

# FINITE-STATE MORPHOLOGICAL TRANSDUCERS FOR THREE KYPCHAK LANGUAGES

Jonathan North Washington Ilnar Salimzyanov Indiana University

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



jonwashi@indiana.edu

Turkic languages (SOV, agglutinative, vowel harmony)

Turkie languages (50 v, aggratinative, vower narmony)			
	Kazakh	Tatar	Kumyk
	/qazaq/	/totar/	/qumuq/
classification	Southern	Northern	Western
population of speakers			
number	8M-12M	5.4M	430K
primary	Kazakhstan	Tatarstan	Dagestan
secondary	China, Mongolia	Bashqortostan	_
external influences			
Mongolic	moderate	light	light
Oghuz	<u>—</u>	light	moderate
Persian	heavy	heavy	heavy
Russian	heavy	heavy	heavy

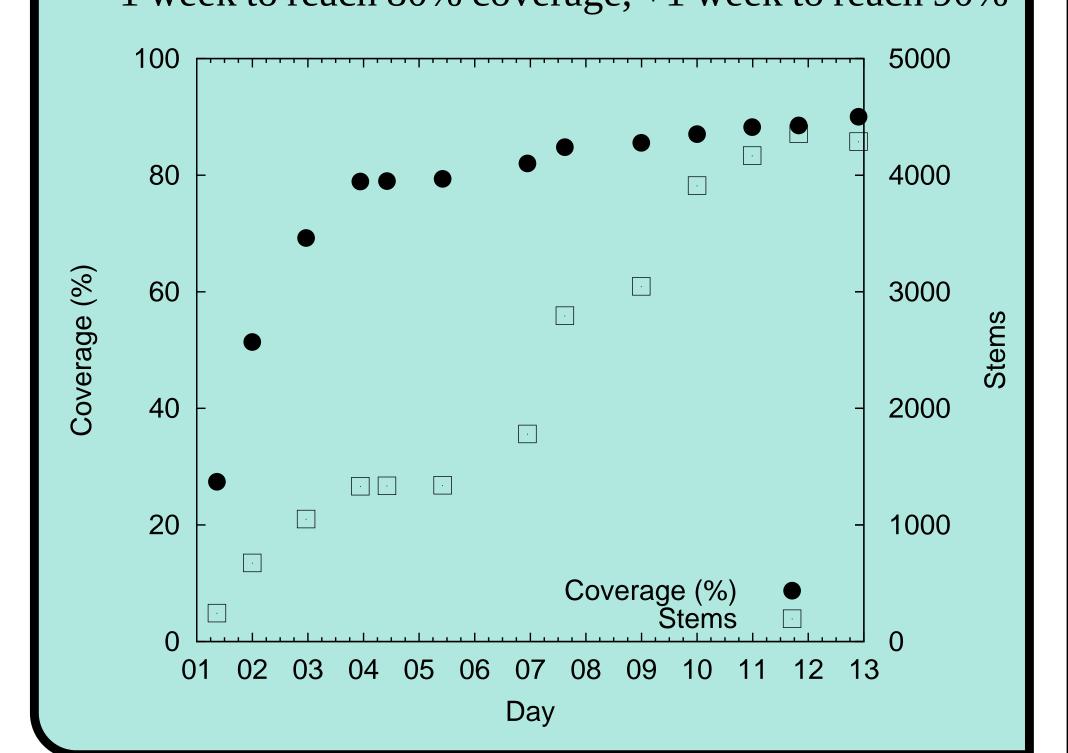
# Morphological transducers

# ...... 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><pl><sg>, алд<n><pxlsg><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

#### Kumyk transducer based on Kazakh, Tatar transducers

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



### Categorisation

'basic'

- Other Turkic transducers: 0-derivation (overgenerates)
- Our approach: categorisation (e.g., adjectives, below) Type Gloss <adj>(<comp>)<subst> <adj>(<comp>)<advl> <adj>(<comp>) яхшы (яхшырак) яхшы (яхшырак) 'good яхшы (яхшырак) иске (искерэк) иске (искерэк) 'dead' үле (—) үле (—)

— (—)

### Further information

төп (—)

- 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)

<sent>

<tv>

Құдай Өзінің жаратқандарының бәріне Аллаһ Үзе яраткан Аллагь Оьзю яратгъан

нәрсәләргә затлагъа

қарап, карап, къарап,

аларның бик яхшы икәнен күрде. олар [everything/thing-s]-to looked.at, they/their very good

Kumyk (kum)

бек яхшы экенин гёрген. being saw.

өте жақсы екенін көрді.

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

#### Tatar (tat)

Gloss.

