



Language Manual

Czech

Sabrina and Eliska

Language Manual
Czech
Sabrine and Eliska
20 June 2007

This document was produced by Acapela Group. We welcome and consider all comments and suggestions. Please send them to:

Acapela Group
Box 1328
SE-171 26 Solna
Sweden

Phone +46 (0) 8 799 86 00
Fax + 46 (0) 8 799 86 01

Acapela Group
33, Boulevard Dolez
7000 Mons
Belgium

Tel: +32 (0)65 37 42 75
Fax: +32 (0)65 37 42 76

Acapela Group
3939, la Lauragaise
BP 58309
F-31683 Labège cedex
France

Tel: +33 (0)5 62 24 71 00
Fax: +33 (0)5 62 24 71 01

www.acapela-group.com

© Copyright Acapela Group 2007. All rights reserved.

List of contents

1	General	4
2	Letters in orthographic text	5
3	Punctuation characters	6
3.1	Comma, colon and semicolon	6
3.2	Quotation marks	6
3.3	Full stop	6
3.4	Question mark	6
3.5	Exclamation mark	6
3.6	Parentheses	6
4	Other non-alphanumeric characters	7
4.1	Non-punctuation characters	7
4.2	The ² and ³ signs	7
4.3	Symbols whose pronunciation varies depending on the context	8
4.3.1	Hyphen	8
4.3.2	Asterisk	8
5	Number processing	9
5.1	Full number pronunciation	9
5.2	Leading zero	10
5.3	Decimal numbers	10
5.4	Monetary amounts	10
5.5	Ordinal numbers	10
5.6	Arithmetic operators	11
5.7	Mixed digits and letters	11
5.8	Time of day	11
5.9	Dates	11
5.10	Phone numbers	12
5.10.1	Ordinary phone numbers	12
5.10.2	International phone numbers	13
6	How to change pronunciation errors	14
7	Czech Phonetic Text	15
7.1	Consonants	15
7.1.1	Symbols for the Czech consonants	15
7.2	Vowels	16
7.2.1	Symbols for the Czech vowels	16
7.3	Lexical accent	16
7.4	Pause	16
7.5	Glottal stops	16
8	Abbreviations	17
9	Web-addresses and email	19

1 General

This document discusses certain aspects of text-to-speech processing for the Czech text-to-speech system, in particular the different types of input characters and text that are allowed.

This version of the document corresponds to the High Density voice Sabine and the High Quality voice Eliska.

2 Letters in orthographic text

Characters from A-Z, a-z (as well as á, Á, ä, Ä, å, Å, ð, Þ, é, É, ê, Ê, í, Í, ï, Ï, ó, Ó, ö, Ö, ø, Ø, š, Š, Ÿ, Ź, ú, Ú, ü, Ü, ý, Ý, ž, Ž) may constitute a word. Certain other characters are also considered as letters, notably those used as letters in other European languages, i.e. “ł, ò, ç”. These letters are not pronounced as in their native languages though, they are pronounced as regular “l, o, a, c” when occurring in a word.

Characters outside of these ranges, i.e. numbers, punctuation characters and other non-alphanumeric characters are not considered as letters.

3 Punctuation characters

Punctuation marks appearing in a text affect both rhythm and intonation of a sentence. The following punctuation characters are permitted in the normal input text string:

, : ; “ ” . ? ! () '

3.1 Comma, colon and semicolon

Comma < , >, colon < : > and semicolon < ; > cause a brief pause to occur in a sentence, accompanied by a small rising intonation pattern just prior to the character.

3.2 Quotation marks

Quotes < “ ” > appearing around a single word or a group of words cause a brief pause before and after the quoted text.

3.3 Full stop

A full stop < . > is a sentence terminal punctuation mark which causes a falling end-of-sentence intonation pattern and is accompanied by a somewhat longer pause. A full stop may also be used as a decimal marker in a number (see chapter 5) and in abbreviations (see chapter 8).

3.4 Question mark

A question mark < ? > ends a sentence and causes question-intonation, first rising and then falling.

3.5 Exclamation mark

The exclamation mark < ! > behaves in a similar manner to the full stop, causing a falling intonation pattern followed by a pause.

3.6 Parentheses

Parentheses < () > appearing around a single word or a group of words cause a brief pause before and after the bracketed text.

4 Other non-alphanumeric characters

4.1 Non-punctuation characters

The characters listed below are processed as non-letter, non-punctuation characters. Some are pronounced at all times and others are only pronounced in certain contexts, which are described in the following sections of this chapter.

