

FINITE-STATE MORPHOLOGICAL TRANSDUCERS FOR THREE KYPCHAK LANGUAGES

Jonathan North Washington Ilnar Salimzyanov

Казан (Идел буе) федераль университеты ilnar.salimzyan@gmail.com

Francis M. Tyers

Special thanks to UiT Norgga Árktalaš Universitehta Aida Sundetova A. Sultanmuradov francis.tyers@uit.no



Kypchak languages



Indiana University

jonwashi@indiana.edu

Turkic languages (SOV, agglutinative, vowel harmony)

classif'tion	Kazakh	Tatar	Kumyk
	/qazaq/	/tɒtɑɾ/	/qumuq/
	S Kypchak	N Kypchak	W Kypchak
population of	f speakers		
number	8M-12M	5.4M	430K
primary	Kazakhstan	Tatarstan	Dagestan
secondary	China, Mongolia	Bashqortostan	—
external influ	iences		
Mongolic	moderate	light	light
Oghuz	—	light	moderate
Persian	heavy	heavy	heavy

heavy

heavy

Morphological transducers

heavy

Russian

Morphological transducers

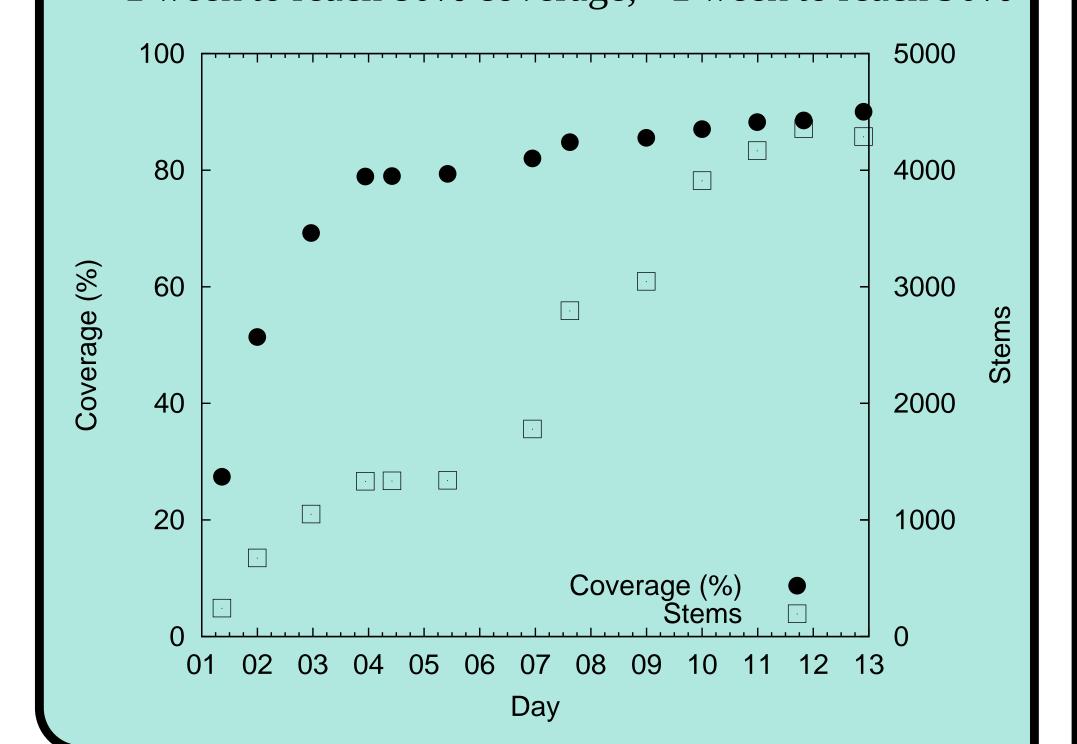
- Take a surface form, and produce valid lexical form(s)
- Take a lexical form, and produce valid surface form(s) 'алдым' ↔ ал<v><tv><ifi><pl><sg>, алд<n><px1sg><nom> Transducers for Turkic languages.....
- Turkish (Çöltekin, 2010 & 2014; Öflazer, 1994)
- Crimean Tatar (Altıntaş, 2001)
- Turkmen (Tantuğ et al., 2006)
- Kyrgyz (Washington et al., 2012)
- Kazakh (Бекманова & Махимов, 2013)
- our Kazakh, Tatar, Kumyk: all GPL (=free and open)! Framework: HFST.....

Reimplements Xerox FST formalisms (lexc & twol)

 Also provides a wrapper around popular free/opensource FST toolkits: SFST, OpenFST, and Foma

...... Development effort...... Kumyk transducer based on Kazakh, Tatar transducers

• \sim 1 week to reach 80% coverage, +1 week to reach 90%



Categorisation

- Other Turkic transducers: 0-derivation (overgenerates)
- Our approach: categorization (e.g., adjectives, below)

our approach categorization (e.g., aajeetives, serow)					
Type	Gloss	<adj>(<comp>)</comp></adj>	<adj>(<comp>)<subst></subst></comp></adj>	<adj>(<comp>)<advl></advl></comp></adj>	
A1	'good'	яхшы (яхшырак)	яхшы (яхшырак)	яхшы (яхшырак)	
A2	ʻold'	иске (искерэк)	иске (искерэк)		
A3	'dead'	үле (—)	үле (—)		
A4	'basic'	төп (—)	— (—)	— (—)	

Further information

- 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!

