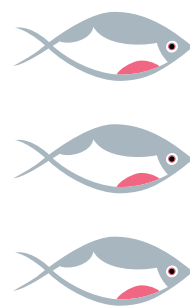


Table of Contents

Braxen 1.0 Documentation	2
Structure	3
Entries	3
Field structure	3
Field 0: Orthography	3
Field 1: Pronunciation	3
Stress	4
Boundaries	4
/e/ Sounds	5
Xenophones (foreign speech sounds)	5
Field 2: Part of speech and morphology	5
Field 3: Language code	6
Field 16: Case	6
Field 26: ID	6
References	7
Appendix A: Field information	8
Appendix B: Phoneme table - Base and IPA	9
Appendix C: Part of speech	12
Appendix D: Language codes	13



Braxen 1.0 Documentation

Braxen 1.0 is a Swedish pronunciation dictionary for speech technology, developed by Swedish Agency for Accessible Media (MTM). It is made public in cooperation with Spr  kbanken Tal.

Structure

The underlying data for Braxen 1.0 is a straight-forward two-dimensional table, with entries for full-form words, abbreviations, and acronyms, and columns (Fields) for various attributes of these entries.

Entries

An entry in Braxen 1.0 is defined by its full-form orthography, and several entries with the same orthography is permissible, but only if they differ in another field.

Field structure

An entry in Braxen 1.0 consists of 27 fields. In its raw, tab-separated form, each field occupies a column, and the order is meaningful, so that columns 0 to 26 have semantics.

In the current version, only a subset of these fields are shared publicly. The structure of the public dictionary remains the same, however, resulting in a number of empty columns (see [Appendix A: Field information](#)).

The most important fields are described in the following.

Field 0: Orthography

The orthography field shows the spelling of the word.

A considerable amount of pragmatics goes into the orthography field. For example, the most frequent casing is used, as different casings can sometimes indicate different pronunciations (see [Field 16: Case](#)).

Table 1. Examples

Orthography	Part of Speech	Transcription
bjšrn	Noun	/b j ũe: rn/
Bjšrn	Proper name	/b j ũe: rn/
BrB	Proper name (abbreviation)	Brottsbalken
BRB	Proper name (acronym)	B R B

Field 1: Pronunciation

The format for the phonetic-phonological transcriptions covers Swedish and common foreign phonemes. Here, this format is referred to as Base.

[Appendix B: Phoneme table - Base and IPA](#) provides a conversion table between Base and IPA. Tools for automatic conversion are included in the [/p5m/scripts](#) directory).

Table 2. Conversion Tools

Script	Task	Call
braxen_conversion.pl	Braxen to IPA(encode) or IPA to Braxen (decode)	perl p5m/scripts/braxen_conversion.pl <encode decode> <infile> <outfile>
validate_braxen.pl	Validates Braxen	perl p5m/scripts/validate_braxen.pl <infile> <outfile>
validate_ipa.pl	Validates IPA	perl p5m/scripts/validate_ipa.pl <infile> <outfile>

Note that the conversion tools are written in Perl. Using e.g. Docker, they can be used without installing a Perl interpreter.

Example of conversion using Docker

Use "encode" for Braxen to IPA format and "decode" for IPA to Braxen format.

```
cd /braxen/repo/root
docker run -v "$PWD":/work -it perl -c "/work/p5m/scripts/braxen_conversion.pl
ENCODE|DECODE /work/INFILE /work/OUTFILE"
```

Example of validation using Docker

```
cd /braxen/repo/root
docker run -v "$PWD":/work -it perl -c "/work/p5m/scripts/validate_braxen.pl
/work/INFILE /work/
OUTFILE"
```

Stress

All Braxen 1.0 words have exactly one main stress, which can be either accent 1 or accent 2. Words with accent 2 also have secondary stress. Stress is marked immediately before the stressed vowel.

Table 3. Stress notation examples

Stress	Notation	Example	Transcription
Main stress accent 1	ō	boll	/b ō l/
Main stress accent 2	ó	bollar, dalbana	/b óo . l ,a r/, /d óa: l - b ,a: . n a/
Secondary stress	,	bollar, dalbana	/b ōo . l ,a r/, /d óa: l - b ,a: . n a/

Boundaries

Word boundaries allow multiple main stresses within an expression.

Table 4. Word boundary examples

Boundary	Notation	Example	Transcription
Word		berg- och dalbana	/b œ rj œ : d óa: l - b , a: . n a/
Compound	-	dalbana	/d óa: l - b ,a: . n a/
Morpheme*	~	transalpin	/t r a n s ~ a l . p 'i: n/
Syllable	.	alpin	/a l . p œ : n/

*The morpheme boundary is optional but can be included when needed.

/e/ Sounds

Braxen 1.0 distinguishes between four /e/ sounds, mainly reflecting a central Swedish pronunciation.

Table 5. /e/ pronunciations

Phoneme	Description	Example
Transcription	e	Semi-open /e/ sound
sett (/s œ t/)	š	Semi-open /e/ sound (dialectal variation)
sštt	/s œ t/	eh
Before stressed syllable in unstressed, open syllable	betona	/b eh . t œ : . n a/
ex	Schwa, used in unstressed syllables	bollen

Xenophones (foreign speech sounds)

The phoneme inventory includes foreign phonemes, so-called xenophones, primarily from English.