Symbol	Reading
/	Lomeno
+	Plus
\$	Dolar
€	Euro
<	menší než
>	větší než
%	procento
^	vokáň
	vertikální závorka
~	Tilda
@	Zavináč
=	Rovná se
²	see below
³	see below
-	see below
*	see below

Table 1 Non-punctuation characters

4.2 The ² and ³ signs

The reading of expressions with ² and ³ is:

Expression	Reading
mm ²	čtvereční milimetr
cm ²	čtvereční centimetr
m ²	čtvereční metr
km ²	čtvereční kilometr
mm ³	kubický milimetr
cm ³	kubický centimetr
m ³	kubický metr
km ³	kubický kilometr

4.3 Symbols whose pronunciation varies depending on the context

4.3.1 Hyphen

A hyphen < - > is pronounced “mínus” in two cases:

- if followed by a digit and no other digit is found in front of the hyphen
- if followed by a digit and an equals sign. If there is no equals sign, it is pronounced “pomlčka”.

In certain date formats, in between days or years, the hyphen is pronounced “až”.
In compounds or between words, the hyphen is not pronounced. Example: T-Mobile.

Expression	Reading
-3	mínus tři
44-3	čtyřicet čtyři pomlčka tři
44-3=41	čtyřicet čtyři mínus tři rovná se čtyřicet jeden
15-20 října	patnáctého až dvacátého října
6-10 listopadu	šestého až desátého listopadu
1998-2004	tisíc devjetset devadesát osum až dva tisíce čtyři
02-02-2002	Druhého února dva tisíce dvě

4.3.2 Asterisk

Asterisk < * > is only pronounced as “krát” if enclosed by digits and followed by equals sign. In other cases it is pronounced as “hvězdička”.

Expression	Reading
2*3	dva hvězdička tři
2*3=6	dva krát tři rovná se šest
*bc	hvězdička b c

5 Number processing

Strings of digits that are sent to the text-to-speech converter are processed in several different ways, depending on the format of the string of digits and the immediately surrounding punctuation or non-numeric characters. To familiarise the user with the various types of formatted and non-formatted strings of digits that are recognised by the system, a brief description of the basic number processing is provided below, along with examples.

Number processing is subdivided into the following categories:

Full number pronunciation

Leading zero

Decimal numbers

Currency amounts

Ordinal numbers

Arithmetic operators

Mixed digits and letters

Time of day

Year

Dates

Phone numbers

5.1 Full number pronunciation

Full number pronunciation is given for the whole number part of the digit string.

Example

2425	full number
2.425	full number
2 425	full number
24,25	24 is a full number, 25 is the decimal part

Numbers denoting thousands, millions and billions (numbers larger than 999) may be grouped using space or full stop. In order to achieve the right pronunciation the grouping must be done correctly.

The rules for grouping of numbers are the following:

- Numbers are grouped in groups of three starting at the end.
- The first group in a number may consist of one, two, or three digits.
- If a group, other than the first, does not contain exactly three digits, the sequence of digits is not interpreted as a full number.
- The highest number read is 999999999999 (twelve digits). Numbers higher than this are read as separate digits.

Number	Reading
2580	dva tisíce p t set osumdesát
2 580	“
2.580	“
25800	dvacet p t tisíc osum set
25 800	“
25.800	“
2580350	dva milióny p t set osumdesát tisíc tři sta padesát
2 580 350	“
2.580.350	“
1000000000	jedna miliarda
1234567890123	jedna dva tři čtyři pět šest sedm osum devět nula jedna dva tři

23 456 789 012

dvacet tři miliarda čtyři sta padesát šest milióny sedum set osumdesát devět tisíc dvanáct

5.2 Leading zero

Numbers that begin with 0 (zero) are read as a whole number, with a zero preceding it.

Number	Reading
09253	nula devět tisíc dv st padesát tři
020	nula dvacet

5.3 Decimal numbers

Comma or full stop may be used when writing decimal numbers.

The full number part of the decimal number (the part before comma or full stop) is read according to the rules in 5.1. The decimals (the part after comma or full stop) are read as separate digits if there are more than 3 digits after the comma. When the decimals are read as a whole number, the words “desetin”, “setin”, and “tisícin” are added after the decimal, depending on the number of digits in the decimal. Note: A number containing a period followed by exactly three digits is not read as a decimal number but as a full number, following the rules in 5.1.

