





moderate

heavy





(1) Кудай Өзү жаратканынын баарына

Jonathan North Washington Indiana University jonwashi@indiana.edu

карап,

Ilnar Salimzyanov Kaзaн (Идел буе) федераль университеты ilnar.salimzyan@gmail.com

өтө жакшы экенин көрдү.

DESIGNING FINITE-STATE MORPHOLOGICAL TRANSDUCERS

FOR KYPCHAK LANGUAGES

Francis M. Tyers UiT Norgga Árktalaš Universitehta francis.tyers@uit.no

tat

Special thanks to: Tolgonay Kubatova Aida Sundetova Ağarahim Sultanmuradov

. Ambiguous characters

қиюда 'chopping down'

дәресләр 'lessons'

ёнкюлер 'darlings

Назарбаевтың 'Nazarbayev's'

артистлар 'artists'

самолётлар 'airplanes'

еллар 'years'

егетләр 'boys'

Have front- and back-vowel readings in native words

/ø, y/ / C _

adjust rules for harmony-forcing characters

sent "back" vowels in Russian words

kaz елдің 'country's'

kum сёзлер 'words'

tat галимнәр 'scientists'

• solution: hairy twol rules cover majority of examples

unaccounted-for words get a harmony-forcing character

Letters that represent front vowels in native words may repre-

native word example Russian word example







•	Turkic lang	guages (SO	OV, ag	glutinative	, vowel	harmony)

	quibul	/quzuq/	/ tbtd1/	'quilluq'
classification	Eastern	Southern	Northern	Western
population of s	speakers			
number	3M	8M-12M	5.4M	430K
primary	Kyrgyzstan	Kazakhstan	Tatarstan	Dagestan
secondary	China, etc.	China, Mongolia	Bashqortostan	<u> </u>
external influe	nces			
Mongolic	moderate	moderate	light	light

. 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 \leftrightarrow an<v><tv><ifi><p1><sg>, anд<n><px1sg><nor 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
- morphotactics implemented in lexc
- morphophonology implemented in twol
- compiled separately; compose-intersected to single transducer алдым \leftrightarrow aл>{D}{I}>м \leftrightarrow aл<v><tv><ifi><p1><sg>
- алдым \leftrightarrow алд> $\{I\}$ м \leftrightarrow алд<n><px1sg><nom>
- 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 papers in LREC proceedings (2012: Kyrgyz, 2014: Kazakh, Tatar, Kumyk)
- And feel free to contact the authors any time!

