

Grafem-til-fonem-modeller for norsk bokmål

Versjon: 2.0 (2024-02-09)

Dette repoet inneholder G2P-modeller for norsk bokmål, som produserer fonemiske transkripsjoner for talenære uttalevarianter (som i vanlige samtaler: **spoken**) eller skriftnære uttalevarianter (som i tekstopplasing **written**) for 5 forskjellige dialektområder:

1. Østnorsk (**e**)
2. Sørvestnorsk (**sw**)
3. Vestnorsk (**w**)
4. Trøndersk (**t**)
5. Nordnorsk (**n**)

Oppsett

Følg installasjonsinstruksjoner fra [Phonetisaurus](#). Du trenger kun å kjøre koden i disse avsnittene:

- "Next grab and install OpenFst-1.7.2"
- "Checkout the latest Phonetisaurus from master and compile without bindings"

Data

Uttaleleksikonene som ble brukt til å trene G2P-modellene er fritt tilgjengelige fra Språkbanken ressurskatalog: [NB Uttale](#). De kan også genereres opp med data og kode fra Github-repoet [Sprakbanken/nb_uttale](#).

Innhold

- **models/**: inneholder modellene, samt andre filer Phonetisaurus bruker
 - **nb_*.fst**: Modellene som kan kjøres med **phonetisaurus-apply**. Utfyllingen til ***** er en streng med dialekt og uttalevariant, f.eks. **e_spoken** eller **t_written**.
 - **nb_*.o8.arpa**: 8-gram-modeller for fonemsekvenser som Phonetisaurus bruker i treningsprosessen.
 - **nb_*.corpus**: Sammenstilte (eng. *aligned*) grafemer og fonemer fra leksikonene.
- **data/**: inneholder leksika for trening og testing, samt utdata fra modellene på testsettet.
 - **NB-uttale*_train.dict** er treningssettene for **models/nb_*.fst**. De inneholder 543 495 ortografi-transkripsjon-par (OTP), og utgjør 80% av alle unike OTP-er i leksikonet.
 - **NB-uttale*_test.dict** er testsettene for **models/nb_*.fst**. De inneholder de gjenstående 20% av OTP-ene i leksikonet, altså 135 787 OTP-er.
 - **predicted_nb_*.dict** er testsettet med transkripsjoner predikert av modellene.
 - **wordlist_test.txt** er ordlisten til testsettet, som modellen kjøres på og genererer transkripsjonsforslag for.
- **evaluate.py**: evalueringsskript som er implementert på nytt for å evaluere disse modellene.
- **g2p_stats.py**: evalueringsskript fra V1.0, som kan brukes for å sammenligne disse modellene med NST-modellene (med og uten markering av trykk og tone) i versjon 1.
- **LICENSE**: Lisensteksten for CC0, som denne ressursen distribueres med.

Bruk

```
phonetisaurus-apply --model models/nb_e_spoken.fst --word_list
data/wordlist_test.txt -n 1 -v > output.txt
```

- Inndata (`--word_list`) burde være en tekstfil med linjesepererte ord. Se for eksempel `data/wordlist_test.txt`.
- De trente modellene er `.fst`-filer i `models`-mappen. I samme mappe ligger også `.corpus`-filer og 8-gram-modeller for fonemsekvenser (`.arpa`-filer) som kommer fra treningsprosessen til Phonetisaurus.
- `-n`-argumentet angir hvor mange av de mest sannsynlige transkripsjonsforslagene som skal skrives til utdata.

Evaluerings

Det er to skript i repoet for å regne ut Word Error Rate (WER) og Phoneme Error Rate (PER), og som gir litt ulike svar pga. implementasjonen av utregningene.

`evaluate.py`

Regner statistikk for alle modellene dersom ingen argumenter er gitt.

```
python evaluate.py
```

Med `-l`-argumentet kan du få statistikk for spesifikke modeller, f.eks. `-l e_spoken`.

- **WER** er antall feil transkripsjoner fordelt på totalt antall ord i testdata. 1 feil = 1 ord.
- **PER** er antall feil fonemer for alle transkripsjonene, fordelt på totalt antall fonemer i testdata. 1 feil = 1 fonem.

Model	Word Error Rate	Phoneme Error Rate
<i>nb_e_written.fst</i>	13.661238654564869	1.9681178920293207
<i>nb_e_spoken.fst</i>	13.72501038144391	1.9832518286152074
<i>nb_sw_written.fst</i>	13.240048644480037	1.8396612197218096
<i>nb_sw_spoken.fst</i>	16.422702734768936	2.426312206336983
<i>nb_w_written.fst</i>	13.240048644480037	1.8396612197218096
<i>nb_w_spoken.fst</i>	16.892833837574894	2.5064155890730686
<i>nb_t_written.fst</i>	13.736133357062347	1.98774986044724
<i>nb_t_spoken.fst</i>	16.47992288013051	2.5809178688066843
<i>nb_n_written.fst</i>	13.736133357062347	1.98774986044724

Model	Word Error Rate	Phoneme Error Rate
<i>nb_n_spoken.fst</i>	17.22590930999963	2.8209779677747715

g2p_stats.py

Regner ut WER og PER for to inputfiler:

1. Referansefila /testdata, f.eks. `data/NB-uttale_e_spoken_test.dict`
2. Modellens transkripsjonsforslag, f.eks. `output.txt` fra kommandoen i [Bruk](#), eller `data/predicted_nb_e_spoken.dict`.

