



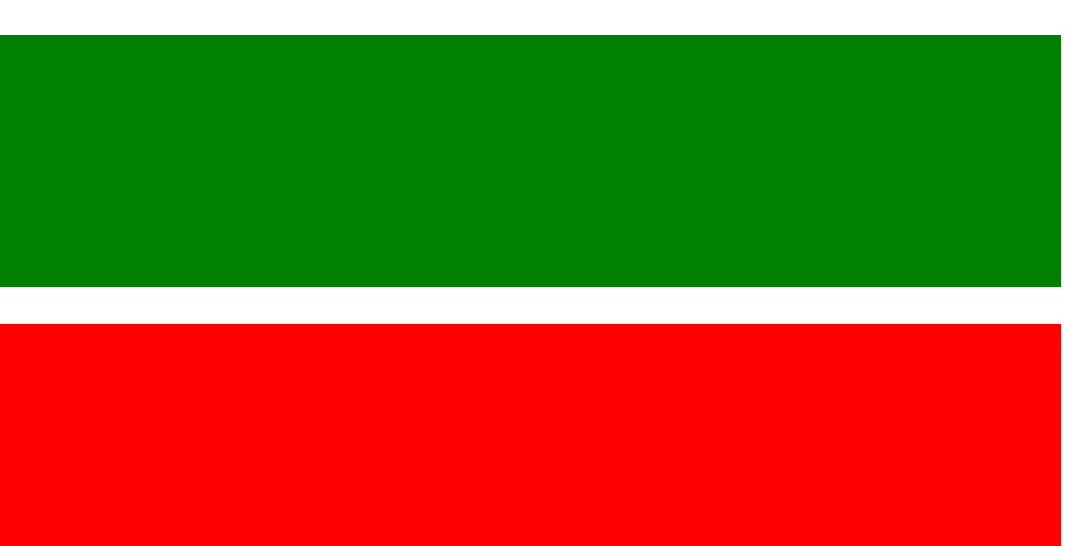
DESIGNING FINITE-STATE MORPHOLOGICAL TRANSDUCERS FOR KYPCHAK LANGUAGES

Jonathan North Washington
Indiana University
jonwashi@indiana.edu

Ilmar Salimzyanov
Казан (Идел буе) федераль университеты
ilmar.salimzyanov@gmail.com

Francis M. Tyers
UIT Norgga Árkalaš Universitehta
francis.tyers@uit.no

Special thanks to:
Tolgonay Kubatova
Aida Sundetova
Ağarahim Sultanmuradov



Kypchak languages



- Turkic languages (SOV, agglutinative, vowel harmony)

	Kyrgyz	Kazakh	Tatar	Kumyk
	/qɯɾɣɯz/	/qɑzɑq/	/tɑtɑr/	/qumɯq/
classification	Eastern	Southern	Northern	Western
population of speakers				
number	3M	8M-12M	5.4M	430K
primary	Kyrgyzstan	Kazakhstan	Tatarstan	Dagestan
secondary	China, etc.	China, Mongolia	Bashkortostan	—
external influences				
Mongolic	moderate	moderate	light	light
Oghuz	—	—	light	moderate
Persian	heavy	heavy	heavy	heavy
Russian	heavy	heavy	heavy	heavy

Morphological transducers

- Efficient (in speed & size) models of a language's morphology

- 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 & 2014; Oflazer, 1994)
- Crimean Tatar (Altıntaş, 2001)
- Turkmen (Tantuğ et al., 2006)
- Kazakh (Бекманова & Махимов, 2013)
- our Kyrgyz, Kazakh, Tatar, Kumyk: all GPL (=free and open)!

Framework: HFST

- Reimplements Xerox FST formalisms (lexc & twol)
- Also provides a wrapper around popular free/open-source FST toolkits: SFST, OpenFST, and Foma

Approach

- morphotactics implemented in lexc
- morphophonology implemented in twol
- compiled separately; compose-intersected to single transducer

алдым ↔ ал<D><I>>м ↔ ал<v><tv><ifi><p1><sg>

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
apertium-turkic@lists.sourceforge.net
- Feel free to post in any language!
- See our papers in LREC proceedings
(2012: Kyrgyz, 2014: Kazakh, Tatar, Kumyk)
- And feel free to contact the authors any time!