Field 2: Part of speech and morphology

Part of speech and morphological data largely follow the Stockholm-Umeå Corpus (SUC) principles (Gustafson-Capkov† & Hartmann, 2006). The UO (foreign word) tag is rarely used, as language codes indicate whether a word is Swedish or foreign.

Different parts of speech or morphological information for the same orthographic form can lead to different pronunciations:

Table 6. Part of speech and pronunciation examples

Orthography	Part of Speech	Transcription
slutet	NN	/s l œ u: . t ex t/
slutet	JJ	/s l óuu: . t ,ex t/

Orthography	Part of Speech	Transcription
planet	NN UTR SIN IND NOM	/p l a . n ɛ: t/
planet	NN NEU SIN DEF NOM	/p l ɛ: . n ex t/

[Appendix C: Part of speech](#) provides a list of selected PoS codes.

Field 3: Language code

Language codes follow the ISO 639-2 standard (Library of Congress, 2017). The language code indicates the intended language of the orthography at the time of pronunciation creation.

Table 7. Examples of different language codes for the same word

Orthography	Language Code	Transcription
Anne	swe	/ɛ̃ n/
Anne	eng	/ɛ̃e n/

[Appendix D: Language codes](#) provides a list of selected language codes.

Field 16: Case

This field indicates case sensitivity:

¥ 1 = Case-sensitive

¥ 0 = Not case-sensitive

Field 26: ID

An internal identifier for each entry.

References

- ¥ Gustafson-Capkov‡, S., & Hartmann, B. (2006). Manual of the Stockholm UmeË Corpus version 2.0.
- ¥ Library of Congress. (2017). ISO 639-2 Language Code List. https://www.loc.gov/standards/iso639-2/php/code_list.php

Appendix A: Field information

Bold fields are shared publicly.

Field	Name	Example
0	orth	bjŠrornas
1	pron	b j Óæ: . r ,u . rn a s
2	posmorph	NN UTR PLU DEF GEN
3	lang	swe
4-15	internal	-
16	case	0
17-25	internal	-
26	id	0060097

Appendix B: Phoneme table - Base and IPA

Base	IPA	Example
Base	IPA	Exempel
p	p	pil
b	b	bil
t	t	bot
rt	!	bort
d	d	bod
rd	"	bord
k	k	kal
g	#	gal
f	f	fal
v	v	val
s	s	sal
rs	\$	fors
h	h	hal
x	%	sjal
c	&	tjat
m	m	matt
n	n	natt
rn	'	barn
ng	(ring
r	r	riv
l	l	liv
rl)	sorl
j	j	jag
w	w	way (eng)
sh	*	she (eng)
zh	+	measure (eng)
z	z	zebra (eng)
dh	,	this (eng)
th	-	thick (eng)
rh	.	road (eng)

Base	IPA	Example
r0	-	father (br. eng)
rx	/	rouge (fre)
tc	t0*	chicken (eng)
dj	d0+	fudge (eng)
xx	x	Bach (ger)
i:	i1	sil
i	2	sill
ih	23	radio
y:	y1	syl
y	4	syll
e:	e1	sel
e	e	sett
eh	e5	betona
ex	6	papper
Š:	71	sŠl
Š	7	sŠtt
ae:	¾1	nŠr
ae	¾	spŠrr
š:	č1	ršn
š	č	ršnn
oe:	ĭ 1	fšr
oe	ĭ	fšrr
u:	u1	bot
u	u	bott
oh	o	bohem
o:	o1	sŀt
o	8	sŀtt
uu:	91	sur
uu	:	surr
uuh	9	butik
uw:	;1	you
uw	;	would
a:	<1	mat

Base	IPA	Example
a	a	matt
aa:	a1	Zlatan
au	a;	paus
eu	7;	euro
ei	e2	way (eng)
ai	a2	lie (eng)
oi	82	boy (eng)
ou	6;	low (eng)
eex	e6	where (eng)
iex	26	here (eng)
uex	; 6	pure (eng)
an	<	sans (fre)
en	7=	pain (fre)
on	o=	pardon (fre)
un	İ =	lundi (fre)

Appendix C: Part of speech

Following (Gustafson-Capkov† & Hartmann, 2006).

Tag	Description	Example
AB	Adverb	inte
DT	Determiner	denna
HA	Interrogative/relative adverb	nŠr
HD	Interrogative/relative determiner	vilken
HP	Interrogative/relative pronoun	som
HS	Interrogative/relative possessive pronoun	vars
IE	Infinitive	att
IN	Interjection	ja
JJ	Adjective	glad
KN	Conjunction	och
NN	Noun	pudding
PC	Participle	utsŠnd
PL	Verb particle	ut
PM	Proper noun	Mats
PN	Pronoun	hon
PP	Preposition	av
PS	Possessive pronoun	hennes
RG	Numeral	tre
RO	Ordinal	tredje
SN	Subjunction	att
UO	Foreign word	the
VB	Verb	kasta

Appendix D: Language codes

Code	Language
swe	Swedish
nob	Norwegian bokmål
nno	Norwegian nynorsk
dan	Danish
ice	Icelandic
fin	Finnish
eng	English
ger	German
fre	French
spa	Spanish
por	Portuguese
dut	Dutch
gre	Greek
rus	Russian
cze	Czech
chi	Chinese
jap	Japanese
kor	Korean
tha	Thai
swa	Swahili
ara	Arabic