Құдай Өзінің жаратқа	андарының бәріне	қарап,	өте	жақсы екенін	көрді.
Аллаh Үзе яраткан	н нәрсәләргә	карап,	аларның бик	яхшы икәнен	күрде.
Аллагь Оьзю яратгъа	н затлагъа	къарап,	олар бек	яхшы экенин	гёрген.
God own-his created	[everything/thi	ing-s]-to looked.at,	they/their very	good being	saw.
'God looked at everything	he had created and saw the	at it was verv aood.'	(Bible, Genesis	s 1:31)	
Kyrgyz (kir)	Kazakh (kaz)	Tatar (tat)		Kumyk (kum)	
Кудай Өзү жаратканынын баарына карап, өтө жакшы экенин көрдү.	Құдай Өзінің жаратқандарының бә қарап, өте жақсы екенін көрді.	эріне Аллаh Үзе яра	ткан нәрсәләргә кара кшы икәнен күрде.	п, Аллагь Оьзю яратг	тьан затлагъа яхшы экенин гёрген.
Кудай <n><nom> 03<prn><ref><px3sp><nom> жарат<v><tv><ger_past><px3sp><gen> баары<prn><qnt><px3sp><dat> кара<v><tv><gna_perf> ,<m> — 0т0<adv> жакшы<adj> э<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></m></gna_perf></tv></v></dat></px3sp></qnt></prn></gen></px3sp></ger_past></tv></v></nom></px3sp></ref></prn></nom></n>	Құдай <n><nom> 03<prn><ref><px3sp><gen> жарат<v><tv><ger_past><pl><py3se for="" prop<="" properties="" td=""><td>нәрсә<n><pl><pre>kapa<v><tv><gre>cm> ,<cm> anap<pre>prn><per fuk<adv=""> яхшы<adj></adj></per></pre></cm></gre></tv></v></pre></pl></n></td><td>ox3sp><nom> or_past> dat> dat> na_perf> s><p3><pl>><per> st><px3sp><acc></acc></px3sp></per></pl></p3></nom></td><td>Аллагь<n><nom> Oьз<prn><ref><px3 ярат<v=""><tv><gpr_p зат<n=""><pl><pre>зат<n><pl><pre>sat<n><pl><pre>cm> oлар<pre>onap<pre>prn><pers><pl>бек<adv> яхшы<adj><pre>э<cop><ger_past><pre>rëp<v><tv><past><pre><sent></sent></pre></past></tv></v></pre></ger_past></cop></pre></adj></adv></pl></pers></pre></pre></pre></pl></n></pre></pl></n></pre></pl></gpr_p></tv></px3></ref></prn></nom></n></td><td>perf> oast> perf> o3><pl><nom> opx3sp><acc></acc></nom></pl></td></py3se></pl></ger_past></tv></v></gen></px3sp></ref></prn></nom></n>	нәрсә <n><pl><pre>kapa<v><tv><gre>cm> ,<cm> anap<pre>prn><per fuk<adv=""> яхшы<adj></adj></per></pre></cm></gre></tv></v></pre></pl></n>	ox3sp> <nom> or_past> dat> dat> na_perf> s><p3><pl>><per> st><px3sp><acc></acc></px3sp></per></pl></p3></nom>	Аллагь <n><nom> Oьз<prn><ref><px3 ярат<v=""><tv><gpr_p зат<n=""><pl><pre>зат<n><pl><pre>sat<n><pl><pre>cm> oлар<pre>onap<pre>prn><pers><pl>бек<adv> яхшы<adj><pre>э<cop><ger_past><pre>rëp<v><tv><past><pre><sent></sent></pre></past></tv></v></pre></ger_past></cop></pre></adj></adv></pl></pers></pre></pre></pre></pl></n></pre></pl></n></pre></pl></gpr_p></tv></px3></ref></prn></nom></n>	perf> oast> perf> o3> <pl><nom> opx3sp><acc></acc></nom></pl>
		Tagset			
<n> Noun <iv></iv></n>	Intransitive <nom></nom>		sent> Sentence	<pre><gna_perf></gna_perf></pre>	
<v> Verb <tv></tv></v>	Transitive <gen></gen>		past> Past (Gen	•	(Perfect)
<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>	_ *		ifi> Past	<pre><gpr_past></gpr_past></pre>	
<pre><det> Determiner <pl></pl></det></pre>	Plural <dat></dat>	Dative	` _ 0	ness/Recent)	(Past)
<adj> Adjective <ref></ref></adj>			·	1 poss. <ger_past></ger_past>	
<adv> Adverb <pers< td=""><td>> Personal <cm></cm></td><td>Comma</td><td>(Singula</td><td>I/Piulal)</td><td>(Past)</td></pers<></adv>	> Personal <cm></cm>	Comma	(Singula	I/Piulal)	(Past)

..... Desonorisation (kaz & kir)......

- {N} desonorises to д after a consonant алма- $\{N\}\{I\}$ \rightarrow алманы 'apple-ACC'
- сыр- $\{N\}\{I\}$ → сырды 'secret-ACC' • келатсаң 'if you come' $\{L\}$ desonorises to д after cons. of sonority $\leq l$
- сыр- $\{L\}\{A\}$ р → сырлар 'secret—PL'
- кыз- $\{L\}\{A\}$ р → кыздар 'girl—PL' "L Desonorisation"
- "N Desonorisation" %{N%}:д <=> :VoicedCns %>: __ ;

%{L%}:д <=> :VoicedLowSonCns %>: __ ;

- • 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 ;

...... Morphological & orthographical words......