Example output

	Gloss
(1) Кудай Өзү жаратканынын баарына карап, өтө жакшы экенин көрдү.	
Кудай Өзінің жаратқандарының бәріне карап, өте жақсы екенін көрді.	
Аллаһ Үзе яратқан нәрсәләргә карап, аларның бик яхшы икәнән күрдө.	
Аллаһь Озью яратгъан затлагъа къарап, олар бек яхшы экенин гёрген.	
God own-his created [everything/thing-s]-to looked.at, they/their very good being saw.	
‘God looked at everything he had created and saw that it was very good.’ (Bible, Genesis 1:31)	

Output			
Kyrgyz (kir)	Kazakh (kaz)	Tatar (tat)	Kumyk (kum)
Кудай Өзү жаратканынын баарына карап, өтө жакшы экенин көрдү.	Кудай Өзінің жаратқандарының бәріне карап, өте жақсы екенін көрді.	Аллаһ Үзе яратқан нәрсәләргә карап, аларның бик яхшы икәнән күрдө.	Аллаһь Өзью яратгъан затлагъа къарап, олар бек яхшы экенин гёрген.
Кудай<n><nom> Өз<prn><ref><px3sp><nom> жарат<v><tv><ger_past><px3sp><gen> баары<prn><qnt><px3sp><dat> кара<v><tv><gna_perf> ,<cm> өтө<adv> жакшы<adj> э<cop><ger_past><px3sp><acc> көр<v><tv><ifi><p3><sg> ,<sent>	Кудай<n><nom> Өз<prn><ref><px3sp><gen> жарат<v><tv><ger_past><pl><px3sp><gen> бәрі<prn><qnt><px3sp><dat> қара<v><tv><gna_perf> ,<cm> өте<adv> жақсы<adj> е<cop><ger_past><px3sp><acc> көр<v><tv><ifi><p3><sg> ,<sent>	Аллаһ<n><nom> Үз<prn><ref><px3sp><nom> ярат<v><tv><gpr_past> нәрсә<n><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><sg> ,<sent>	Аллаһь<n><nom> Өз<prn><ref><px3sp><nom> ярат<v><tv><gpr_past> зат<n><pl><dat> къара<v><tv><gna_perf> ,<cm> олар<prn><pers><p3><pl><nom> бек<adv> яхшы<adj> э<cop><ger_past><px3sp><acc> гёр<v><tv><past><p3><sg> ,<sent>

	Tagset
<n>	Noun
<v>	Verb
<prn>	Pronoun
<det>	Determiner
<adj>	Adjective
<adv>	Adverb
<iv>	Intransitive
<tv>	Transitive
<p3>	Third person
<pl>	Plural
<ref>	Reflexive
<pers>	Personal
<nom>	‘Nominative’
<sent>	Sentence
<gen>	Genitive
<acc>	Accusative
<dat>	Dative (Eyewitness/Recent)
<qnt>	Quantifier
<cm>	Comma
<gna_perf>	Verbal adverb (Perfect)
<past>	Past (General)
<ifi>	Past (Past)
<gpr_past>	Verbal adjective (Past)
<px3sp>	3rd person poss. (Singular/Plural)
<ger_past>	Verbal noun (Past)

Morphophonology

Desonorisation (kaz & kir)

- {N} desonorises to д after a consonant
алма-{N}{I} → алманы ‘apple-ACC’
сыр-{N}{I} → сырды ‘secret-ACC’
- {L} desonorises to д after cons. of sonority ≤ /l/
сыр-{L}{A}p → сырлар ‘secret-PL’
кыз-{L}{A}p → кыздар ‘girl-PL’

“L Desonorisation”
%{L%}:д <=> :VoicedLowSonCns %>: __ ;
“N Desonorisation”
%{N%}:д <=> :VoicedCns %>: __ ;