Құдай Өзінің жаратқандарының Аллаh Үзе яраткан нәрсәләргә карап, аларның бик яхшы икәнен күрде. бәріне қарап, өте жақсы екенін көрді. Құдай<n><nom> Аллаh<n><nom>  $\theta$ 3prn><ref><px3sp><gen> Y3<prn><ref><px3sp><nom> жapaт<v><tv><ger past><pl><px3sp><gen> ярат<v><tv><gpr past> нәрсә<n><pl><dat>

<sent>

6əpiprn><qnt><px3sp><dat> қара<v><tv><gna perf> өте<adv> жақсы<adj>

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

көр<v><tv><ifi><p3><sg>

own-his created

kapa<v><tv><qna perf> аларprn><pers><p3><pl><gen> бик<adv> яхшы<adj> и<cop><ger past><px3sp><acc> κγp<v><tv><past><p3><sg>

къарап, олар бек яхшы экенин гёрген. Аллагь<n><nom> 0ьз<prn><ref><px3sp><nom> ярат<v><tv><gpr past> зат<n><pl><dat> къapa<v><tv><qna perf> олаponapon бек<adv>

Аллагь Оьзю яратгъан затлагъа

яхшы<adj> э<cop><ger past><px3sp><acc> rëp<v><tv><past><p3><sg> <sent>

Tagset ..... <per>> Personal Noun Third person 3rd person poss. <px3sp> Verb Plural (Singular/Plural) Comma <gna\_perf> Verbal adverb <sent> Sentence 'Nominative' Pronoun Determiner Genitive <past> Past (General) (Perfect) <gen> Adjective <gpr past> Verbal adjective Past Accusative <ifi> Adverb Dative (Eyewitness/Recent) (Past) <adv> <dat> Quantifier <ger past> Verbal noun (Past) Intransitive Reflexive Transitive <ref>

# Orthography-phonology mapping issues

Have front- and back-vowel readings in native words

		letters	values	examples
	kaz	и, у, ю	/wej, we, jew/ /wej, we, jev/	қиюд <mark>а</mark> 'chopping down' киюд <mark>е</mark> 'getting dressed'
_	tat	e	э / С _ /j/+ы /j/+э	дәресләр 'lessons' еллар 'years' егетләр 'boys'
	kum	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 cover majority of examples
- unaccounted-for words marked harmony-forcing char
- adjust rules for harmony-forcing characters ..... Loanwords .......

## Letters that represent front vowels in native words may

represent back vowels in Russian words native word example Russian word example елдің 'country's' Назарбаевтың 'Nazarbayev's' галимнәр 'scientists'

артистлар 'artists' самолётлар 'airplanes' сёзлер 'words' solution: separate continuation lexicon (messy rules)

LEXICON N1-RUS :%{\angle \cdots\} N1 ;

LEXICON Nouns

артист:apтист N1-RUS ; ! "artist" галим:галим N1 ; ! "scientist"

- twol rules handle phonology for spelt-out words отыздан 'from thirty', бестен 'from five'
- no phonological triggers available in numerals (etc.) 30-дан 'from 30', 5-тен 'from 5'
- solution: phonology-triggering characters

4:4%{3%}%{c%} NUM-DIGIT; ! "τθρτ" 5:5%{3%}%{c%} NUM-DIGIT; ! "6ec" 3%0:3%0%{a%}%{3%} NUM-DIGIT ; ! "отыз"

..... A resulting messy twol rule......

RdYotVow =  $\ddot{e}$  ω  $\ddot{E}$  Ю ; AbstractVow =  $%{a%}$  % ${y%}$  % ${y%}$  % ${o%}$  ; "A front unrounded harmony" [ [: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 ] \_ ;  $[[%{3\%}:0|%{\gamma\%}:0] : Cns*]/[[:0 - AbstractVow:] | %-:]* _ ;$ :RdYotVow :Cns\* %{¾%}:0 :Cns\* ]/[ :0 - %{¾%}:0 ] \_ ; [ :Cns :p %{,2%}: %>: :Cns\* ]/:0 \_ ; [ [ :Vow - :RdYotVow ] :RdYotVow :Cns :Cns\* ]/:0 \_ ; [:Vow]/[[:0 - й:0]|%>:]\_;

# Evaluation

.....Number of stems.

Part of speech	Number of stems		
r are 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	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  $\geq 1$  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......

	Corpus	Tokens	Coverage (%)	Amb.
Kazakh	Wikipedia News Religion	25.6M 3.8M 851K	$85.61 \pm 1.37$ $92.12 \pm 2.72$ $92.49 \pm 1.66$	2.43 2.88 2.63
(r50547)	Average		$90.07 \pm 1.91$	2.64
Tatar	Wikipedia News Religion	159K 5.2M 382K	$86.35 \pm 2.17 89.75 \pm 0.07 91.25 \pm 2.55$	2.24 2.30 2.24
(r50260)	Average		$89.12 \pm 1.60$	2.26
Kumyk	News Religion	286K 227K	$91.10 \pm 0.86$ $92.47 \pm 1.03$	1.53 1.53
(r50300)	Average		$91.78 \pm 0.94$	1.53
• solocted &	proofed unique	random sur	rface forms from no	MC COrpor

selected & proofed unique random surface forms from news corpora Language Forms Drecision (%)  $\mathbf{p}_{\mathbf{q}}$ 

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

### Ongoing and future work

- Disambiguation, more stems
- Machine translation between these languages
- Other Turkic lgs.: Nogay, Bashqort, Uzbek, Chuvash