## 1.3 Kody alfanumeryczne

## 1. Uwagi wstępne

Jeśli zbiorem obiektów kodowanych  $V_1$  jest zbiór znaków to kod  $f:V_1 \to V_2^*$  (dla ustalonego alfabetu  $V_2$  na ogół binarnego) nazywamy kodem alfanumerycznym (ang. alphanumeric code).

W tym podrozdziale poznamy szereg binarnych kodów alfanumerycznych. Najpopularniejsze binarne kody alfanumeryczne to ośmiobitowy kod ASCII (American Standard Code for Information Interchange) i szesnastobitowy kod Unicode. Są to oczywiście kody o stałej długości słowa kodowego.

Znany każdemu telegrafiście i harcerzowi kod Morse'a jest również kodem alfanumerycznym binarnym ale kodem zmiennej długości. Kody alfanumeryczne ASCII zmiennej długości mają specyficzne zalety opisane bliżej w podrozdziale o kompresji danych.

Kod ASCII został wprowadzony w USA w 1963 roku jako kod 7 bitowy. Kod ASCII w swej podstawowej wersji jest więc w zasadzie 7-mio bitowym kodem ale uzupełnionym z reguły bitem parzystości. Jako kod 8 bitowy ma jednak szereg odmian tzw. wersji narodowych.

Kod ASCII jest równoważny z kodem ISO-7 (ISO to skrót od International Organization for Standardization). Kod ISO-7 został opisany w normie ISO 646 wydanej w roku 1973. Norma ISO 646 przewiduje wprowadzenie znaków narodowych co wiąże się z rozszerzeniem kodu ASCII do kodu 8 bitowego

ISO-8859-2 to ogólnie już przyjęty standard kodowania polskich liter wg polskich norm PN stanowiący rozszerzenie ASCII.

Kod alfanumeryczny UNICODE jest standardem ISO/IEC-10646. Jest to kod 16 bitowy (pierwsza wersja) lub 31 bitowy (druga wersja). UNICODE umożliwia zapisanie wszystkich znaków z alfabetów narodowych (również cyrylicy, alfabetu chińskiego i alfabetu japońskiego)

Kod ASCII w wersji podstawowej koduje 128 znaków. Pierwsze 33 znaki (uporządkowane wg wartości w kodzie NKB słowa kodowego to tzw. znaki sterujące takie jak CR (Carriage Return czyli "powrót karetki") lub LF (Line Feed czyli "od nowego wiersza"). Służą one do sterowania systemem drukowania lub wyświetlania znaków. Pozostałe znaki to m.in. małe i duże litery alfabetu angielskiego, cyfry, znaki przestankowe oraz kilka symboli matematycznych takich jak nawiasy i znak równości.

## 2. Kod ASCII

W technice cyfrowej i wywodzącej się z niej inżynierii komputerowej dominują obecnie kody alfanumeryczne oparte na standardzie ASCII (American Standard Code for Information Interchange), znanym także jako ANSI X3.4 oraz ISO-646-US i US-ASCII. W standardzie tym 128 podstawowym znakom, takim jak litery alfabetu łacińskiego, cyfry, znaki przestankowe oraz kody sterujące, przyporządkowano liczby z zakresu od 0 do 127. Z uwagi na oczywiste problemy z sortowaniem tak różnorodnych znaków tabele kodowe zestawia się według wartości kodów (dziesiętnych i szesnastkowych) - jak w poniższej tabeli.

Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char
0	00	NUL null	32	20	space	64	40	(a)	96	60	`
1	01	SOH start of heading	33	21	!	65	41	A	97	61	a
2	02	STX start of text	34	22	"	66	42	В	98	62	b
3	03	ETX end of text	35	23	#	67	43	C	99	63	c
4	04	EOT end of transmission	36	24	\$	68	44	D	100	64	d
5	05	ENQ enquiry	37	25	%	69	45	Е	101	65	e
6	06	ACK acknowledge	38	26	&	70	46	F	102	66	f
7	07	BEL bell	39	27	•	71	47	G	103	67	g
8	08	BS backspace	40	28	(	72	48	Н	104	68	h
9	09	TAB horizontal tabulation	41	29	)	73	49	I	105	69	i
10	0A	LF line feed	42	2A	*	74	4A	J	106	6A	j
11	0B	VT vertical tabulation	43	2B	+	75	4B	K	107	6B	k
12	0C	FF form feed	44	2C	,	76	4C	L	108	6C	1
13	0D	CR carriage return	45	2D	-	77	4D	M	109	6D	m
14	0E	SO shift out	46	2E		78	<b>4</b> E	N	110	6E	n
15	0F	SI shift in	47	2F	/	79	4F	O	111	6F	o
16	10	DLE data link escape	48	30	0	80	50	P	112	70	p
17	11	DC1 device control one	49	31	1	81	51	Q	113	71	q
18	12	DC2 device control two	50	32	2	82	52	R	114	72	r
19	13	DC3 device control three	51	33	3	83	53	S	115	73	S
20	14	DC4 device control four	52	34	4	84	54	T	116	74	t
21	15	NAK negative acknowledge	53	35	5	85	55	U	117	75	u
22	16	SYN synchronous idle	54	36	6	86	56	V	118	76	v
23	17	ETB end of transmission block	55	37	7	87	57	W	119	77	W
24	18	CAN cancel	56	38	8	88	58	X	120	78	X
25	19	EM end of medium	57	39	9	89	59	Y	121	79	y
26	1A	SUB substitute	58	3A	:	90	5A	Z	122	7A	Z
27	1B	ESC escape	59	3B	,	91	5B	[	123	7B	{
28	1C	FS file separator	60	3C	<	92	5C	\	124	7C	
29	1D	GS group separator	61	3D	=	93	5D	]	125	7D	}
30	1E	RS record separator	62	3E	>	94	5E	^	126	7E	~
31	1F	US unit separator	63	3F	?	95	5F		127	7F	DEL

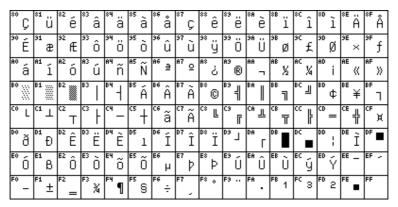
Tab. 1. US-ASCII

ASCII jest kodem 7-bitowym, przy czym oryginalny standard nie definiuje roli ósmego bitu, który w czysto sprzętowych realizacjach cyfrowych układów sterowania ciągle bywa wykorzystywany do kontroli parzystości lub do przechowywania dodatkowego atrybutu (np. podświetlenia). Najczęściej jednak ósmy bit służy rozszerzeniu podstawowego kodu ASCII o niezbędne znaki alfabetów narodowych, symbole matematyczne, znaki semigraficzne itp.. Niegdyś próbowano dopasować kod ASCII do warunków lokalnych poprzez odmienne wykorzystanie niektórych wartości kodów; np. w ISO-646-DE pod wartościami z zakresu 123-126 (7B-7E) kryją się odpowiednio znaki ä, ö, ü i ß. W miarę, jak rosły potrzeby, pojawiały się kolejne kody o rozmaitych sposobach i zakresach implementacji dodatkowych znaków, a wielu producentów sprzętu i oprogramowania wprowadzało własne "standardy". Z tych najbardziej znane jest rozszerzenie kodu ASCII wprowadzone przez firmę IBM, czyli CP437 (IBM PC Extended ASCII czyli DosLatinUS),

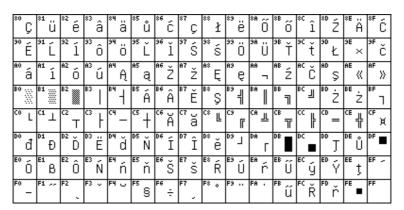
które zainicjowało całą serię tzw. stron kodowych, począwszy od CP850 (DosLatin1) i CP852 (DosLatin2), obecnych w kartach graficznych, drukarkach i w MS-DOS.