Epenthesis

- Turn {y} into a harmonised high vowel when a vowel doesn’t follow the following consonant:
мур{y}н → мурун ‘nose’
мур{y}н+{I}м → мурдум ‘my nose’

%{y%}:Vy <=> [:LastVowel :Cns* :Cns]/[:0] __ [:Cns [:#. :Cns]]/[:0 %>:] ; where Vy in (и у и и у у у у у у) LastVowel in (и у е э я а ё о у ю) matched ;
--

Morphotactics

Morphological & orthographical words

- өнүктүрөбүзбү ? ‘will we develop [it]?’
өнүк<v><tv><caus><aor><p1><pl>+бы<qst>
- келатсаң ‘if you come’
кел<v><iv><prc_impf>+жат<vaux><gna_cnd><p2><sg>
- Irregular [noun + possessive + case] forms
- Some combinations of possessive + case morphemes are unpredicted (i.e., not formed simply by concatenation and application of phonology):

case	form	1SG	2SG	3SP
nominative	—	-(I)м	-(I)н	-(c)I
accusative	-NI	-(I)мдI	-(I)ндI	-(c)Iн
genitive	-NIн	-(I)мдIн	-(I)ндIн	-(c)IндA
locative	-DA	-(I)мдAн	-(I)ндAн	-(c)IндAн
ablative	-DAn	-(I)мдAn	-(I)ндAn	-(c)IндAn
dative	-GA	-(I)мAн	-(I)нAн	-(c)IнAн

A,I,N,D,G have various allophones; (I) null after vowels; (c) null after cons.

- underlying <px3sp> form used: {s}{I}{n}
- phonological rules delete {s}, {n} by context
- morphophonology more complicated, morphotactics simpler

Noun-noun compounds

- a N-N compound type: N2 has <px3> and related morphology
e.g., аба ырайы<n><loc>: аба ырайында, *аба ырайыда

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

Orthography-phonology mapping issues

Ambiguous characters

- Have front- and back-vowel readings in native words

	letters	values	examples
kaz	и, у, ю	/əj, əw, jəw/ /əj, əw, jəw/	киюда ‘chopping down’ киюде ‘getting dressed’
tat	е	э / C _ /j/+ы /j/+э	дәресләр ‘lessons’ еллар ‘years’ егетләр ‘boys’
kum	ё, ю	/jə, jy/ /jo, ju/	гюнләр ‘days’ гезләр ‘eyes’ юрекләр ‘hearts’ ёнколлар ‘darlings’ юлдузлар ‘stars’ ёллар ‘roads’

- solution: hairy twol rules cover majority of examples
- unaccounted-for words get a harmony-forcing character
- 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
kaz	елдің ‘country’s’	Назарбаевтың ‘Nazarbayev’s’
tat	галимнәр ‘scientists’	артистлар ‘artists’
kum	сөзләр ‘words’	самолётлар ‘airplanes’

- solution: separate continuation lexicon (messy rules)

LEXICON N1-RUS
:%{a%} N1 ;
LEXICON Nouns
артист:артист N1-RUS ; ! "artist"
галим:галим N1 ; ! "scientist"

Acronyms and numerals

- twol rules handle phonology for spelt-out words
отыздан ‘from thirty’, бесден ‘from five’
- no phonological triggers available in numerals
- (incorrect phonological triggers in acronyms)
30-дан ‘from 30’, 5-тен ‘from 5’
- solution: phonology-triggering characters
- simplified: e.g., {c} for all voiceless ostruents

4:4%{ə%}%{c%} NUM-DIGIT ; ! "төрт"
5:5%{ə%}%{c%} NUM-DIGIT ; ! "бес"
3%0:3%0%{a%}%{ə%} NUM-DIGIT ; ! "отыз"