- өнүктүрөбүзбү? 'will we develop [it]?'
- өнүк<v><tv><caus><aor><pl><pl>+бы<qst>
- 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 appli-

case	form	1SG	2SG	3SP
nominative		-(I)M	-(І)ң	-(c)I
accusative	-NI	-(Í)мдI	-(I)ндI	-(c)IH
genitive	-NIH	-(I)мдIн	-(I)́ндІн	-(c)ІнІн
locative	-DA	-(І)мдА	-(Ï)ңдA	-(с)ІндА
ablative	-DAн	-(Ï)мдAн,	-(I)ндAн ,	-(с)ІнАн
		-(І)мАн	-(І)ңАн	
dative	-GA	-(I) MA	-(І)ңА	-(c)IнA

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

- underlying <px3sp> form used: {s}{I}{n}
- {s} and {n} default to c and н; rules map to null by context
- morphophonology more complicated, morphotactics simpler
-Noun-noun compounds...... a N-N compund 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

cation of phonology):

аба% ырайы:аба% ырай N-INFL-3PX-COMPOUND "weather"

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

solution: separate continuation lexicon (messy rules) LEXICON N1-RUS :%{~%} 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', бестен '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%{∋%}%{c%} NUM-DIG	IT ; ! "төрт"
5:5%{9%}%{c%} NUM-DIG	IT ; ! "бес"
3%0:3%0%{a%}%{3%} NUM	I-DIGIT ; ! "отыз"

..... + vowel letters..... • [a o y] become [яёю] after й and й deletes

- й incorporated into the context of many rules
- additional rules to change the characters and delete original й

"Deletion of й before yoticised vowels" й:0 <=> __ [:YotVow]/[:0 | %>:] ;

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

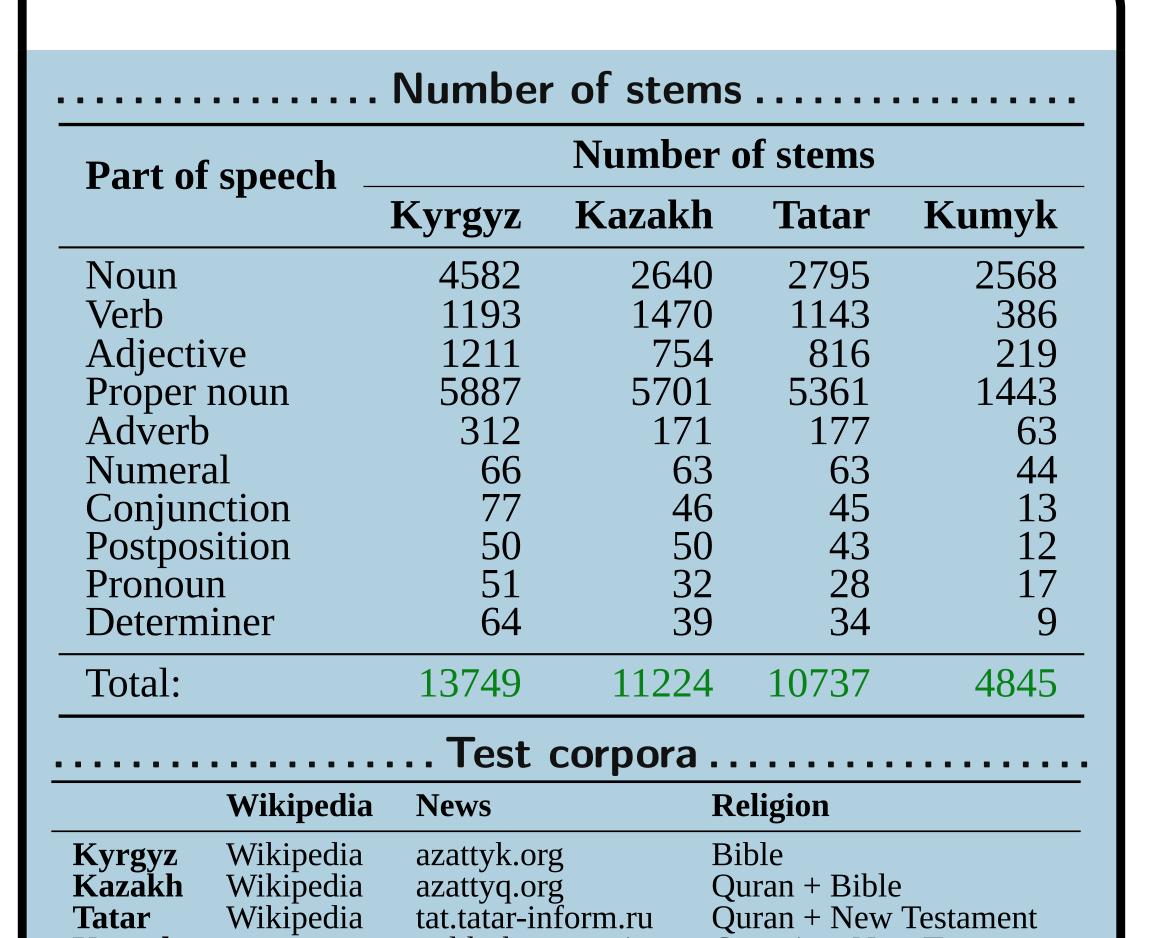
RdYotVow = ë ω Ë Ю ; AbstractVow = %{a%} %{9%} %{γ%} %{o%} ; "A front unrounded harmony" led Harmony
[[:FrontVow | [:Vow :ь]] :Cns :Cns*]/:0 _ ;
[[:RdVow :ь] :Cns :Cns*]/:0 _ ;
[[[\[.#. | :Vow]] :RdYotVow] :Cns :Cns*]/:0 _ ;
[:RdYotVow й:0 :RdYotVow :Cns :Cns*]/[:0 - й:0] _ ; [[%{3%}:0|%{γ%}:0] :Cns*]/[[:0 - AbstractVow:] | %-:]* _ ; except [:Cns :p %{¾%}: %>: :Cns*]/:0 _ ; [[:Vow - :RdYotVow] :RdYotVow :Cns :Cns*]/:0 ; [:Vow]/[[:0 - й:0] | %>:] _ ;

