

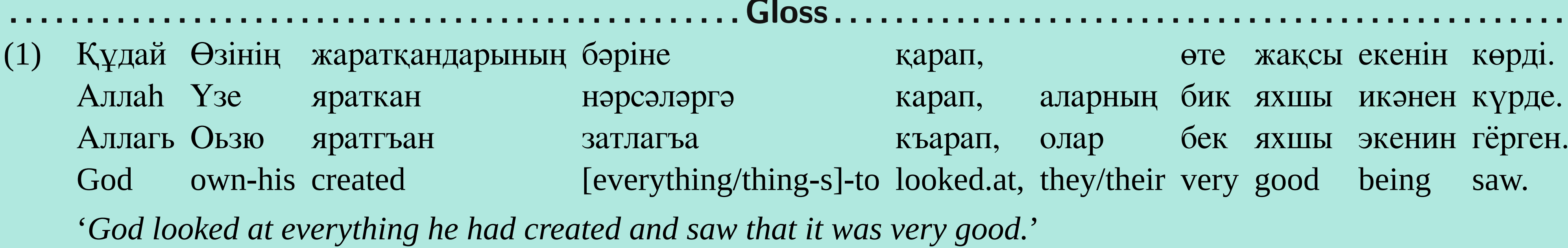


Francis M. Tyers

Special thanks to
Aida Sundetova
sun27aida@gmail.com



Example output



- ### Output

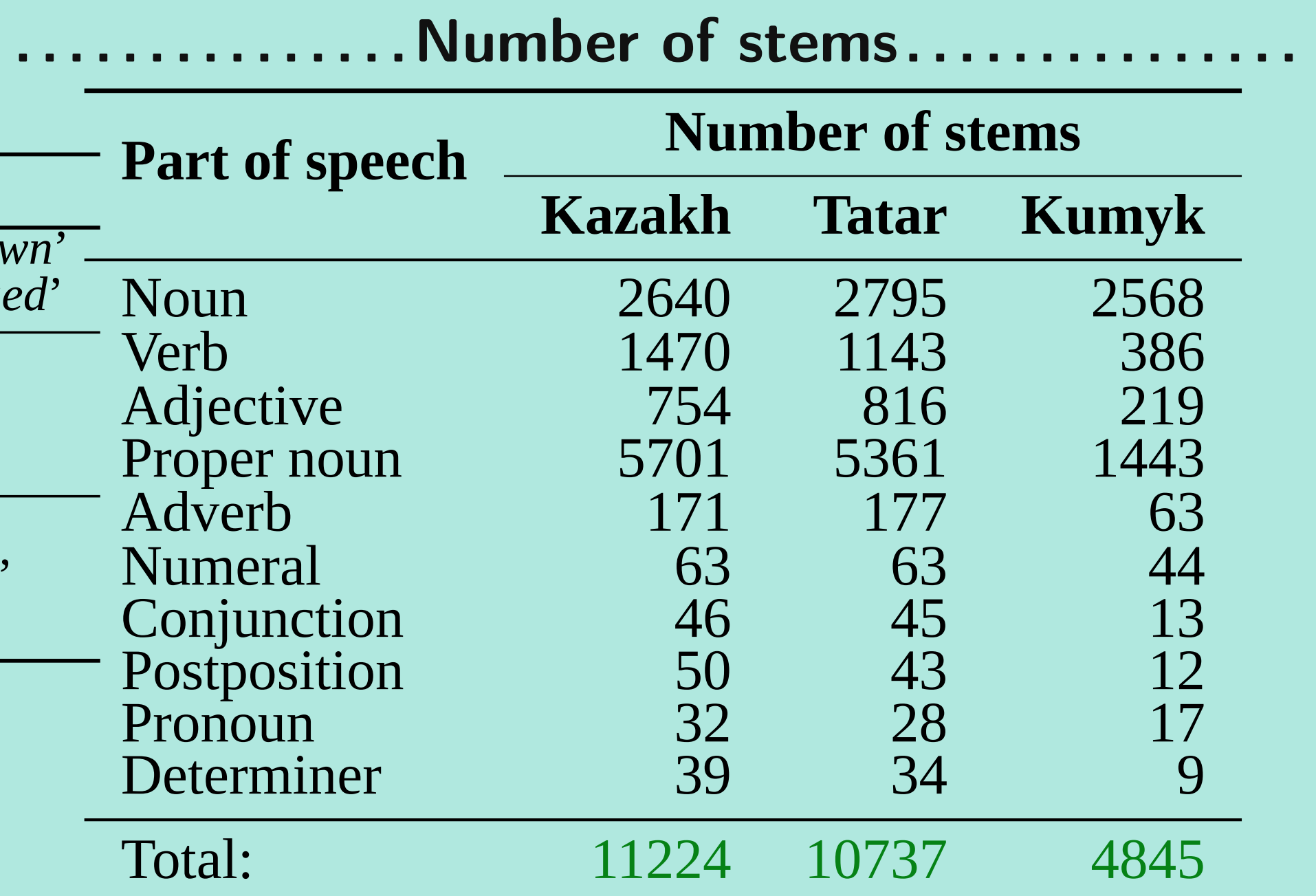
Kazakh (kaz)	Tatar (tat)	Kumyk (kum)
Құдай Өзінің жаратқандарының бәріне қарап, өте жақсы екенін көрді.	Аллаһ Үзе яраткан нәрсәләргә карап, аларның бик яхшы икәнән күрде.	Аллагъ Обзю яратгъан затлагъа къарап, олар бек яхшы экенин гёрген.
Құдай<n><nom> Өз<prn><ref>px3sp<gen> жарат<v><tv><ger_past>pl<px3sp><gen> бәрі<prn><qnt><px3sp><dat> қара<v><tv><gna_perf> ,<cm> — өте<adv> жақсы<adj> е<cop><ger_past>px3sp<acc> көр<v><tv><ifi>p3<sg> .<sent>	Аллаһ<n><nom> Үз<prn><ref>px3sp<nom> ярат<v><tv><gpr_past> нәрсә<n><pl><dat> кара<v><tv><gna_perf> ,<cm> алар<prn><pers>p3<pl><gen> бик<adv> яхшы<adj> и<cop><ger_past>px3sp<acc> күр<v><tv><past>p3<sg> .<sent>	Аллагъ<n><nom> Обз<prn><ref>px3sp<nom> ярат<v><tv><gpr_past> зат<n><pl><dat> къара<v><tv><gna_perf> ,<cm> олар<prn><pers>p3<pl><nom> бек<adv> яхшы<adj> э<cop><ger_past>px3sp<acc> гёр<v><tv><past>p3<sg> .<sent>

