

Kodiranje slovenskih besedil

Informacija in Kodi – Vaja 2

Asistent: as. dr. Klemen Grm

Kodiranje vira

- Vir informacije: besedilo (≈ 100 znakov)
- Komunikacijski kanal: računalniški pomnilnik oz. mrežna povezava (2 znaka)
- Kodiranje slovenskih besedil

Kodiranje vira

- Uporabni/neuporabni kodi
- Enakomerni/neenakomerni kodi
- Trenutni/netrenutni kodi

Kodiranje besedil

- Uporabni, trenutni kodi
- Enakomerni ali neenakomerni
- Prvi standard: tabela ASCII

ASCII

- 7-bitna kodna tabela – 128 znakov
- Primarni namen: angleška besedila
- Enakomerni kod
- Hex zapis za lažjo predstavo računalniškega pomnilnika

ASCII (1977/1986)

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
0_ 0	NUL 0000	SOH 0001	STX 0002	ETX 0003	EOT 0004	ENQ 0005	ACK 0006	BEL 0007	BS 0008	HT 0009	LF 000A	VT 000B	FF 000C	CR 000D	SO 000E	SI 000F
1_ 16	DLE 0010	DC1 0011	DC2 0012	DC3 0013	DC4 0014	NAK 0015	SYN 0016	ETB 0017	CAN 0018	EM 0019	SUB 001A	ESC 001B	FS 001C	GS 001D	RS 001E	US 001F
2_ 32	SP 0020	! 0021	" 0022	# 0023	\$ 0024	% 0025	& 0026	' 0027	(0028) 0029	* 002A	+ 002B	, 002C	- 002D	. 002E	/ 002F
3_ 48	0 0030	1 0031	2 0032	3 0033	4 0034	5 0035	6 0036	7 0037	8 0038	9 0039	: 003A	; 003B	< 003C	= 003D	> 003E	? 003F
4_ 64	@ 0040	A 0041	B 0042	C 0043	D 0044	E 0045	F 0046	G 0047	H 0048	I 0049	J 004A	K 004B	L 004C	M 004D	N 004E	O 004F
5_ 80	P 0050	Q 0051	R 0052	S 0053	T 0054	U 0055	V 0056	W 0057	X 0058	Y 0059	Z 005A	[005B	\ 005C] 005D	^ 005E	_ 005F
6_ 96	` 0060	a 0061	b 0062	c 0063	d 0064	e 0065	f 0066	g 0067	h 0068	i 0069	j 006A	k 006B	l 006C	m 006D	n 006E	o 006F
7_ 112	p 0070	q 0071	r 0072	s 0073	t 0074	u 0075	v 0076	w 0077	x 0078	y 0079	z 007A	{ 007B	 007C	} 007D	~ 007E	DEL 007F

Windows-1250 (CP1250)

- Razširitev ASCII za srednjeevropske jezike
- 8-bitni enakomerni kod, 256 znakov
- Poljščina, češčina, slovaščina, madžarščina, slovenščina, hrvaščina, srbščina (latinica), romunščina, albanščina

