

FINITE-STATE MORPHOLOGICAL TRANSDUCERS FOR THREE KYPCHAK LANGUAGES

Jonathan North Washington Ilnar Salimzyanov

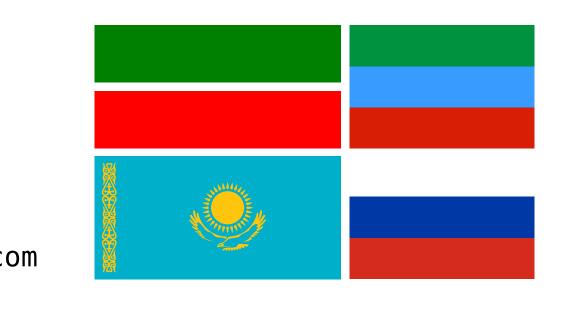
Казан (Идел буе) федераль университеты ilnar.salimzyan@gmail.com Indiana University jonwashi@indiana.edu

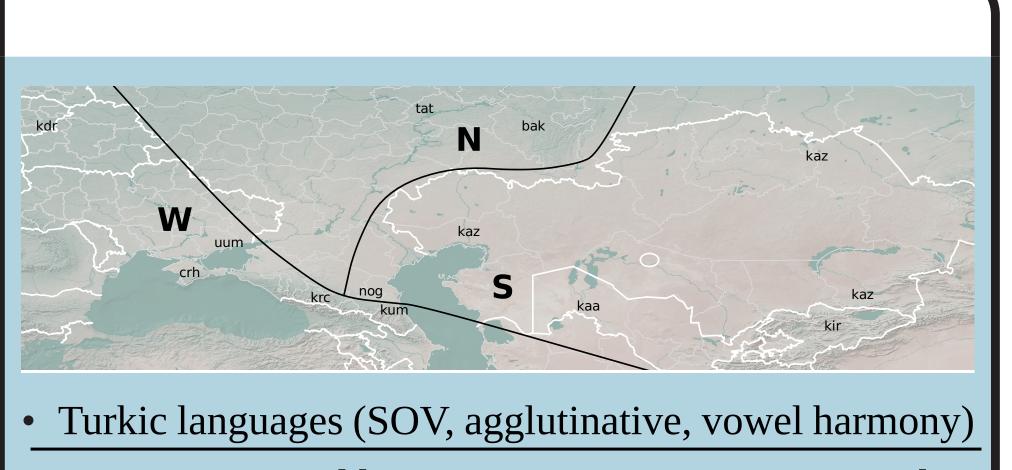
God

Francis M. Tyers

Aida Sundetova UiT Norgga Árktalaš Universitehta sun27aida@gmail.com francis tyers@uit.no

Special thanks to





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

Morphological transducers

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

heavy

heavy

<tv>

<p3>

<pl>

• 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)

heavy

Russian

- 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.....

.... Morphological & orthographical words

- өнүктүрөбүзбү? 'will we develop [it]?' ӨНҮК<v><tv><caus><aor><pl>><pl>+бы<qst>
- келатсаң 'if you come'
- кел<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) І н
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) M A	-(I)ңA	-(S) І н А

- 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 ;

! "weather"

чакыруу% кагазы:чакыруу% кагаз N-INFL-3PX-COMPOUND ; ! "invitation"

Gloss. Аллагь Оьзю

яратгъан затларгъа къарап, олар бек яхшы экенин гёрген. own-his made

<ifi>

<prn>

<qnt>

thing-s-to look-having, they very good being saw.

'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 гёрген.	
Құдай <n><nom> 03<prn><ref><px3sp><gen> жарат<v><tv><ger_past><pl><px3sp><gen> бәрі<prn><qnt><px3sp><dat> қара<v><tv><gna_perf> ,<cm> — 0те<adv> жақсы<adj> e<cop><ger_past><px3sp><acc> көр<v><tv><ifi><p3><sent> .<sent></sent></sent></p3></ifi></tv></v></acc></px3sp></ger_past></cop></adj></adv></cm></gna_perf></tv></v></dat></px3sp></qnt></prn></gen></px3sp></pl></ger_past></tv></v></gen></px3sp></ref></prn></nom></n>		Yз <pr>><gen> ярат</gen></pr>	Aллah <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></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><qdat> къара<v><tv><gna_perf> ,<cm> олар<prn><pers><p3><ql><nom> бек<adv> яхшы<adj> э<cop><ger_past><px3sp><acc> гёр<v><tv><past><p3><sent> </sent></p3></past></tv></v></acc></px3sp></ger_past></cop></adj></adv></nom></ql></p3></pers></prn></cm></gna_perf></tv></v></qdat></pl></n></gpr_past></tv></v></nom></px3sp></ref></prn></nom></n>	
			Tagset			
<n></n>	Noun	<nom></nom>	'Nominative'	<itg> Int</itg>	terrogative	
< V >	Verb	<gen></gen>	Genitive	<pers></pers>		
<det></det>	Determiner	<acc></acc>	Accusative	<pre><ger_past> Ve</ger_past></pre>	rbal noun (Past)	
<adj></adj>	Adjective	<px3sp></px3sp>	3rd person poss.	<pre><gna_perf> Ve</gna_perf></pre>	rbal adverb (Perfect)	
<adv></adv>	Adverb		(Singular/Plural)	<pre><gpr_past> Ve</gpr_past></pre>	rbal adjective (Past)	
<iv></iv>	Intransitive	<past></past>	Past (General)	<cm> Co</cm>	mma	

Pronoun

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

Transitive

Plural

Third person

- $\{L\}$ desonorises to π 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 | %>:] ;

- Part of Apertium Turkic project:
- http://wiki.apertium.org/wiki/Apertium_Turkic
- Transducers available live at turkic.apertium.org
- Source code available from apertium's svn repo
- Turkic RBMT mailing list (>25 subscribers): apertium-turkic@lists.sourceforge.net Feel free to post in any language!
- See our paper in the LREC 2014 proceedings
 - And feel free to contact the authors any time!

	Number (of stems	<u> </u>	
Dart 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	

Sentence

Past (Eyewitness/Recent) <sent>

_				
	type	lang	contents	origin
_	Encyclop	kaz tat kum	wpdump wpdump —	20131006 20130225 —
-	News	kaz tat kum	RFE/RL Татар-информ Ёлдаш	azattyq.org 2010 tat.tatar-inform.ru 2005-2011 yoldash.etnosmi.ru
	Religion	kaz tat kum	quran + bible quran + nt genesis + nt	kkitap.net, kuran.kz ibt.org.ru, tanzil.net ibt.org.ru
• split into 10 equal parts: coverage calc			rage calculated over each	

. Test corpora

- split into 10 equal parts; coverage calculated over each separately; standard deviation of mean calculated
 - Coverage measures
- 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 analysed corpus

.........Coverage results (as of r36739)...... Coverage (%) Corpus **Tokens** Language Wikipedia 25.6M 85.61 ± 1.37 3.8M 92.12 ± 2.72 News Kazakh Religion 851K 92.49 ± 1.66 90.07 ± 1.91 Average 86.35 ± 2.17 Wikipedia 159K 5.2M 89.75 ± 0.07 News Tatar Religion 382K 91.25 ± 2.55 89.12 ± 1.60 Average Wikipedia 286K 91.10 ± 0.86 News Kumyk 227K 92.47 ± 1.03 Religion 91.78 ± 0.94 Average

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%