		01	0	02	0	03	٠	04	*	05	÷	06	÷	07	•	08	٠	09	0	0A	0	0B	ď	oc.	Ŷ	OD	ľ	0E	Ŋ	0F	×
10	٠	11	4	12	\$	13	ij	14	1	15	8	16	-	17	\$	18	<b>†</b>	19	¥	1A	÷	18	+	10	_	1D	↔	1E	•	1F	•
20		21	ļ	22	11	23	#	24	\$	25	%	26	8	27	1	28	(	29	)	2A	*	2B	+	20	,	20	-	2E		2F	7
30	0	31	1	32	2	33	3	34	4	35	5	36	6	37	7	38	8	39	9	3 <b>A</b>	:	3B	;	30	<	3D	=	3E	>	3F	?
40	0	41	Α	42	В	43	С	44	D	45	Ε	46	F	47	G	48	Н	49	Ι	4A	J	48	Κ	40	L	40	М	ЧE	N	ЧF	0
50	Р	51	Q	52	R	53	S	54	Т	55	U	56	٧	57	М	58	Χ	59	Υ	5A	Z	5B	[	5C	\	50	]	5E	^	5F	_
60	,	61	а	62	b	63	С	64	d	65	е	66	f	67	g	68	h	69	i	6A	j	6B	k	ec.	1	6D	m	6E	n	6F	0
70	р	71	q	72	r	73	S	74	t	75	u	76	٧	77	W	78	×	79	y	7A	Z	78	}	7C		70	3	7E	2		
80	Ç	81	ü	82	é	83	â	84	ä	85	à	86	å	87	Ç	**	ê	89	ë	8A	è	88	ï	8C	î	\$D	ì	\$E	Ä	8F	Å
90	É	91	æ	92	Æ	93	ô	94	ö	95	ò	96	û	97	ũ	98	ÿ	99	ö	9A	Ü	9B	¢	9C	£	9D	¥	9E	Pts	9F	f
ΑO	á	A1	í	A2	ó	A3	ú	A4	ñ	A5	Ñ	A6	a	A7	2	A8	ò	A9	_	AA	7	AB	X	AC	14	AD	i	AE	((	AF	>>
BO	<b>W</b>	B1		B2		В3	Π	В4	+	B5	1	B6	$\parallel$	B7	П	B8	٦	В9	1	BA		BB	٦	BC	Ţ	BD	Ш	BE	1	BF	٦
CO	L	C1	Τ	CZ	Т	C3	F	CЧ	_	C5	+	C6	F	C7	$\ \cdot\ $	C8	L	C9	ſſ	CA	T	СВ	╦	СС	ŀ	CD	=	CE	#	CF	⊥
DO	Ш	D1	₹	DZ	Т	D3	L	D4	F	D5	F	D6	Г	D7	⊥	D8	†	D9	J	DA	Γ	DB		DC		DD	I	DE	I	DF	•
ΕO	α	E1	β	E2	Γ	E3	Π	EЧ	Σ	E5	σ	E6	μ	E7	τ	E®	Φ	E9	Θ	ΕA	Ω	EB	δ	EC	00	ED	Ø	EE	€	EF	n
FO	≡	F1	±	F2	≥	F3	≤	FЧ	ſ	F5	J	F6	÷	F7	×	F8	۰	F9		FA	•	FB	√	FC	n	FD	2	FE	•	FF	

Tab. 2. CP437



Tab. 3. CP850



Tab. 4. CP852

ISO 8859 to zestaw kilkunastu tabel wykorzystujących wszystkie 256 wartości 8-bitowego słowa kodowego i rozszerzających właściwą tabelę ASCII (kody od 0 do 127 pozostają bez zmian) o znaki pochodzące z poszczególnych regionów Europy, alfabet grecki, cyrylicę

itp. Podobnie jak w podstawowym kodzie ASCII, pierwsze 32 pozycje każdej rozszerzonej tabeli ISO 8859-*n* to niedrukowalne znaki sterujące. Odpowiednikiem CP850 jest ISO 8859-1, czyli Latin1, dedykowany Europie Zachodniej. Odpowiednikiem CP852 jest ISO 8859-2 (Latin2) zawierający komplet polskich znaków, a np. greckie litery można znaleźć w ISO 8859-7 (Greek)