Windows-1250																
	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
0_0	NUL 0000	SOH 0001	STX 0002	ETX 0003	EOT 0004	ENQ 0005	ACK 0006	BEL 0007	BS 0008	HT 0009	LF 000A	VT 000B	FF 000C	CR 000D	SO 000E	SI 000F
1_16	DLE 0010	DC1 0011	DC2 0012	DC3 0013	DC4 0014	NAK 0015	SYN 0016	ETB 0017	CAN 0018	EM 0019	SUB 001A	ESC 001B	FS 001C	GS 001D	RS 001E	US 001F
2_32	SP 0020	! 0021	" 0022	# 0023	\$ 0024	% 0025	& 0026	' 0027	(0028) 0029	* 002A	+ 002B	, 002C	- 002D	. 002E	/ 002F
3_48	0 0030	1 0031	2 0032	3 0033	4 0034	5 0035	6 0036	7 0037	8 0038	9 0039	: 003A	; 003B	< 003C	= 003D	> 003E	? 003F
4_64	@ 0040	A 0041	B 0042	C 0043	D 0044	E 0045	F 0046	G 0047	H 0048	I 0049	J 004A	K 004B	L 004C	M 004D	N 004E	O 004F
5_80	P 0050	Q 0051	R 0052	S 0053	T 0054	U 0055	V 0056	W 0057	X 0058	Y 0059	Z 005A	[005B	\ 005C] 005D	^ 005E	_ 005F
6_96	` 0060	a 0061	b 0062	c 0063	d 0064	e 0065	f 0066	g 0067	h 0068	i 0069	j 006A	k 006B	l 006C	m 006D	n 006E	o 006F
7_112	p 0070	q 0071	r 0072	s 0073	t 0074	u 0075	v 0076	w 0077	x 0078	y 0079	z 007A	{ 007B	 007C	} 007D	~ 007E	DEL 007F
8_128	€ 20AC		,		„ 201E	… 2026	† 2020	‡ 2021		‰ 2030	Š 0160	< 2039	Ś 015A	Ť 0164	Ž 017D	Ž 0179
9_144		` 2018	/	“ 201C	” 201D	• 2022	— 2013	— 2014		™ 2122	š 0161	> 203A	ś 015B	ť 0165	ž 017E	ž 017A
A_160	NBSP 00A0	˘ 02C7	˘ 02D8	Ł 0141	ł 00A4	Ą 0104	! 00A6	Ś 00A7	“ 00A8	© 00A9	Ş 015E	« 00AB	¬ 00AC	SHY 00AD	® 00AE	Ž 017B
B_176	° 00B0	± 00B1	¸ 02DB	ł 0142	´ 00B4	µ 00B5	¶ 00B6	· 00B7	, 00B8	ą 0105	ş 015F	» 00BB	Ł 013D	˘ 02DD	I 013E	ž 017C
C_192	Ř 0154	Á 00C1	Â 00C2	Ă 0102	Ǻ 00C4	Í 0139	Ć 0106	Ç 00C7	Č 010C	É 00C9	Ę 0118	Ě 00CB	Ě 011A	Í 00CD	Î 00CE	Ď 010E
D_208	Đ 0110	Ñ 0143	Ň 0147	Ó 00D3	Ô 00D4	Õ 0150	Ö 00D6	× 00D7	Ř 0158	Ů 016E	Ú 00DA	Ů 0170	Ü 00DC	Ý 00DD	Ť 0162	ß 00DF
E_224	ř 0155	á 00E1	â 00E2	ă 0103	ǻ 00E4	í 013A	ć 0107	ç 00E7	č 010D	é 00E9	ę 0119	ě 00EB	ě 011B	í 00ED	î 00EE	ď 010F
F_240	đ 0111	ñ 0144	ň 0148	ó 00F3	ô 00F4	õ 0151	ö 00F6	÷ 00F7	ř 0159	ů 016F	ú 00FA	ü 0171	ü 00FC	ý 00FD	ț 0163	· 02D9

IBM-852 (CP852)

- Še ena razširjena ASCII tabela
- Spodnjih 7 bitov: ASCII
- Isti jeziki kot Windows-1250

ASCII (1977/1986)

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
0_0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
	0000	0001	0002	0003	0004	0005	0006	0007	0008	0009	000A	000B	000C	000D	000E	000F
1_16	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
	0010	0011	0012	0013	0014	0015	0016	0017	0018	0019	001A	001B	001C	001D	001E	001F
2_32	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
	0020	0021	0022	0023	0024	0025	0026	0027	0028	0029	002A	002B	002C	002D	002E	002F
3_48	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
	0030	0031	0032	0033	0034	0035	0036	0037	0038	0039	003A	003B	003C	003D	003E	003F
4_64	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
	0040	0041	0042	0043	0044	0045	0046	0047	0048	0049	004A	004B	004C	004D	004E	004F
5_80	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
	0050	0051	0052	0053	0054	0055	0056	0057	0058	0059	005A	005B	005C	005D	005E	005F
6_96	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
	0060	0061	0062	0063	0064	0065	0066	0067	0068	0069	006A	006B	006C	006D	006E	006F
7_112	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL
	0070	0071	0072	0073	0074	0075	0076	0077	0078	0079	007A	007B	007C	007D	007E	007F

