

[Zapis znakov in uporaba korpusov]

Korpusno jezikoslovje / Jezikovne tehnologije
UNG
2009/2010
10. 5. 2010

[Pregled predavanja]

1. Zapis znakov v računalnikih
2. Primeri uporabe korpusov

[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*)

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

[Osnovni pojmi II.]

■ **nabor znakov** (*character set*)

- množica znakov
- vsakemu znaku je pripisana njegova številčna koda

■ **koda znaka** (*character code*)

- 1-1 relacija med znakom iz nekega nabora znakov in številko, npr. A = 26, B = 27, ...
- 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, Ä, Å
 - 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 O P
QRSTUVWXYZ[\]^_`abcdefghijklmnop
ijklmnopqrstuvwxyz{|}~

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

Naberi 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	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï	Š	Š	Ţ	Ž	ŠHY	Ž	Ž
A-	00A0	0104	02D8	0141	00A4	013D	015A	00A7	00A8	0160	015E	0164	0179	00AD	017D	017B								
	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175								
B-	0080	0105	02D8	0142	00B4	013E	015B	02C7	00B8	0161	015F	0165	017A	02D0	017E	017C								
	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191								
C-	0114	00C1	00C2	0102	00C4	0139	0106	00C7	010C	00C9	0118	00CB	011A	00CD	00CE	010E								
	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207								
D-	0110	0143	0147	00D3	00D4	0150	00D6	00D7	0158	0166	00D4	0170	00DC	00DD	0162	00DF								
	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223								
E-	0155	00E1	00E2	0103	00E4	013A	0107	00E7	0100	00E9	0119	00EB	011B	00ED	00EE	010F								
	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239								
F-	0111	0144	0148	00F3	00F4	0151	00F6	00F7	0159	016F	00FA	0171	00FC	00FD	0163	02D9								
	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 pokritja 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	Ď
0x010D	269	LATIN SMALL LETTER C WITH CARON	ď

Nazaj v ASCII

ASCII je včasih še vedno edini varen zapis:

- če so problemi pri vnosu ali izpisu znakov
- če so 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. š = Š = 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"?>
<recept>
  <naslov>Recept za ribano kašo</naslov>
...
```

Nekateri dovoljeni nabori znakov:

- 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 čmrlja. »Kako dolgo smo hodili!« je mu je potožila miška. Mačka pa mu je rekla »Cuj čmrlj, 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 nabore (od tih, ki smo jih omenili v predavanju) ponudi Word? Kaj se zgodi, če besedilo odpremo v privzeti 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 »rl« znak za polglasnik, namesto »č« pa znak iz mednarodne fonetične abecede IPA. Ker je prvi polglasnik naglašen, mu dodajte še ostrivec, torej:

t̪m̩érl̩j

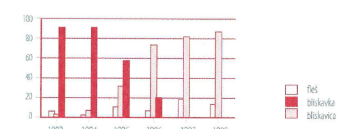
1. Najdi potrebne znake med »Vstavi znak« v urejevalniku Word. Pozor: v okencu levo zgoraj izberi font, ki podpira te fonetične znake, verjetno Arial Unicode MS. Kako so ponujeni znaki znaki urejeni? Kako se znaka imenujeta?
2. Znaka nato poišči še preko spletne strani Unicode. Kateri sta njuni kodni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

Študije iz korpusnega jezikoslovja

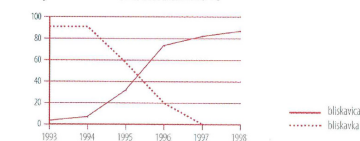
- Knjiga
Vojko Gorjanc: *Uvod v korpusno jezikoslovje*.
Domžale: Izolit, 2005.
163 str.
- prva knjiga o korpusnem jezikoslovju pri nas
- Predstavitev korpusov in korpusnega jezikoslovja
- Četrti del:
Korpusni opisi slovenskega jezika
(uporaba korpusa Fida)
- Dva primera:
 - spreminjanje jezika - terminologija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

Spreminjanje jezika

Razmerje med poimenovanji flis – bliskavica – bliskavica (%)



Razmerje med različnima bliskavica 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 ⁹ vrsta	156/25	20/0	0,15	0
prštejemo 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 kanaldančnega niza iskalnega pogoja imenovan tudi/imenujemo tudi

opisan neposreden način odnosa duškov oksid,	imenovan tudi	snegalni piln, zaradi katerega postane števki omuljen
Vitamin D1,	imenovan tudi	tamin, je verjetno najbolj znan med šestimi vitamini
Vitamin B2,	imenovan tudi	riboflavin, je praznecar deležen najmanj pozornosti
Stopnja dostopa do kode	imenujemo tudi	dosaj procedure
numeričizacijave maroge, ta samostarski lučzar,	imenovan tudi	žicopajni legren, je v preteklosti
posanost moškega življenja, to skupino zveinrk	imenujemo tudi	elnicodrome
že kdaj slišali-a), da žemljo	imenujemo tudi	modi planet?
Zate spletne strani	imenujemo tudi	HTML dokumenti. V osnovi je HTML dokument
Vetpustno osednost	imenujemo tudi	razcepjeno osednost; to je izraz, s katerim
karte naved 1 : 10 000 in 1 : 5 000	imenujemo tudi	delajine gredloške karte, karte v števjuh mestih
Odklajanje bitnih elektronov	imenujemo tudi	sevanje žarkov B, ves pojav pa
Sonci v tistem agrogativnem stanju	imenujemo tudi	tudi:ar. Tudi pr njih naznaniat, kdaj se
da bogata predvsem z železovimi spojinami, jih	imenujemo tudi	železnata tla, ferisol

“znan kot”

Del konkanančnega niza iskalnega pogoja znan kot

vdvajamo z visokomergijskimi rentgenskimi žarki, je proces svetlobo oddajajo galaksij. Zaradi pojave, ki je danes kot jih je imel slavn člen lambda (dolj) esopkuje pred približno 17 milijardami let. Ta dogodek je sesavljen iz zvezd, ki tvorijo acetozem, v observatoriju Mount Palomar, ki je pozneje poslal arabski astronomi je bil Al Bata'i, v Evropi je delo Hertzsprung prikazal grafično. Graf, sedaj amblotike, kot npr. fatalni tip stopotokula A, nad njim in tudi povečane. Pojav je

znan kot rentgenoluminescenca
znan kot rdeči premik, je ahko iz spremembe barvne
znan kot kozmoloska konstanta, ki ga je
znan kot "Veliki bum" (tudi "Veliki po")
znan kot Jobova kista. Dve od teh zvezd imata tudi
znan kot Halov teleskop. Za postavitev tega teleskopa
znan kot Albatrogus. Delovja je v zadnji četrtini
znan kot Hertzsprung-Russellov diagram, prikazuje zvezd
znan kot srčnatostna belterija, za katero je značilno, da je
znan kot fota in izvaja in je pogosto opazovan

Samostalnik (Samostalnik)

Del predščenega konkanančnega niza iskalnega pogoja Sam (Sam) v podkarpisu naravoslovne vede (Cobiss)

Enocelčni plurimediji naučajo rdeča krvna proces stojičanja tedaj, ko vmerlamo v umno sredstvo za vrste lesa, papir, kovine, steklo, acetatno v vodiki in kisl. Modik se nakiba na negativni negativni elektroni (katalizi), kolik pa na pozitivni lastnosti demitk zavre temeljijo na capičnih pojavi zvišamo, olje začne zlojavati in pri tem nastane ogljen pitile bogov in kraljev, ki se v času eurednjega dela ali tekusa nevarna, več kraljev, večjih dreh na zemeljski ekvator (palatnik) ter na oba je pri večjih operacijah uporabljal karbolno Znanstveniki menijo, da ta 35-centimetrska figura prikazuje del motorja, pač pa le kot sestavi deli motornega tega ima sodobna kopija kar 4-krat večji delovni pomnilnik (RAM) in 4-krat večji trajni razpisilni delovni mehanizem deokribonukleinske karnijski postopek, kako bi danes pridobili izvirni razpisilni sta hitro uveljavljene in visoko stopnja postopoma tistega, potem ko so jih prepeljali s polietilen sestava je odbrana od matične karmine, odlašanja

telesca (entrotite)
strijevanje (indilec)
celuloza (celuloid)
elektroni (kanodi)
elektrodi (anodi)
disperzije (razprševanje)
prah (kadi)
ženskiue (opamlad)
toraciklov (dendritov)
pola (tečaja),
kolino (feneti),
vrača (Samama),
bloka (ohlija)
pomnilnik (RAM)
pomnilnik (ROM)
kislina (DNA),
hidrokod (tegl),
stropenosti (kaskošnosti)
gliholoni (PEG),
prsti (terzije)

in ob tem povzročajo silne napade
akrilno steklo (PMMA) in nekatere
kislj pa na pozitroni
Pri gorilni celici je postopek
in absorpcije (vsrkavanja) svetlobe
ki tedaj zraku kot oljni oblik črne barve
v enostoj predelje v pravičnu boga,
in le enega odsega (ozračja) (akosna),
severnega in južnega. Če naprej
da je pospreči zadržati. Krovje on
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, omija in okusa
v vodi topljen polimerni snovki, katere
in Zveča hitj, ki sodelujejo pri nastajanju

Meronomi

- katere vzorce bi uporabili za iskanje meronimov - holonimov?

III. Korpusi v leksikografiji: terminološki sovarček

1. naredimo specializiran korpus s področja, ki nas zanima
2. specializiran korpus primerjamo z referenčnim korpusom, da dobimo ključne besede (wordsmith)
3. ključnim besedam pogledamo kolokacije, in izberemo večbesedne termine (sketch engine)
4. iztočnicam preko konkordanc najdemo podpomene (če so) in poiščemo dobre primere uporabe
5. vnesemo dobljeno v program za izdelovanje terminoloških baz
