

FINITE-STATE MORPHOLOGICAL TRANSDUCERS FOR THREE KYPCHAK LANGUAGES

Jonathan North Washington Indiana University jonwashi@indiana.edu

Ilnar Salimzyanov

Some University ilnar.salimzyan@gmail.com

Francis M. Tyers

Also special thanks to Aida Sundetova UiT Norgga Árktalaš Universitehta francis.tyers@uit.no email@email

Gloss.







• Turkic languages (SOV, agglutinative, vowel harmony)				
	Kazakh	Tatar	Kumyk	
	population of speakers			
number pronunc primary secondary	8M-12M /qazaq/ Kazakhstan China, Mongolia	5.4M /tɒtɑr/ Tatarstan	430K /qumuq/ Dagestan	
	external influences	5		
Mongolic Oghuz Persian	moderate — heavy	light light heavy	light moderate heavy	

Morphological transducers

heavy

Russian

Take a surface form, and produce valid lexical form(s)

...... Morphological transducers

heavy

heavy

- Take a lexical form, and produce valid surface form(s) 'алдым' ал<v><tv><ifi><p1><sg>, алд<n><px1sg><nom>
- Transducers for Turkic languages......
- Turkish (Çöltekin, 2010; Öflazer, 1994)
- Crimean Tatar (Altıntaş, 2001)
- Turkmen (Tantuğ et al., 2006)
- Kyrgyz (Tyers et al., 2012) • GPL (=free and open)!
- Framework: HFST.....
- Reimplementation of Xerox FST formalisms (lexc and twol)
- Also provides a wrapper around popular free/opensource FST toolkits: SFST, OpenFST, and Foma Development effort.....

Morphotactics

- Morphological & orthographical words
- өнүктүрөбүзбү ? 'will we develop [it]?' ӨНҮК<v><tv><caus><aor><pl>><pl>+бы<qst>
- келатсаң 'if you come'
- Keл<v><iv><prt impf>+жат<vaux><gna cnd><p2><sg>
- ...Irregular [noun + possessive + case] forms...
- Some combinations of possessive and case morphemes are distinct (i.e., not formed simply by concatenation):

case	form	1SG	2SG	3SP
nom	<u>—</u>	-(I)M	-(I)ң	-(S)I
acc	-NI	-(І)мдІ	-(I)ңдI	-(S)IH
gen	-NIH	-(І)мдІн	-(І)ңдІн	- (S)ІнІн
loc	-DA	-(І)мдА	-(І)ңдА	-(S) І ндА
abl	-DAн	-(І)мдАн,	-(I)ндAн,	-(S) І нАн
		-(І)мАн	-(І)ңАн	
dat	-GA	-(I) MA	-(І)ңА	-(S)IHA

- Trade-off:
- morphophon. complicateder, morphotactics simpler
- underlying form used: {S}{I}{n}
- phonological rules delete {n}, {S} by context

..... Noun-noun compounds.....

 one type of N-N compunds: N2 has <px3> and related morphology

LEXICON N-INFL-3PX-COMPOUND %<n%>:%>%{S%}%{I%}%{n%} GEN-POS ; LEXICON Nouns аба% ырайы:аба% ырай N-INFL-3PX-COMPOUND ; ! "weather"

чакыруу% кагазы:чакыруу% кагаз N-INFL-3PX-

COMPOUND ; ! "invitation"

Example output

Аллагь Оьзю яратгъан затларгъа къарап, олар бек яхшы экенин гёрген. thing-s-to look-having, they very good being own-his made God

'Go ····· Kazakh	d looked at everything he h		Output		Kumyk	
Құдай Өзінің жаратқандарының бәріне қарап, өте жақсы екенін көрді.		Алла	Аллаһ Үзе яраткан нәрсәләргә карап, аларның бик яхшы икәнен күрде.		Аллагь Оьзю яратгъан затлагъа къарап, олар бек яхшы экенин гёрген.	
жарат <v>бәрі<pre>prr қара<v></v></pre> ,<cm></cm></v>	<ref><px3sp><gen> <tv><qer_past><pl><qnt><px3sp>< tv><qnt><px3sp><dat> <tv><gna_perf> </gna_perf></tv></dat></px3sp></qnt></px3sp></qnt></pl></qer_past></tv></gen></px3sp></ref>	Yз<р ярат нәрс кара , <cm алар бик< яхшы и<со</cm 	<prn><pers><p3><pl>< adv> <adj> p><ger_past><px3sp>< v><tv><past><p3><sg></sg></p3></past></tv></px3sp></ger_past></adj></pl></p3></pers></prn>	gen>	0ьз <pr зат<n="" ярат<v=""> къара< ,<cm> олар<р бек<ad э<cop="" яхшы<а=""></ad></cm></pr>	dj> <ger_past><px3sp><acc> <tv><past><p3><sg></sg></p3></past></tv></acc></px3sp></ger_past>
			Tagset			
<n><n><v><v><det><det><det><det><det><det><det><det< th=""><th>Noun Verb Determiner Adjective Adverb Intransitive Transitive</th><th><nom> <pen><pen><acc><px3sp><past><ifi><prn><prn></prn></prn></ifi></past></px3sp></acc></pen></pen></nom></th><th>'Nominative' Genitive Accusative 3rd person poss. Past Past</th><th><g <g (Singular⁄d <c< th=""><th>na_perf</th><th>t>Verbal noun (Past) f>Verbal adverb (Perfect) t>Verbal adjective (Past) Comma Sentence</th></c<></g </g </th></det<></det></det></det></det></det></det></det></v></v></n></n>	Noun Verb Determiner Adjective Adverb Intransitive Transitive	<nom> <pen><pen><acc><px3sp><past><ifi><prn><prn></prn></prn></ifi></past></px3sp></acc></pen></pen></nom>	'Nominative' Genitive Accusative 3rd person poss. Past Past	<g <g (Singular⁄d <c< th=""><th>na_perf</th><th>t>Verbal noun (Past) f>Verbal adverb (Perfect) t>Verbal adjective (Past) Comma Sentence</th></c<></g </g 	na_perf	t>Verbal noun (Past) f>Verbal adverb (Perfect) t>Verbal adjective (Past) Comma Sentence