Code page 852

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
8_128	Ç	ü	é	â	ã	û	ć	ç	ł	ë	ő	õ	î	ž	ä	ć
	00C7	00FC	00E9	00E2	00E4	016F	0107	00E7	0142	00EB	0150	0151	00EE	0179	00C4	0106
9_144	É	Í	Î	Ô	Õ	Ł	Į	Ś	ś	Ö	Ů	Ť	ţ	Ț	×	Č
	00C9	0139	013A	00F4	00F6	013D	013E	015A	015B	00D6	00DC	0164	0165	0141	00D7	010D
A_160	Á	Í	Ó	Ú	Ā	ā	Ž	ž	Ę	ę	¬	Ż	Č	Ş	«	»
	00E1	00ED	00F3	00FA	0104	0105	017D	017E	0118	0119	00AC	017A	010C	015F	00AB	00BB
B_176	⌘	⌘	⌘		†	Á	Â	Ě	Ş	†	‖	¶	¶	Ž	ž	ŀ
	2591	2592	2593	2502	2524	00C1	00C2	011A	015E	2563	2551	2557	255D	017B	017C	2510
C_192	Ł	ł	Ť	ť	—	†	Ǻ	ǻ	Ł	ŕ	Ł	ŕ	ŕ	=	†	▣
	2514	2534	252C	251C	2500	253C	0102	0103	255A	2554	2569	2566	2560	2550	256C	00A4
D_208	ď	Đ	Ď	Ě	ď	Ň	í	î	ě	Ĵ	ŕ	■	■	Ť	Ů	■
	0111	0110	010E	00CB	010F	0147	00CD	00CE	011B	2518	250C	2588	2584	0162	016E	2580
E_224	Ó	ß	Ô	Ň	ň	ň	Š	š	Ř	Ú	ř	Ů	ý	ý	ţ	´
	00D3	00DF	00D4	0143	0144	0148	0160	0161	0154	00DA	0155	0170	00FD	00DD	0163	00E4
F_240	SHY	~	ˆ	˘	˙	Š	÷	˘	˚	˚	˚	ű	Ř	ř	■	NBSP
	00AD	02DD	02DE	02C7	02D8	00A7	00F7	00B8	00B0	00A8	02D9	0171	0158	0159	25A0	00A0

ISO-8859-2

- Mednarodni standard za *prikaz* besedil
- Nova 8-bitna tabela
- Nima kontrolnih znakov
- Isti jeziki + lužiščina, turkmenščina

ISO/IEC 8859-2 (Latin-2)

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
0_ 0																
1_ 16																
2_ 32	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
	0020	0021	0022	0023	0024	0025	0026	0027	0028	0029	002A	002B	002C	002D	002E	002F
3_ 48	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
	0030	0031	0032	0033	0034	0035	0036	0037	0038	0039	003A	003B	003C	003D	003E	003F
4_ 64	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
	0040	0041	0042	0043	0044	0045	0046	0047	0048	0049	004A	004B	004C	004D	004E	004F
5_ 80	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
	0050	0051	0052	0053	0054	0055	0056	0057	0058	0059	005A	005B	005C	005D	005E	005F
6_ 96	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
	0060	0061	0062	0063	0064	0065	0066	0067	0068	0069	006A	006B	006C	006D	006E	006F
7_ 112	p	q	r	s	t	u	v	w	x	y	z	{		}	~	
	0070	0071	0072	0073	0074	0075	0076	0077	0078	0079	007A	007B	007C	007D	007E	
8_ 128																
9_ 144																
A_ 160	NBSP	À	Á	Â	Ã	Ä	Å	Š	Ś	Š	Œ	Ž	SHY	Ž	Ž	
	00A0	0104	02D8	0141	00A4	013D	015A	00A7	00A8	0160	015E	0164	0179	00AD	017D	017E
B_ 176	°	à	á	â	ã	ä	å	ı	ş	ş	ţ	ž	ˆ	ž	ž	
	00B0	0105	02DB	0142	00B4	013E	015B	02C7	00B8	0161	015F	0165	017A	02DD	017E	017C
C_ 192	Ř	Á	Â	Ã	Ä	Í	Ć	Ç	Č	É	Ê	Ë	Ě	Í	Î	Ď
	0154	00C1	00C2	0102	00C4	0139	0106	00C7	010C	00C9	0118	00CB	011A	00CD	00CE	010E
D_ 208	Ð	Ñ	Ň	Ó	Ô	Õ	Ö	×	Ř	Ů	Ú	Ů	Ů	Ý	Ť	ß
	0110	0143	0147	00D3	00D4	0150	00D6	00D7	0158	016E	00DA	0170	00DC	00DD	0162	00DF
E_ 224	ř	á	â	ã	ä	í	ć	ç	č	é	ê	ë	ě	í	î	ď
	0155	00E1	00E2	0103	00E4	013A	0107	00E7	010D	00E9	0119	00EB	011B	00ED	00EE	010F
F_ 240	đ	ń	ň	ó	ô	õ	ö	÷	ř	ů	ú	ů	ů	ý	ť	·
	0111	0144	0148	00F3	00F4	0151	00F6	00F7	0159	016F	00FA	0171	00FC	00FD	0163	02D9

MacCE (Mac OS Central European)

- Svojevrstna razširitev ASCII tabele
- -Albanščina, hrvaščina, romunščina
- +Estonščina, litovščina, latvijščina
- Razširjena ločila, matematični simboli