A resulting messy twol rule

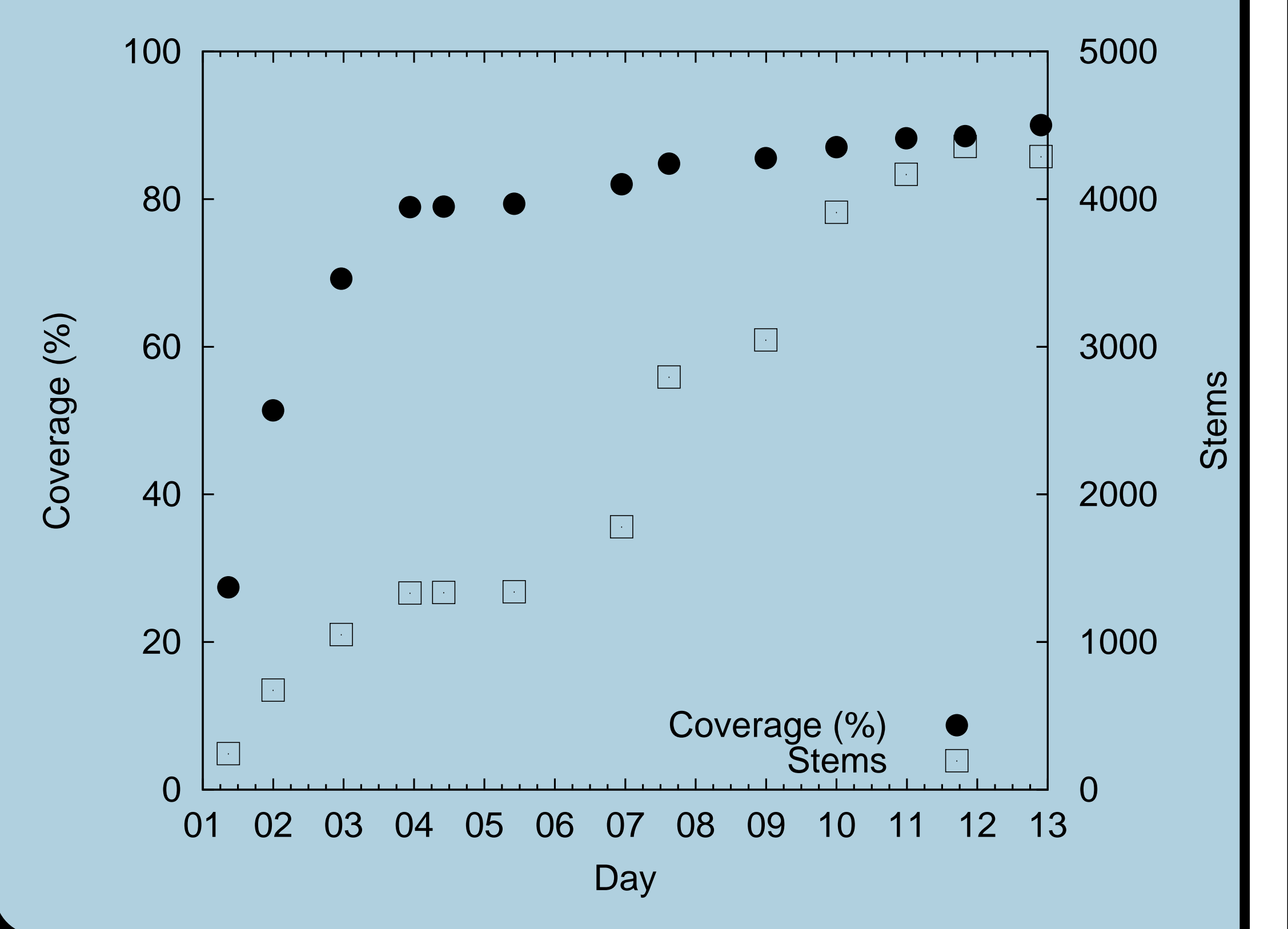
RdYotVow = ё ю ё ю ; AbstractVow = %{a%} %{ə%} %{y%} %{o%} ; “A front unrounded harmony” %{A%}:e <=> [[:FrontVow [[:Vow :ь]] :Cns :Cns*]/:0 _ ; [[:RdYotVow :ь] :Cns :Cns*]/:0 _ ; [[[\ [:#. :Vow]] :RdYotVow] :Cns :Cns*]/:0 _ ; [[:RdYotVow й:0 :RdYotVow :Cns :Cns*]/[[:0 - й:0] _ ; [[%{ə%}:0] [%{y%}:0] :Cns*]/[[[:0 - AbstractVow:]] %:-]* _ ; except [[:RdYotVow :Cns* %[%{a%}:0 :Cns*]/[[:0 - %[%{a%}:0] _ ; [[:Cns :p %[%{a%}:0 %>: :Cns*]/:0 _ ; [[:Vow - :RdYotVow] :RdYotVow :Cns :Cns*]/:0 _ ; [[:Vow]/[[[:0 - й:0]] %>:] _ ;
--