Dec I	Hex	Char		**	C.						
128	80	PAD padding character	Dec		Char	Dec	Hex	Char	Dec	Hex	Char
	81	HOP high octet preset	160			192	C0	Ŕ	224	E0	ŕ
	82	BPH break permitted here	161		Ą	193	C1	Á	225	E1	á
	83	NBH no break here	162		Č	194	C2	Â	226	E2	â
	84	IND index	163	A3	Ł	195	C3	Ă	227	E3	ă
	85	NEL next line	164		¤	196	C4	Ä	228	E4	ä
	86	SSA start of selected area	165	A5	Ľ	197	C5	Ĺ	229	E5	ĺ
	87	ESA end of selected area	166		Ś	198	C6	Ć	230	E6	ć
	88	HTS character tabulation set	167	A7	§	199	C7	Ç	231	E7	ç
		HTJ character tabulation with	168	A8		200	C8	Č	232	E8	č
137	89	justification	169	A9	Š	201	C9	É	233	E9	é
138 8	8A	VTS line tabulation set	170	AA	Ş	202	CA	Ę	234	EA	ę
	8B	PLD partial line forward	171	AB	Ť	203	СВ	Ë	235	EB	ë
140	8C	PLU partial line backward	172	AC	Ź	204	CC	Ě	236	EC	ě
141 8	8D	RI reverse line feed	173	AD	soft	205	CD	Í	237	ED	í
142	8E	SS2 single-shift two	174	۸E	hyphen Ž	206	CE	Î	238	EE	î
143	8F	SS3 single-shift three	175		Ż	207	CF	Ď	239	EF	ď
144	90	DCS device control string		B0	<b>2</b> .	208	D0	Đ	240	F0	đ
145	91	PU1 private use one	177	В0 В1	a	209	D1	Ń	241	F1	ń
146	92	PU2 private use two		B2	ą	210	D2	Ň	242	F2	ň
147	93	STS set transmit state		B3	ł	211	D3	Ó	243	F3	ó
148	94	CCH cancel character	180		1	212	D4	Ô	244	F4	ô
149	95	MW message waiting		B5	ľ	213	D5	Ő	245	F5	ő
150	96	SPA start of guarded area		B6	ś	214	D6	Ö	246	F6	ö
151	97	EPA end of guarded area		В7	\$ •	215	D7	×	247	F7	÷
152	98	SOS start of string	184			216	D8	Ř	248	F8	ř
153	99	SGCI single graphic character introducer	185		š	217		Ů	249		ů
154	9A	SCI single character introducer	186		ş	218	DA	Ú	250	FA	ú
155		CSI control sequence introducer	187		ť	219	DB	Ű	251	FB	ű
156		ST string terminator	188		ź	220	DC	Ü	252	FC	ü
157		OSC operating system command	189		"	221	DD	Ý	253	FD	ý
158		PM privacy message	190		ž	222	DE	Ţ	254	FE	ţ
159		APC application program command	191		Ż	223	DF	В	255	FF	•

Tab. 5. ISO 8859-2 (Latin2)

Wprowadzenie standardu ISO 8859 miało w założeniu uporządkować zasady kodowania rozszerzeń ASCII. Niestety, nie tylko nie zapobiegło definitywnie dalszemu mnożeniu się stosowanych kodów alfanumerycznych, lecz wręcz zapoczątkowało kolejną ich falę,

tak w ramach samego standardu (wystarczy prześledzić ewolucję obsługi języków krajów nadbałtyckich od ISO 8859-4 (Latin4) poprzez ISO 8859-10 (Latin6) do ISO 8859-13 (Latin7 Baltic Rim)), jak i incjatyw zewnętrznych (na przykład Microsoft zmodyfikował tabelę ISO 8859-2 i wprowadził ją do Windows jako stronę kodową CP1250 (WinLatin2)). Główną tego przyczyną są obiektywne ograniczenia wynikające ze zbyt małej ilości słów kodu 8-bitowego.

80	€			82	,			84	,,	85		86	#	87	#			89	%	8A	š	88	<	8C	Ś	8D	Ť	\$E	ž	8F _	2
		91	٠	92	,	93	"	94	"	95	•	96	-	97	-			99	тн	9A	š	98	>	90	ś	9D	ť	9E	ž	9F _	2
AO		A1	Ÿ	A2	Ü	A3	Ł	A4	Ħ	A5	Ą	A6	1	A7	8	A8		A9	0	AA	Ş	AB	((			AD	-	AE	®	AF .	2
BO	۰	B1	±	B2	Ţ	В3	ł	вч	-	B5	μ	B6	1	B7		B8	,	В9	ą	BA	Ş	BB	<b>&gt;&gt;</b>	ВС	Ľ	BD	~	BE	ĭ	BF .	2
CO	Ŕ	C1	Á	CZ	Â	C3	Ă	СЧ	Ä	C5	Ĺ	C6	ć	C7	Ç	C8	č	C9	É	CA	Ę	СВ	Ë	СС	Ě	CD	Í	CE	Î	CF C	
DO	Đ	D1	Ń	DZ	Ň	D3	Ó	D4	ô	D5	ő	D6	ö	D7	×	D8	Ř	D9	Ů	DA	Ú	DB	Ű	DC	Ü	DD	Ý	DE	Ţ	DF E	3
ΕO	ŕ	E1	á	E2	â	E3	ă	EЧ	ä	E5	ĺ	E6	ć	E7	Ç	E\$	č	E9	é	ΕA	ę	EB	ë	EC	ě	ED	í	EE	î	EF (	
F0	đ	F1	ń	F2	ň	F3	ó	F4	ô	F5	ő	F6	ö	F7	÷	F8	ř	F9	ů	FA	ú	FB	ű	FC	ü	FD	ý	FE	ţ	FF .	