..... Tagset

- | | | | | | | | |
|-------|--------------|-------|--------------|--------|-----------------------------|------------|----------------------------|
| <n> | Noun | <p3> | Third person | <pers> | Personal | <px3sp> | 3rd person poss. |
| <v> | Verb | <pl> | Plural | <cm> | Comma | | (Singular/Plural) |
| <prn> | Pronoun | <nom> | ‘Nominative’ | <sent> | Sentence | <gna_perf> | Verbal adverb |
| <det> | Determiner | <gen> | Genitive | <past> | Past (General) | | (Perfect) |
| <adj> | Adjective | <acc> | Accusative | <ifi> | Past
(Eyewitness/Recent) | <gpr_past> | Verbal adjective
(Past) |
| <adv> | Adverb | <dat> | Dative | | | <ger_past> | Verbal noun (Past) |
| <iv> | Intransitive | <qnt> | Quantifier | | | | |
| <tv> | Transitive | <ref> | Reflexive | | | | |

Evaluation

- | | letters | values | examples |
|------------|---------|--------------------------------------|--|
| kaz | и, у, ю | /əj, əw, jəw/
/əj, əw, jəw/ | киюда 'in the process of chopping'
киюде 'in the process of getting' |
| tat | е | ə / C _
/j/+ы
/j/+ə | дәресләр 'lessons'
еллар 'years'
егетләр 'boys' |
| kum | ё, ю | /ø, y/ / C _
/jø, jy/
/jo, ju/ | гёзлер 'eyes', гюнлер 'days'
юреклер 'hearts', ёнкюлер 'dar'
юлдузлар 'stars', ёллар 'roads' |
- solution: hairy twol rules for majority of cases
 - unaccounted-for words marked harmony-forcing char
 - adjust rules for harmony-forcing characters



Loanwords

- Kazakh - елді / Назарбаевты Tatar - галимнр, артистлар Kumyk - сёзлер, самолётлар - separate continuation lexicon - result: super messy twol rules
- **Acronyms and numerals** отыздан, бестен handled by twol 30-дан, 5-тен not handled by twol

$$5\{a\}\{c\} > -\{D\}\{A\}_H$$

..... **Test corpora**

- | | Wikipedia | News | Religion |
|------------|-----------|---------------------|-------------------------|
| kaz | Wikipedia | azattyk.org | Quran + Bible |
| tat | Wikipedia | tat.tatar-inform.ru | Quran + New Testament |
| kum | — | yoldash.etnosmi.ru | Genesis + New Testament |

- ### Evaluation measures
- **Naïve coverage** - percentage of surface forms in a given corpus receiving ≥ 1 analysis
 - **Mean ambiguity** - average number of analyses for each surface form found in analysed corpus
 - **Precision** - of a form's analyses, % correct
 - **Recall** - % of analyses provided by transducer that are correct for a form, by comparing against a gold standard

Evaluation results				
Language	Corpus	Tokens	Coverage (%)	Amb.
Kazakh	Wikipedia	25.6M	85.61 ± 1.37	0.00
	News	3.8M	92.12 ± 2.72	0.00
	Religion	851K	92.49 ± 1.66	0.00
(r50547)	Average		90.07 ± 1.91	0.00
Tatar	Wikipedia	159K	86.35 ± 2.17	0.00
	News	5.2M	89.75 ± 0.07	0.00
	Religion	382K	91.25 ± 2.55	0.00
(r50260)	Average		89.12 ± 1.60	0.00
Kумык	News	286K	91.10 ± 0.86	0.00
	Religion	227K	92.47 ± 1.03	0.00
(r50300)	Average		91.78 ± 0.94	0.00

- selected & proofed unique random surface forms from news corpora

Language	Forms	Precision (%)	Recall (%)
Kazakh	1000	98.61	57.98
Tatar	1000	95.03	85.65
Kumyk	500	96.57	69.11

- Disambiguation (already exists for Kazakh)
- More stems (especially Kumyk)
- Machine translation between these languages