		Kyrgyz	Kazakh	Tatar	Kumyk
begun 80% co	ov.	Apr. 2011 Aug.? 2011	Dec. 2010 Aug. 2012	Dec. 2011 Aug. 2012	Oct. 2013 Oct. 2013
time		4 months	19 months	7 months	1 week
• Kumy	k tra	ansducer base ansducer base o reach 80% (d on Kazakh,	, Tatar transd	
100) [. , , , , , .		. , ,	5000
80) -		•		4000
age (%)) -				- 3000 Stems
Coverage 7) -				- 2000
20) -				- 1000
()	00 00 04 05		erage (%) • Stems □	0
	01	02 03 04 05	06 07 08 Day	09 10 11 12	2 13

..... Adjectives

- morphologically distinct adjective classes
- most sources: adjectives can be used substantively and adver-
- Other Turkic transducers: 0-derivation (overgenerates)
- but not all adjectives have all of the following: comparative forms, substantive readings, adverbial readings
- Our approach: categorisation
- only correct forms are analysed and generated

A4 'basic' төп (—) — (—) — (—)	A1	ʻgood'	яхшы (яхшырак)	яхшы (яхшырак)	яхшы (яхшырак)
	A2	ʻold'	иске (искерәк)	иске (искерәк)	— (—)
	A3	ʻdead'	үле (—)	үле (—)	— (—)
	A4	ʻbasic'	төп (—)	— (—)	— (—)



. Evaluation measures

• Naïve coverage - percentage of surface forms in a given corpus receiving ≥ 1 analysis

voldash.etnosmi.ru Genesis + New Testamen

- Mean ambiguity average number of analyses for each surface form found in analysed corpus
- Precision probability that a provided analysis is accurate
- Recall probability that a certain correct analysis is among those provided by the transducer

..... Evaluation results

	Corpus	Tokens	Coverage (%)	Amb.
	Wikipedia	5.3M	84.51 ± 2.27	3.56
Kyrgyz	News	4.1M	91.43 ± 0.51	4.19
Tty15y2	Religion	215K	91.66 ± 1.81	3.99
(r54474)	Average		89.20 ± 3.48	3.91
	Wikipedia	25.6M	85.61 ± 1.37	2.43
Kazakh	News	3.8M	92.12 ± 2.72	2.88
Kazakh	Religion	851K	92.49 ± 1.66	2.63
(r50547)	Average		90.07 ± 1.91	2.64
	Wikipedia	159K	86.35 ± 2.17	2.24
Tatar	News	5.2M	89.75 ± 0.07	2.30
latai	Religion	382K	91.25 ± 2.55	2.24
(r50260)	Average		89.12 ± 1.60	2.26
	News	286K	91.10 ± 0.86	1.53
Kumyk	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 (%) Kazakh 98.61 95.03

- 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.