Example output

God

Kazakh (kaz)

Құдай<n><nom>

өте<adv>

<sent>

жақсы<adj>

Құдай Өзінің жаратқандарының

 θ 3prn><ref><px3sp><gen>

қара<v><tv><qna perf>

6əpiprn><qnt><px3sp><dat>

e<cop><ger past><px3sp><acc>

көр<v><tv><ifi><p3><sg>

бәріне қарап, өте жақсы екенін көрді.

mapar<v><tv><ger past><pl><px3sp><gen>

Gloss. Құдай Өзінің жаратқандарының бәріне Аллаһ Үзе яраткан

нәрсәләргә Аллагь Оьзю яратгъан затлагъа own-his created

Tatar (tat)

карап, къарап, олар [everything/thing-s]-to looked.at, they/their very good

қарап,

аларның бик яхшы икәнен күрде. бек яхшы экенин гёрген. being saw.

өте жақсы екенін көрді.

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

Kumyk (kum) Аллагь Оьзю яратгъан затлагъа Аллаh Үзе яраткан нәрсәләргә карап, аларның бик яхшы икәнен күрде. къарап, олар бек яхшы экенин гёрген.

Аллаh<n><nom> Y3<prn><ref><px3sp><nom> ярат<v><tv><gpr past> нәрсә<n><pl><dat> kapa<v><tv><qna perf>

аларprn><pers><p3><pl><gen> бик<adv> яхшы<adj>

ярат<v><tv><gpr past> зат<n><pl><dat> къapa<v><tv><qna perf> oлapconsprint of the original origin бек<adv> яхшы<adj> ><cop><ger past><px3sp><acc>

Oьз<prn><ref><px3sp><nom>

Аллагь<n><nom>

и<cop><ger past><px3sp><acc> κγp<v><tv><past><p3><sg> rëp<v><tv><past><p3><sg> <sent> .<sent>

			T	agset.			
<n></n>	Noun		Third person		Personal		3rd person poss.
<v></v>	Verb	<pl><</pl>	Plural	<cm></cm>	Comma		(Singular/Plural)
prn>	Pronoun	<nom></nom>	'Nominative'	<sent></sent>	Sentence	<gna_perf></gna_perf>	Verbal adverb
<det></det>	Determiner	<gen></gen>	Genitive	<past></past>	Past (General)		(Perfect)
adj>	Adjective	<acc></acc>	Accusative	<ifi></ifi>	Past	<pre><gpr_past></gpr_past></pre>	Verbal adjective
<adv></adv>	Adverb	<dat></dat>	Dative		(Eyewitness/Recent		(Past)
<iv></iv>	Intransitive	<qnt></qnt>	Quantifier			<ger_past></ger_past>	Verbal noun (Past)
<tv></tv>	Transitive	<ref></ref>	Reflexive				

Orthography-phonology mapping issues

Have front- and back-vowel readings in native words

	letters	values	examples
ka	z и, у, ю	/wej, aw, jaw/ /wej, aw, jaw/	қиюд <mark>а</mark> 'chopping down' киюд <mark>е</mark> 'getting dressed'
ta	t e	э / С _ /j/+ы /j/+э	дәресләр 'lessons' еллар 'years' егетләр 'boys'
ku	m ё,ю	/ø, y/ / C _ /jø, jy/ /jo, ju/	гюнлер 'days' гёзлер 'eyes' юреклер 'hearts' ёнкюлер 'darlings' 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

Letters that represent front vowels in native words may represent back vowels in Russian words

	native word	borrowed word
kaz tat kum	елдің 'country's' галимнәр 'scientists' сёзлер 'words'	Назарбаевтың 'Nazarbayev's' артистлар 'artists' самолётлар 'airplanes'

solution: separate continuation lexicon (messy rules)Acronyms and numerals.....

бестен handled twol by отыздан, handled by 30-дан, 5-тен twol not

5{э}{c}>-{D}{A}н

Resulting messy twol rules RdYotVow = ёюЁЮ;

AbstractVow = $%{a%}^{'}%{3}% %{y} %{0}% ;$
"A front unrounded harmony"
%{A%}:e <=> [[:FrontVow [:Vow :ь]] :Cns :Cns*]/:0 _ ;
[[:RdVow :ь] :Cns :Cns*]/:0 _ ;
[[[\[.#. :Vow]] :RdYotVow] :Cns :Cns*]/:0 _ ;
[:RdYotVow й:0 :RdYotVow :Cns :Cns*]/[:0 - й:0] _ ;
[[%{э%}:0 %{γ%}:0] :Cns*]/[[:0 - AbstractVow:] %-:]* _ ;
except
[:RdYotVow :Cns* %{¾%}:0 :Cns*]/[:0 - %{¾%}:0] _ ;
[:Cns :p %{¾%}: %>: :Cns*]/:0 _ ;
[[:Vow - :RdYotVow] :RdYotVow :Cns :Cns*]/:0 _ ;
[:Vow]/[[:0 - й:0] %>:] _ ;

Evaluation

................Number of stems.

Part of speech	Number of stems				
r art or specen	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

	Wikipedia	News	Religion
kaz tat kum	Wikipedia Wikipedia —	azattyk.org tat.tatar-inform.ru yoldash.etnosmi.ru	Quran + Bible Quran + New Testament 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	Language Corpus Tokens Coverage (%)				
Kazakh	Wikipedia News Religion	25.6M 3.8M 851K	85.61 ± 1.37 92.12 ± 2.72 92.49 ± 1.66	2.43 2.88 2.63	
(r50547)	Average		90.07 ± 1.91	2.64	
Tatar	Wikipedia News Religion	159K 5.2M 382K	86.35 ± 2.17 89.75 ± 0.07 91.25 ± 2.55	2.24 2.30 2.24	
(r50260)	Average		89.12 ± 1.60	2.26	
Kumyk	News Religion	286K 227K	91.10 ± 0.86 92.47 ± 1.03	1.53 1.53	
(r50300)	Average		91.78 ± 0.94	1.53	

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

Future Work

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