

FINITE-STATE MORPHOLOGICAL TRANSDUCERS FOR THREE KYPCHAK LANGUAGES

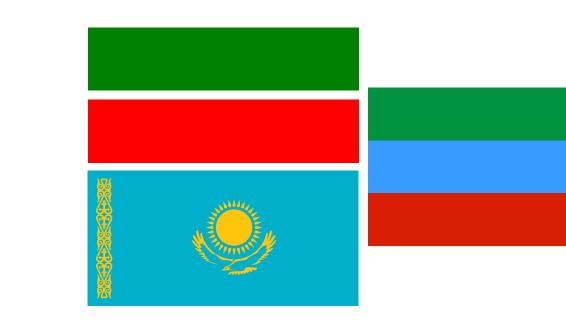
Jonathan North Washington Indiana University

Ilnar Salimzyanov

Francis M. Tyers

Also special thanks to Aida Sundetova

олар бек яхшы экенин гёрген.



jonwashi@indiana.edu

Some University ilnar.salimzyan@gmail.com

UiT Norgga Árktalaš Universitehta francis.tyers@uit.no

email@email



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			
Mongolic Oghuz Persian Russian	moderate — heavy heavy	light light heavy heavy	light moderate heavy heavy	

Morphological transducers

..... Morphological transducers

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

- 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):

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

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

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 ;

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

! "weather"

COMPOUND ; ! "invitation"

Example output

Gloss

Аллагь Оьзю яратгъан затларгъа къарап, thing-s-to look-having, they very good being own-his made God

'God looked at everything he had made and saw that it was very good.'

Kazakh	Tatar	Kumyk
Құдай Өзінің жаратқандарының бәріне қарап, өте жақсы екенін көрді.	Аллаһ Үзе яраткан нәрсәләргә карап, аларның бик яхшы икәнен күрде.	Аллагь Оьзю яратгъан затлагъа къарап, олар бек яхшы экенин гёрге
<pre>Kyдай<n><nom> 03<prn><ref><px3sp><gen></gen></px3sp></ref></prn></nom></n></pre>	Аллаh <n><nom> Y3<prn><ref><px3sp><nom> ярат<v><tv><gpr_past> нәрсә<prn><itg><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><sent> Tagset</sent></p3></past></tv></v></acc></px3sp></ger_past></cop></adj></adv></gen></pl></p3></pers></prn></cm></gna_perf></tv></v></dat></pl></itg></prn></gpr_past></tv></v></nom></px3sp></ref></prn></nom></n>	Аллагь <n><nom> Oьз<prn><ref><px3sp><nom> ярат<v><tv><gpr_past> зат<n><pl><pl><pl><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<n><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom<<pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom<<pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom><pre>sat<nom<<pre>sat<nom><pre>sat<nom><pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<pre>sat<nom<<< td=""></nom<<<></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></nom<<pre></pre></nom></pre></nom></nom<<pre></pre></nom></pre></nom></pre></nom></pre></nom></nom<<pre></pre></nom></pre></nom></pre></nom></pre></nom></nom<<pre></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></nom></pre></n></pre></n></pre></n></pre></n></pre></n></pre></n></pre></n></pre></n></pre></n></pre></n></pre></n></pre></n></pre></n></pre></pl></pl></pl></n></gpr_past></tv></v></nom></px3sp></ref></prn></nom></n>

			Tagset		
<n></n>	Noun	<p2></p2>	Second person	<px3sg></px3sg>	3rd person poss. (Singular)
<np></np>	Proper noun	<p3></p3>	Third person	<px3pl></px3pl>	3rd person poss. (Plural)
<v></v>	Verb	<ant></ant>	Anthroponym	<neg></neg>	Negative
<det></det>	Determiner	<dem></dem>	Demonstrative	<aor></aor>	Aorist
<cnjcoo></cnjcoo>	Coord. conjunct.	<m></m>	Masculine	<imp></imp>	Imperative
<cnjadv></cnjadv>	Adv. conjunct.	<sg></sg>	Singular	<gna_perf< td=""><td>>Verbal adverb (Perfect)</td></gna_perf<>	>Verbal adverb (Perfect)
<adv></adv>	Adverb	<pl><</pl>	Plural	<pre><pre_impf< pre=""></pre_impf<></pre>	>Participle (Imperfect)
<vaux></vaux>	Auxiliary verb	<nom></nom>	'Nominative'	<prc_irre< td=""><td>>Participle (Irrealis)</td></prc_irre<>	>Participle (Irrealis)
<cop></cop>	Copula	<gen></gen>	Genitive	<pre><preal< pre=""></preal<></pre>	>Participle (Realis)
<iv></iv>	Intransitive	<acc></acc>	Accusative	<cm></cm>	Comma
<tv></tv>	Transitive	<loc></loc>	Locative		

Morphophonology

Desonorisation .

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

	TTUTTIO	 		
Part of speech		Number of stems		
	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				
Encyclop kaz	wpdum	ıp	20131006	
tat	wpdum		20130225	
kıın	n —	_		

Number of stems...

Encyclop	kaz tat kum	wpdump wpdump —	20131006 20130225 —	
News	kaz tat kum	RFE/RL Татар-информ Ёлдаш	azattyq.org 2010 tat.tatar-inform.ru yoldash.etnosmi.r	
Religion	kaz tat kum	quran + bible quran + nt genesis + nt	kkitap.net, kuran.l ibt.org.ru, tanzil.n 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
IXaZaKII	Religion	851K	92.49 ± 1.66
	Average	_	90.07 ± 1.91
	Wikipedia	159K	86.35 ± 2.17
Tatar	News	5.2M	89.75 ± 0.07
Talai	Religion	382K	91.25 ± 2.55
	Average	_	89.12 ± 1.60
	Wikipedia	_	
Kumyk	News	286K	91.10 ± 0.86
ixumyk	Religion	227K	92.47 ± 1.03
	Average	_	91.78 ± 0.94

 selected 1000 surface forms at random from RFE/RL corpus, proof read analyses

Precision & recall......

- **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