

 README.md

Grapheme-to-Phoneme models for Norwegian

Version

20200601: 1.0

Introduction

This resource contains Grapheme-to-Phoneme (G2P) models for Norwegian, to be used with the G2P engine [Phonetisaurus](#). The G2P models can be used to generate pronunciation lexica from word lists. For more information on how to do that, consult the [Phonetisaurus repo](#).

The models are trained on the Norwegian pronunciation lexicon for ASR, originally made by the defunct company Nordisk språkteknologi (NST), currently distributed by the [National Library of Norway](#).

Two models have been developed. One is trained on a full version of the lexicon, including phones, marking of primary and secondary stress, and tone. The other is trained on a simplified version where tonal markings and markings of secondary stress are removed.

Content

- *train/*: contains the models, as well as auxiliary files used by Phonetisaurus
 - *model-wtone-nob.fst* contains full tone and stress specifications
 - *model-notone-nob.fst* lacks tone and secondary stress
 - *lexica/*: contains various lexica used for training and testing
 - *NST-total_train.dict* is the training set for *model-wtone-nob.fst*. It contains 612 366 word-transcription pairs (WTP) and constitutes 90% of the unique WTPs in the NST lexicon.
 - *NST-total_test.dict* is the test set for *model-wtone-nob.fst*. It consists of the remaining 10% of the unique WTPs in the NST lexicon, which have been randomly selected
 - *NST-total-notone_nosecstress_train.dict* is the training set for *model-notone-nob.fst*. It is equal to *NST-total_train.dict*, but markings of tone and secondary stress have been removed
 - *NST-total-notone_nosecstress_test.dict* is the test set for *model-notone-nob.fst*. It is equal to *NST-total_test.dict*, but markings of tone and secondary stress have been removed
 - *NST-total_test_predicted.dict* is the test set with tones and secondary stress with transcriptions predicted by the G2P system
 - *NST-total_test_notone_predicted.dict* is the test set without tones and secondary stress with transcriptions predicted by the G2P system
- g2p_stats.py* is the evaluation script used in this project.

Transcription standard

Although the original NST lexicon uses X-SAMPA as a transcription standard, an equivalent standard is used in this project, which is easier to read by humans, *NoFAbet*. NoFAbet is in part based on [2-letter ARPAbet](#) and is made by [Nate Young](#) for the National Library of Norway in connection with the development of *NoFA*, a forced aligner for Norwegian, soon to be released in Språkbankens resource catalogue.

X-SAMPA-NoFAbet equivalence table

X-SAMPA	NoFAbet	Example
A:	AA0	bad
{:	AE0	vær
{	AEH0	vært
{*l	AEJ0	sei

X-SAMPA	NoFAbet	Example
E* _u 0	AEW0	sau
A	AH0	hatt
A* _l	AJ0	kai
@	AX0	behage
b	B	bil
d	D	dag
e:	EE0	lek
E	EH0	penn
f	F	fin
g	G	gul
h	H	hes
l	IH0	sitt
i:	II0	vin
j	J	ja
k	K	kost
C	KJ	kino
l	L	land
l=	LX0	
m	M	man
m=	MX0	
n	N	nord
N	NG	eng
n=	NX0	
o:	OA0	rå
O	OA _H 0	gått
2:	OE0	løk
9	OE _H 0	høst
9* _Y	OE _J 0	køye
U	OH0	f*ort
O* _Y	OJ0	konvoy
u:	OO0	bod
@ _U	OU0	show
p	P	pil
r	R	rose
d`	RD	rekord
l`	RL	perle
l`=	RLX0	

X-SAMPA	NoFAbet	Example
n`	RN	barn
n`=	RNX0	
s`	SJ	pers
t`	RT	stort
r=	RX0	
s	S	sil
S	SJ	sju
s=	SX0	
t	T	tid
u0	UH0	russ
u0 j	UH0_J	Anhui
}:	UU0	hus
v	V	vase
w	W	Washington
Y	YH0	nytt
y:	YY0	ny

Unstressed syllables are marked with a 0 after the vowel or consonant syllable nucleus. The nucleus is marked with a 1 for tone 1 and a 2 for tone 2. Secondary stress is marked with 3. In the material without tone and stress marking, all 3s are replaced by zeros and all 2s with 1s.

For compatibility with NoFA, retroflex s is rendered as SJ instead of RS, which means that there is no distinction between postalveolar and retroflex s in the transcriptions.

Evaluation

Model	Word Error Rate	Phoneme Error Rate
model-wtone-nob.fst	14.29	2.76
model-notone-nob.fst	10.44	2.00

The PER calculation is borrowed from [this tutorial](#).

Usage

The models created in this project are public property with the license [CC0](#). See also [Phonetisaurus' license](#).