Tab. 6. CP1250 (WinLatin2)

## 3. Unicode

W zapisie UTF-8 każdy znak UCS/Unicode jest przedstawiany w postaci sekwencji od jednego do sześciu 8-bitowych bajtów, zależnie od wartości samego znaku. Poniższa tabela obrazuje zasadę, na jakiej to się odbywa:

Liczba bajtów	1. bajt	2. bajt	3. bajt	4. bajt	Liczba bitów	_	ksymalna vartość
UTF-8					Unicode	Hex	Dec
1	0xxxxxxx				7	7F	127
2	110xxxxx	10xxxxxx			11	7FF	2047
3	1110xxxx	10xxxxxx	10xxxxxx		16	FFFF	65535
4	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx	21	1FFFFF	2097151
5	111110xx	10xxxxxx	itd.		26	3FFFFFF	67108863
6	11111110x	10xxxxxx	itd.		31	7FFFFFFF	2147483647

**Tab. 7. UTF-8** 

```
0xxxxxxx : pierwszy i jedyny bajt sekwencji
10xxxxxx : drugi i dalsze bajty sekwencji
110xxxxx : pierwszy bajt sekwencji 2-bajtowej
1110xxxx : pierwszy bajt sekwencji 3-bajtowej
11110xxx : pierwszy bajt sekwencji 4-bajtowej
111110xx : pierwszy bajt sekwencji 5-bajtowej
1111110x : pierwszy bajt sekwencji 6-bajtowej
```

Tab. 8. Znaczenie bajtu w zapisie UTF-8

W wielobajtowej sekwencji bity oznaczone 'x'-ami czytane od 1-szego, najstarszego bajtu tworzą właściwą wartość znaku UCS/Unicode. Z kolei wartość bieżącego bajtu wskazuje

na jego miejsce w sekwencji UTF-8. W ten sposób zapis UTF-8 jest kompatybilny z 7-bitowym US-ASCII i zachowuje względną zwartość tekstów o niewielu znakach rozszerzonych, pozwalając jednocześnie na zapis nawet 31-bitowych wartości i na łatwą do realizacji synchronizację i interpretację przetwarzanych sekwencji.

**Przykład:** Oto kilka przykładów 16-bitowych wartości Unicode w zapisie UTF-8:

Unicode	UTF-8
Hex	Hex
0001	01
007F	7F
0080	C2 80
07FF	DF BF
0800	E0 A0 80
OFFF	EO BF BF
1000	E1 80 80
FFFF	EF BF BF

Zapis UTF-8 jest nadmiarowy, bowiem np. wartościom 16-bitowym przyporządkowuje wartości 24-bitowe (3-bajtowe). W efekcie istnieje wiele ciągów bajtów, które nie są legalnymi sekwencjami UTF-8, nawet jeżeli wykonanie przekształcenia odwrotnego jest technicznie możliwe.

**Przykład:** Ciąg bajtów C2 00 nie jest sekwencją UTF-8, ponieważ jego drugi bajt nie jest zgodny z wzorcem 10xxxxxx.

**Przykład:** Ciąg bajtów C0 80 nie jest legalną sekwencją UTF-8, ponieważ jest dłuższy, niż trzeba: po zastosowaniu do niego maski binarnej 110xxxxx 10xxxxxx uzyskujemy wynik 00000 000000, czyli zerową wartość kodu, którą w UTF-8 zapisuje się poprawnie w jednym bajcie.

sekwencja	C0 80	11000000	10000000	
maska		110xxxxx	10xxxxxx	
wynik	00	00000	000000	

Oprócz zapisu UTF-8 spotyka się czasem UTF-7 wykorzystujący tylko 7 bitów używanego słowa oraz UTF-16 będący po prostu zapisem kolejnych kodów.

A oto reprezentacja polskich "ogonków" w Unicode:

Znak	Kod	Znak	Kod
ą	0105	Ą	0104
ć	0107	Ć	0106
ę	0119	Ę	0118
ł	0142	Ł	0141
ń	0144	Ń	0143
Ó	00F3	Ó	00D3
Ś	015B	Ś	015A
Ź	017A	Ź	0179
Ż	017C	Ż	017B

Tab. 9. Polskie znaki w Unicode