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
8_ 128	Ä 00C4	Å 0100	ā 0101	É 00C9	Ą 0104	Ö 00D6	Ü 00DC	á 00E1	ą 0105	Č 010C	ă 00E4	č 010D	Ć 0106	ć 0107	é 00E9	Ž 0179
9_ 144	ž 017A	Đ 010E	ı 00ED	đ 010F	Ě 0112	ē 0113	Ê 0116	ó 00F3	è 0117	ô 00F4	ö 00F6	õ 00F5	ú 00FA	Ě 011A	ě 011B	ü 00FC
A_ 160	† 2020	° 00B0	₺ 0118	£ 00A3	§ 00A7	• 2022	¶ 00B6	ß 00DF	® 00AE	© 00A9	™ 2122	₹ 0119	“ 00AB	≠ 2260	ğ 0123	İ 012E
B_ 176	ı 012F	Ī 012A	≤ 2264	≥ 2265	ī 012B	ķ 0136	∂ 2202	Σ 2211	ł 0142	Ł 013B	Į 013C	Ľ 013D	ŕ 013E	Ĺ 0139	Í 013A	Ŋ 0145
C_ 192	ņ 0146	Ň 0143	¬ 00AC	√ 221A	ñ 0144	Ñ 0147	Δ 2206	« 00AB	» 00BB	… 2026	NBSP 00A0	ň 0148	Ŏ 0150	Õ 00D5	ö 0151	Ō 014C
D_ 208	— 2013	— 2014	“ 201C	” 201D	‘ 2018	’ 2019	÷ 00F7	◊ 25CA	ō 014D	Ř 0154	í 0155	Ř 0158	‹ 2039	› 203A	ř 0159	Ŕ 0156
E_ 224	ŕ 0157	Š 0160	, 201A	„ 201E	š 0161	Ś 015A	ś 015B	Á 00C1	Ť 0164	ť 0165	Í 00CD	Ž 017D	ž 017E	Ů 016A	Ó 00D3	Ô 00D4
F_ 240	ū 016B	Ů 016E	Ú 00DA	ů 016F	Ů 0170	ů 0171	Ů 0172	ų 0173	Ý 00DD	ý 00FD	ķ 0137	Ž 017B	Ł 0141	ž 017C	Ÿ 0122	˘ 02C7

Unicode

Situacija 1995:

- Ločene kodne tabele za različne jezike (skupine jezikov)
- Implementacije, specifične operacijskim sistemom in proizvajalcem
- Pomankljiva standardizacija

Unicode

Ideja:

- En kod za vsa besedila
- Enoten standard
- Neodvisnost od strojne, programske opreme

Unicode

- 154 abeced, do 10000 znakov / abecedo
- Potrebe:
 - Veliko število kodnih mest ($N > 10^6$)
 - Neenakomerno kodiranje za varčevanje s pomnilnikom
 - Trenutni kod
- Ena Unicode tabela, 1.112.064 znakov
- Več različnih transformacijskih formatov (UTF)

UTF-16

- Neenakomeren kod, 1 ali 2 16-bitna znaka
- LE in BE izvedbe
- Uporaba: (interno) Windows, Java
- Odsvetovana uporaba na spletu
- Izrazito negospodarno kodiranje – vsaj 16 bitov/znak

UTF-8

- De-facto standard za spletno kodiranje besedila
- Uporaba:
 - Unix-like sistemi (Linux, Mac OS, Android,...)
 - Splet
 - Python
- Neenakomeren kod: 1, 2, 3 ali 4 8-bitni znaki
- ASCII znaki: identičen zapis
- Struktura predpon

Number of bytes	Bits for code point	First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
1	7	U+0000	U+007F	0xxxxxxx			
2	11	U+0080	U+07FF	110xxxxx	10xxxxxx		
3	16	U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	
4	21	U+10000	U+10FFFF ^[12]	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

Laboratorijska vaja 2

- Naloga 1 – izpis kodnih zamenjav šumnikov slovenske abecede
 - Č,Š,Ž,č,š,ž
- IBM-852,
- ISO-8859-2
- Windows-1250,
- MacCE
- UTF-16LE, UTF-16BE
- UTF-8
- Decimalni, hex, binarni zapis

Laboratorijska vaja 2

- Naloga 2 – kodiranje in dekodiranje UTF-8
 - Vhod: zaporedje kodnih mest
 - Izhod: v utf-8 kodirana besedilna datoteka
1. Beri vhodno datoteko, pretvori kodna mesta v številske podatke
 2. Pretvorba kodnih mest v Unicode znake
 3. Zapis izhodne datoteke
 4. Zapis unikatnih znakov in kodnih zamenjav