Morphophonology

Plural

<p3>

<pl><pl></pl>

Desonorisation.

<qnt>

<itg>

• {N} desonorises to д after a consonant алма- $\{N\}\{I\}$ \rightarrow алманы 'apple-ACC' $cыp-{N}{I} → сырды 'secret-ACC'$

Third person

• $\{L\}$ desonorises to μ after cons. of sonority $\leq l$ сыр- $\{L\}\{A\}$ р → сырлар 'secret—PL' кыз- $\{L\}\{A\}p \rightarrow$ кыздар 'girl–PL'

"L Desonorisation"

%{L%}:д <=> :VoicedLowSonCns %>:

"N Desonorisation"

%{N%}:д <=> :VoicedCns %>: ;

• Turn {y} into a harmonised high vowel when a vowel doesn't follow the following consonant: $myp{y}H \rightarrow mypyh 'nose'$ $мур{y}H+{I}M \rightarrow мурдум 'my nose'$

%{y%}:Vy <=> [:LastVowel :Cns* :Cns]/[:0] __ [:Cns [.#. | :Cns]]/[:0 | %>:] ; where Vy in (иүииүыыууыуу) LastVowel in (иүеэөяаёоыюу) matched ;

......й+vowel letters.....

- [a o y] become [яёю] after й and й deletes
- й incorporated into the context of many rules
- + separate rules to change the characters
- + a rule to delete the original й

"Deletion of й before yoticised vowels" й:0 <=> __ [:YotVow]/[:0 | %>:] ;

Further information

- The transducer is available from apertium's svn repo: info at http://wiki.apertium.org/wiki/apertium-kir
- Turkic RBMT mailing list (>25 subscribers): apertium-turkic@lists.sourceforge.net Feel free to post in any language!
- See our paper in the LREC 2012 proceedings
- And feel free to contact the authors any time!

Evaluation

Part of speech		Number of stems			
Part of speech	Kazakh	Tatar	Kumyk		
Noun	2640	2795	2568		
Verb	1470	1143	386		
Adjective	754	816	219		
Proper noun	5701	5361	1443		
Adverb	171	177	63		
Numeral	63	63	44		
Conjunction	46	45	13		
Postposition	50	43	12		
Pronoun	32	28	17		
Determiner	39	34	9		
Total:	11224	10737	4845		
Test corpora					
Г	. 1		2012100		

Encyclop	kaz	wpdump	20131006
	tat	wpdump	20130225
	kum	—	—
News	kaz	RFE/RL	azattyq.org 2010
	tat	Татар-информ	tat.tatar-inform.ru
	kum	Ёлдаш	yoldash.etnosmi.ru
Religion	kaz	quran + bible	kkitap.net, kuran.k
	tat	quran + nt	ibt.org.ru, tanzil.ne
	kum	genesis + nt	ibt.org.ru

- split into 10 equal parts; coverage calculated over each separately; standard deviation of mean calculated
- Naïve coverage percentage of surface forms in a given corpus receiving ≥ 1 analysis (surface forms may have missing analyses)
- **Mean ambiguity -** average number of analyses for each surface form found in analyed corpusCoverage results (as of r36739)......

Language	Corpus	Tokens	Coverage (%)
	Wikipedia	25.6M	85.61 ± 1.37
Kazakh	News	3.8M	92.12 ± 2.72
NdZdKII	Religion	851K	92.49 ± 1.66
	Average	_	90.07 ± 1.91
Tatar	Wikipedia	159K	86.35 ± 2.17
	News	5.2M	89.75 ± 0.07
	Religion	382K	91.25 ± 2.55
	Average	<u> </u>	89.12 ± 1.60
Kumyk	Wikipedia	_	
	News	286K	91.10 ± 0.86
	Religion	227K	92.47 ± 1.03
	Average	_	91.78 ± 0.94
Precision & recall			

- selected 1000 surface forms at random from RFE/RL corpus, proof read analyses
- Precision (of a form's analyses % correct): 97.32%
- **Recall** (percentage of analyses provided by the transducer that are correct for a form, by comparing against a gold standard): 94.56%