Related issues

- Stem-final й combines with following а, о, у to form я, ё, ю
-

Development effort

- Kumyk transducer based on Kazakh, Tatar transducers
- ~1 week to reach 80% coverage, +1 week to reach 90%



Categorisation

- Example: morphologically distinct adjective classes (not documented elsewhere)
- Other Turkic transducers: 0-derivation (overgenerates)
- Our approach: categorisation (generates/analyses only correct forms)

Type	Gloss	<adj> (<comp>)	<adj> (<comp>) <subst>	<adj> (<comp>) <advl>
A1	‘good’	яхшы (яхшырак)	яхшы (яхшырак)	яхшы (яхшырак)
A2	‘old’	искке (искерек)	искке (искерек)	— (—)
A3	‘dead’	үлө (—)	үлө (—)	— (—)
A4	‘basic’	төп (—)	— (—)	— (—)

Evaluation

Number of stems

Part of speech	Kyrgyz	Kazakh	Tatar	Kumyk
Noun	4582	2640	2795	2568
Verb	1193	1470	1143	386
Adjective	1211	754	816	219
Proper noun	5887	5701	5361	1443
Adverb	312	171	177	63
Numeral	66	63	63	44
Conjunction	77	46	45	13
Postposition	50	50	43	12
Pronoun	51	32	28	17
Determiner	64	39	34	9
Total:	13749	11224	10737	4845

Test corpora

	Wikipedia	News	Religion
Kyrgyz	Wikipedia	azattyk.org	Bible
Kazakh	Wikipedia	azattyq.org	Quran + Bible
Tatar	Wikipedia	tat.tatar-inform.ru	Quran + New Testament
Kumyk	—	yoldash.etosmi.ru	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

	Corpus	Tokens	Coverage (%)	Amb.
Kyrgyz	Wikipedia	5.3M	84.51 ± 2.27	3.56
	News	4.1M	91.43 ± 0.51	4.19
	Religion	215K	91.66 ± 1.81	3.99
(r54474)	Average		89.20 ± 3.48	3.91
Kazakh	Wikipedia	25.6M	85.61 ± 1.37	2.43
	News	3.8M	92.12 ± 2.72	2.88
	Religion	851K	92.49 ± 1.66	2.63
(r50547)	Average		90.07 ± 1.91	2.64
Tatar	Wikipedia	159K	86.35 ± 2.17	2.24
	News	5.2M	89.75 ± 0.07	2.30
	Religion	382K	91.25 ± 2.55	2.24
(r50260)	Average		89.12 ± 1.60	2.26
Kumyk	News	286K	91.10 ± 0.86	1.53
	Religion	227K	92.47 ± 1.03	1.53
(r50300)	Average		91.78 ± 0.94	1.53

- selected & proofed unique random surface forms from news corpora

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

Ongoing and future work

- Disambiguation, more stems, clean up transducers
- Machine translation between these languages
- Bring other Kypchak transducers to comparable performance: Qaraqalpaq, Bashqort, Nogay, Crimean Tatar
- Other Turkic lgs: Uzbek, Uyghur, Chuvash, Yakut, Tuvan, etc.