

LESMEG.md

Grafem-til-fonem-modeller for norsk

Versjon

20200601: 1.0

Innledning

Denne ressursen innegolder *grafem-til-fonem*-modeller (GtP-modeller) for norsk som er tilpassa GtP-systemet [Phonetisaurus](#). GtP-modellene brukes til å generere uttaleleksika fra ordlister. Mer informasjon om hvordan man gjør det, finner man på sidene til Phonetisaurus.

Modellene er trent på the norske uttaleleksikonet for talegjenkjenning som ble utvikla av det forenævrende selskapet Nordisk språkteknologi (NST). Dette leksikonet distribueres nå av [Nasjonalbiblioteket](#).

Vi har utvikla to modeller. Den første er trent på en full versjon av leksikonet, og inkluderer segmenter, markering av primær- og sekundærtrykk og tonemer. Den andre er trent på en forenkla versjon av leksikonet der toner og sekundærtrykk ikke er fjerna.

Innhold

- *train/*: inneholder modellene, samt andre filer Phonetisaurus bruker
 - *model-wtone-nob.fst* er trent på fullversjonen av leksikonet
 - *model-notone-nob.fst* er trent på en versjon av leksikonet uten toner og sekundærtrykk
- *lexica/*: inneholder leksika for trening og testing
 - *NST-total_train.dict* er treningssettet for *model-wtone-nob.fst*. Det inneholder 612 366 ortografi-transkripsjon-par (OTP), og utgjør 90% av alle unike OTP-er i NST-leksikonet
 - *NST-total_test.dict* er testsettet for *model-wtone-nob.fst*. Det inneholder de gjenstående 10% av OTP-ene i NST-leksikonet, og er tilfeldig plukka ut
 - *NST-total-notone_nosecstress_train.dict* er treningssettet for *model-notone-nob.fst*. Det er likt *NST-total_train.dict*, men markeringa av toner og sekundærtrykk er fjerna
 - *NST-total-notone_nosecstress_test.dict* er testsettet for *model-notone-nob.fst*. Det er likt *NST-total_test.dict*, men markeringa av toner og sekundærtrykk er fjerna
 - *NST-total_test_predicted.dict* er testsettet med transkripsjoner predikert av Phonetisaurus. Transkripsjonene har toner og sekundærstress
 - *NST-total_test_notone_predicted.dict* er testsettet med transkripsjoner predikert av Phonetisaurus. Transkripsjonene har ikke toner og sekundærstress *g2p_stats.py* er evalueringsskriptet som er brukt i evalueringa under.

Transkripsjonsstandard

Selv om den originale versjonen av NST-leksikonet bruker transkripsjonsstandarden X-SAMPA, har vi valgt å bruke en annen, ekvivalent standard som er lettere å lese og skrive for mennesker, *NoFAbet*. NoFAbet er til dels basert på [ARPAbet](#) og er laga av [Nate Young](#) på oppdrag for Nasjonalbiblioteket i forbindelse med utviklinga av *NoFA*, en norsk modell for forced alignment, som snart vil være tilgjengelig i Språkbankens ressurskatalog.

X-SAMPA-NoFAbet-ekvivalenstabell

X-SAMPA	NoFAbet	Eksempel
A:	AA0	bad
{:	AE0	vær
{	AEH0	vært

X-SAMPA	NoFAbet	Eksempel
{*I	AEJ0	sei
E*u0	AEW0	sau
A	AH0	hatt
A*I	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
I	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	OAH0	gått
2:	OE0	løk
9	OEH0	høst
9*Y	OEJ0	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

X-SAMPA	NoFAbet	Eksempel
l`=	RLX0	
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

Trykklette stavelser er merka med 0 etter stavelseskjernen. Stavelseskjernen er merka med 1 ved tonem 1 og 2 ved tonem 2. Sekundærstress merkes med 3. I materialet uten toner og sekundærstress er alle 3-ere bytta ut med 0 og alle 2-ere med 1.

For å være kompatibel med NoFA er retrofleks s skrevet som SJ, ikke RS, noe som innebærer at det ikke er noen kontrast mellom retrofleks og postalveolar s-lyd i dette materialet.

Evaluering

Modell	Word Error Rate	Phoneme Error Rate
<i>model-wtone-nob.fst</i>	14.29	2.76
<i>model-notone-nob.fst</i>	10.44	2.00

Utregninga av Phoneme Error Rate er henta [herfra](#).

Bruk

Modellene utvikla i dette prosjektet er offentlig eiendom med lisensen [CC0](#). De kan brukes til et hvilket som helst formål. Se for øvrig [lisensen til Phonetisaurus](#).