Number	Reading
16,234	šestnáct celých dvě stě třicet čtyři
3,1415	tři celé jedna čtyři jedna pět
1251,04	tisíc dvě stě padesát jedna celých nula čtyři
1.251,04	tisíc dvě stě padesát jedna celých nula čtyři
2.50	dva tečka padesát
2,50	dva celé padesát
3.141	tři tisíce sto čtyřicet jedna

5.4 Monetary amounts

The following principles are followed for monetary amounts:

- Numbers with zero or two decimal places preceded or followed by the currency markers \$, €, CZK, SKK, HUF, or PLN are read as monetary amounts.
- Numbers with zero or two decimal places preceded or followed by the words “dolar”, “euro”, “libra”, “lib.”, “yen”, “česká koruna”, “slovenská koruna”, “forint” or “zlotý” (singular or plural) are read as monetary amounts.
- Accepted decimal markers are comma and full stop.
- The decimal part (consisting of two digits) in monetary amounts is read as “i nn cent”.
- If the decimal part is “00” it will not be read.

Example	Reading
\$15.00	patnáct dolar
15.00CZK	patnáct českých korun
15.00 euro	patnáct euro
€ 200.50	dvě stě euro i padesát cent
1 000 000 \$	jeden milión dolar

There is also the possibility of writing large amounts as follows:

\$ 1 milión jeden milión dolar

5.5 Ordinal numbers

Numbers are read as ordinals in the following cases:

- The number is followed by a full stop and a space, with no capital letter after the space. Examples: 1. , 12.

- The number is followed by a period, then a month name or one of the month name abbreviations and the number is smaller or equal to 31. The number may be preceded by a day or an abbreviation for a day. Examples: 3. Ledna, 3. Led., Sobota 3. Led.

5.6 Arithmetic operators

Numbers together with arithmetical operators are read according to the examples below.

Expression	Reading
-12	minus dvanáct
+12	plus dvanáct
2*8	dva hvězdička osm
2*8=16	dva krát osm rovná se šestnáct
2/8	dva osmini
8/2=4	osm děleno dva rovná se čtyři
25%	dvacet pět procento

5.7 Mixed digits and letters

If a letter appears within a sequence of digits, the groups of digits will be read as numbers according to the rules above. The letter marks the boundary between the numbers. The letter will also be read.

Examples:

Expression	Reading
77BB84	77 B B 84
0092BC87	0 0 92 B C 87

5.8 Time of day

The colon is used to separate hours, minutes and seconds.

Possible patterns are:

- hh:mm (or h:mm)
- hh:mm:ss (or h:mm:ss)
- hhHmm (or hHmm) ex: 2H45

h = hour, m = minute, s = second.

In pattern a): If the “mm”-part is equal to “00”, this part will not be read. “Hodini” is read after the hours, whether there be minutes or not. “Minut” is added after the “mm” part.

In pattern b): “Sekunt” will be added after the seconds part. If the “ss”-part is equal to “00”, this part will not be read.

Pattern (c) follows the rules for pattern (a).

Expression	Reading
8:15	osum hodin patnáct minut

5.9 Dates

The valid formats for dates are:

- dd.mm.yyyy, dd-mm-yyy and dd/mm/yyyy
- dd.mm.yy, dd-mm-yy and dd/mm/yy

“yyyy” is a four-digit number, “yy” is a two-digit number, “mm” is a month number between 1 and 12 and “dd” a day number between 1 and 31.

Full stop, hyphen and slash may be used as delimiters.

In all formats, one or two digits may be used in the “mm” and “dd” part. Zeros may be used in front of numbers below 10.

Examples of valid formats and their readings:

Type 1: dd.mm.yyyy, dd-mm-yyy and dd/mm/yyyy

10.02.2003 or 10.2.2003 desátéhot únorat dva tisíce tři

10-02-2003 or 10-2-2003 “

10/02/2003 or 10/2/2003 “

Type 2: dd.mm.yy, dd-mm-yy and dd/mm/yy

10.02.03 or 10.2.03 the desátéhot únorat dva tisíce tři

10-02-03 or 10-2-03 “

10/02/03 or 10/2/03 “

Ranges of days and years are also supported.