- **WER** er antall feil transkripsjoner fordelt på totalt antall ord i modellens forslag. 1 feil = 1 ord.
- **PER** er summen av PER for hver transkripsjon, fordelt på totalt antall ord i modellens forslag. 1 feil = 1 fonem.

OBS: Denne metoden tar ikke hensyn til ulik lengde på transkripsjonene. En transkripsjon med 2 fonemer der 1 er feil får en 50% PER, mens et ord på 10 fonemer med 1 feil får 10% PER, og gjennomsnittet blir 35%. Om man regner antall feil på totalt antall fonemer, blir gjennomsnittet for de to 16,7%.

```
python g2p_stats.py data/NB-uttale_e_spoken_test.dict
data/predicted_nb_e_spoken.dict
# WER: 14.049209423582523
# PER: 2.714882650391985
```

Model	Word Error Rate	Phoneme Error Rate
<i>nb_e_written.fst</i>	13.97114598599277	2.7038190765903214
<i>nb_e_spoken.fst</i>	14.049209423582523	2.714882650391985
<i>nb_sw_written.fst</i>	13.541060631724685	2.5423757844377284
<i>nb_sw_spoken.fst</i>	16.729141964989285	3.34063477772742
<i>nb_w_written.fst</i>	13.541060631724685	2.5423757844377284
<i>nb_w_spoken.fst</i>	17.186475877661337	3.4137304874392114
<i>nb_t_written.fst</i>	14.059519688924565	2.7190289235234104
<i>nb_t_spoken.fst</i>	--	--
<i>nb_n_written.fst</i>	14.059519688924565	2.7190289235234104
<i>nb_n_spoken.fst</i>	--	--

OBS: Modellforslagene fra *t_spoken* og *n_spoken* har ikke samme antall ord som testdata, som gjør at skriptet krasjer.

Transkripsjonsstandard

G2P-modellene har blitt trent på data med transkripsjonsstandarden NoFabet , som er lettere å lese for mennesker enn X-SAMPA. NoFabet er delvis basert på [ARPAbet](#) og ble utviklet for Nasjonalbiblioteket av [Nate Young](#) i forbindelse med utviklingen av [NoFA](#), en forced alignment-modell for norsk. Tabellen nedenfor viser ekvivalente symboler for notasjonene X-SAMPA, IPA og NoFabet.

X-SAMPA-IPA-NoFabet konverteringstabell

X-SAMPA	IPA	NoFabet	Example
A:	ɑ:	AA0	bad
{:	æ:	AE0	vær
{	æ	AEH0	vært
{*I	æI	AEJ0	sei
E*u0	æu	AEW0	sau
A	ɑ	AH0	hatt
A*I	ɑI	AJ0	kai
@	ə	AX0	behage
b	b	B	bil
d	d	D	dag
e:	e:	EE0	lek
E	ɛ	EH0	penn
f	f	F	fin
g	g	G	gul
h	h	H	hes
I	ɪ	IH0	sitt
i:	i:	IIO	vin
j	j	J	ja
k	k	K	kost
C	ç	KJ	kino
l	l	L	land
l=	ɫ	LX0	
m	m	M	man
n	n	N	nord
N	ŋ	NG	eng

X-SAMPA	IPA	NoFAbet	Example
n=	ŋ	NX0	
o:	o:	OA0	rå
O	ɔ	OAH0	gått
2:	ø:	OE0	løk
9	œ	OEH0	høst
9*Y	œy	OEJ0	køye
U	u	OH0	f*ort
O*Y	ɔy	OJ0	konvoy
u:	u:	OO0	bod
@U	oʊ	OU0	show
p	p	P	pil
r	r	R	rose
d`	d̥	RD	rekord
l`	l̥	RL	perle
l`=	l̥	RLX0	
n`	ŋ̥	RN	barn
n`=	ŋ̥	RNX0	
s`	ʃ̥	RS	pers
t`	t̥	RT	stort
s	s	S	sil
S	ʃ	SJ	sju
t	t	T	tid
u0	ʊ	UH0	russ
}:	ʊ:	UU0	hus
v	u	V	vase
w	w	W	Washington
Y	y	YH0	nytt
y:	y:	YY0	ny

Trykklette stavelser er markert med 0 etter stavelseskjernen. Stavelseskjernen er markert med 1 ved tonem 1 og 2 ved tonem 2. Sekundærstress merkes med 3.

Lisens

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