dreh na zemeljski ekvator (palutrik) ter na oba je pri svojih operacijah uporabljal karbolno Znanstveniki menijo, da ta 15-cerimetska figura prikazuje del motaja, pač pa le kot sestavni del modernega tega ima sodobna kopija kar 8-krat večji delovni pomnilnik (RAM) in 4-krat večji trajni različne dejavnostne molekule deoksiribonudeinle-ameriški postopek, kako ta stanje pekeloboliti nutirjev lastnosti sta hitro učilovanje in visoka močnja postopoma stišja, potem ko so jih prepajali s polistilten sestava je odvisna od matične karnine, odhajanja</p>	<p><b>telesca</b> (retrocite) <b>strjevanje</b> (tridilec), <b>celuloza</b> (celulozi), <b>elektrodi</b> (katodi), <b>elektrodi</b> (anodi), <b>disperzije</b> (razpisevanja) <b>prah</b> (čad), <b>žentive</b> (spamifidi) <b>izrasikov</b> (dendritov) <b>soja</b> (trčaja), <b>kislino</b> (fenol), <b>vrača</b> (Samana), <b>bloka</b> (ohišja), <b>pomnilnik</b> (RAM) <b>pomnilnik</b> (ROM) <b>kislina</b> (DNA), <b>hidroksoid</b> (veg), <b>strepnesci</b> (toksičnost), <b>glikolom</b> (PEG), <b>prsti</b> (erocije)</p>	<p>in ob tem povezočaje silne napade akivno steklo (PIMMA) in nekatere kisk pa na pozitivni Pri gorih celici je postopek in absorpcije (vsrkavanja) svetlobe ki stihiti zraklu kot stihni oboki. Črne barve v mestni prelovo v praviljuna lanja. In le enoga obsega izrasita (akivna), severnega in južnega. Če naprej da je papreči zastripitve. Kacenje so ki vstopa v drugo stanje. Zaradi čim lažjega in 4-krat večji trajni pomnilnik (ROM) od originala. Sicer pa je vse sestavne dele Poskus je trajal nekaj mesecev. ki je za izdelavo mila neprimerno boljji so brez barve, vonja in okusa. v vodi tojtno polimerno snaki, katere in živih bitij, ki sodekujejo pri nastajanju</p>
---	---	---

---

---

---

---

---

---

---

---

---

---

---

---

## Meronomi

- katere vzorce bi uporabili za iskanje meronomov - holonomov?

---

---

---

---

---

---

---

---

---

---

---

---

### III. Korpusi v leksikografiji: terminološki sovarček

Vaja

1. specializiran korpus s področja primerjamo z referenčnim korpusom, da dobimo ključne besede
2. ključnim besedam pogledamo kolokacije, da dobimo večbesedne termine (iztočnice)
3. iztočnicam preko konkordanc najdemo podpomene (če so) in poiščemo dobre primere uporabe

---

---

---

---

---

---

---

---