Examples:

1980-1990 tisíc devět set osmudesát až tisíc devět set devadesát

1939-45 tisíc devět set třicet devět ažčtyřicet pět

14-15 Duben čtrnáct až patnáct Duben

14 až15 Duben čtrnáct až patnáct Duben

6. do 10. února šestého do desátého února

6-10.1 šest až deset Ledna

Other possible formats include:

Sobota, 12. Srpen

Sobota, 12. Srpen 2000

30. Leden 1980

Duben 1999

Valid abbreviations for months: led., ún., bř., břez., dub., květ., červ., červen., srp., zář., říj., list., pros.

Valid abbreviations for days: pon., ut., stř., čtv., pát., sob., ned.

5.10 Phone numbers

In this section the patterns of digits that are recognized as phone numbers are described. In the pronunciation of phone numbers, groups of two and three digits are read as normal numbers, groups of four digits are read out digit by digit, with a pause between the groups.

5.10.1 Ordinary phone numbers

Sequences of digits in the following formats are treated as phone numbers.

The following sequences of digits can be separated by a space, a hyphen or a full stop. The first separation can be a space, a hyphen, or a backslash :

x xxxx xxxx
x xx xx xx xx
xx xxx xxxx
xx xxxx xxx
xxx xx xx xx

The preceding sequences can also all be preceded by a “0”.

The sequence xxx xxx xxx is recognized as a phone format only if preceded by "tel, mob, telefon, mobil, fax", with eventual “:” between these words and the phone number.

5.10.2 International phone numbers

All preceding formats can be recognised if preceded by international prefix (space is optional):

00 x	+ x
00 xx	+ xx
00 xxx	+ xxx

These international prefixes can be followed by an optional (0) before the phone number.

Other recognized formats, that must be preceded by “00” or “+” with optional separator :

xx xxxx xxxx
xxx xxx xxxx
xxx xxxx xxx
xxxx xx xx xx
xxx xxxx xxxx
xxxx xxx xxxx
xxxx xxxx xxx
xxxxx xx xx xx
xxxx xxxx xxxx
xxxxx xxx xxxx
xxxxx xxxx xxx
xxxxxx xx xx xx
xxx xxx xxx xxx

6 How to change pronunciation errors

Words that are not pronounced correctly by the text-to-speech converter can be entered in the user lexicon (see User's guide). In this lexicon, the user enters a phonetic transcription of the word (see chapter 7). Phonetic translations can also be entered directly in the text, using a PRN-tag (see User's guide).

7 Czech Phonetic Text

The Czech uses the Czech subset of the SAMPA phonetic alphabet (Speech Assessment Methods Phonetic Alphabet). The symbols are written with a space between each phoneme.

Only SAMPA may be used in phonetic transcriptions. Symbols not listed here are not valid in phonetic transcriptions and will be ignored if included in the user lexicon or in a PRN tag.

7.1 Consonants

7.1.1 Symbols for the Czech consonants

Symbol	Word	Phonetic text	Comment
b	Bota	\b o1 t a \	
d	Dal	\d a1 l \	
d'	Di	\d' i1 \	
dz	Leckdy	\l e1 dz g d l \	
dZ	Lečbych	\l e1 dZ b i x \	
f	Forma	\f o1 r m a \	
g	Kde	\g d e \	
H	Had	\H a1 t \	
j	Jas	\j a1 s \	
k	Krk	\k r=1 k \	
l	Led	\l e1 t \	
m	Mák	\m a:1 k \	
n	Noc	\n o1 ts \	
n'	Nic	\n' i1 ts \	
N	Banka	\b a1 N k a \	
p	Pes	\p e1 s \	
r	Ret	\r e1 t \	
r'	Tu	\t u1 r' \	
s	Sen	\s e1 n \	
S	Šabach	\S a1 b a x \	
t	Tam	\t a1 m \	
t'	Tati	\t a1 t' i \	
tS	Tavi	\t a1 v i tS \	
ts	Cíl	\ts i:1 l \	
v	Vak	\v a1 k \	
x	Chata	\x a1 t a \	
z	Zub	\z u1 p \	
Z	Žal	\Z a1 l \	
R	Buřeň	\b u1 R e n' \	
?	Glottal stop	\ ? o1 s u m= \	

Table 2 Czech consonants

7.2 Vowels

7.2.1 Symbols for the Czech vowels

Symbol	Word	Phonetic text	Comment
a	Pas	\p a1 s \	Short
a:	Rád	\r a:1 t \	Long
e	Les	\l e1 s \	Short
e:	Lék	\l e:1 k \	Long
i	Myš	\m i1 S \	Short
i:	Píbl	\p i:1 b l= \	Long
o	Rok	\r o1 k \	Short
o:	Móda	\m o:1 d a \	Long
u	Kus	\k u1 s \	Short
u:	Ku v	\k u1 tS u: f \	Long
ou	Bambous	\b a1 m b ou s \	Short
au	Braum	\b r au1 m \	Short
eu	Breu	\b r eu1 \	Short
l=	Bágl	\b a:1 g l= \	Syllabic l
m=	Cozm	\ts o1 z m= \	Syllabic m
r=	Crha	\ts r=1 H a \	Syllabic r
E:	äke	\E:1 k e \	Foreign vowel
2:	Mälmö	\m E:1 l m 2: \	Foreign vowel
y:	Blüml	\b l y:1 m l= \	Foreign vowel
ou:	Kouďův	\k ou:1 d' u: f \	Long
eu:	Kubeův	\k u1 b eu: f \	Long

Table 3 Czech vowels

7.3 Lexical accent

A lexical accent is used to indicate the level of prominence (or emphasis) of a syllable in a word. Practically all words in Czech have a lexical accent even if it does not always serve to differentiate between two different words. It is therefore important to include stress marks when writing phonetic transcriptions.

In the phonetic transcriptions, the lexical accent is indicated by the symbol “1” placed directly after (no space) the accented vowel.

7.4 Pause

An underscore < _ > in a phonetic transcription generates a small pause.

7.5 Glottal stops

A glottal stop, represented by the phonetic symbol /ʔ/, is a small sound which is often used to separate two words when the second word starts with a vowel. This sound can be inserted in a transcription in order to improve the pronunciation.

8 Abbreviations

In the current version of the Czech text-to-speech system, the abbreviations in table 4 below are recognized in all contexts. These abbreviations are mostly case-insensitive (except for those indicated below by “*”).

As previously mentioned, there are also abbreviations for the days of the week and the months, see chapter 5.9.

Abbreviation	Reading
ap.	Apodob e
atd.	Atakdale
cal	Kalorie
ccm	kubický centimetr
cdrom	C D rom
cm	Centimetr
g	Gram
hl.	Hlasovi
hlas.	Hlasové
ič	Infra ervení
jedn.	Jednotka
kb	Kilobajti
kc	Kilokalorie
kg	Kilogram
Kilometry/h	Kilometry za hodinu
km	Kilometr
krychl (.)	Krichlové
kW	Kilowatt
kWh*	Kilovat hodin
l	Litr
m	Metr
metry/s	metry za sekundu
mm	Milimetr
multimed.	Multimédiál e
míle/h	míle za hodinu
paměť.	Pam ova
přík (.)	Příkas
spr.	Spravce
stopy/s	stopy za sekundu
synch	Synchronizace
t	Tuna
zpr.	Zpráva
zrychl.	Zrychlena
tver (.)	tverečné
°C	stupeň Celsia
°F	stupeň Fahrenheita
°K	stupeň Kelvina
cad	kanackí dolar
czk	česká koruna
dkk	danská koruna

eur	Euro
gbp	brická libra
huf	mařarský forint
jpy	Jen
pln	polský zlotí
skk	slovenská koruna
usd	americký dolar

Table 4 Abbreviations

“m”, “g”, “s” and “t” are expanded only when appearing after a number.

Examples

25 m
30 s
45 g

Readings

dvacet pět metri
třicet sekund
čtyřicet pět gramí

9 Web-addresses and email

Web-addresses and email-addresses are read as follows:

- “www” is read as three w’s spelled letter by letter.
- Full stops are read as “tečka”, hyphens as “pomlčka”, underscore (“_”) as “podtržník”, slash (“/”) as “lomeno”.
- “cz, fr” and all the other abbreviations for countries are spelled out letter by letter.
- The “@” is read “zaviná ”.
- Words/strings (including “org”, “com” and “edu”) are pronounced according to the normal rules of pronunciation in the system and in accordance with the lexicon.

String

www.acapela-group.com

<http://www.acapela-group.com>

novak@yahoo.cz

milos_novak@yahoo.cz

Reading

w w w tečka acapela pomlčka group tečka com

h t t p dvojtečka lomeno lomeno w w w tečka acapela
pomlčka group tečka com

novak zaviná yahoo tečka c z

milos dlouha pomlčka novak zaviná yahoo tečka c z