









жаратканынын баарына

Құдай Өзінің жаратқандарының бәріне

Jonathan North Washington Indiana University jonwashi@indiana.edu

қарап,

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

DESIGNING FINITE-STATE MORPHOLOGICAL TRANSDUCERS

FOR KYPCHAK LANGUAGES

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

Special thanks to: Tolgonay Kubatova Aida Sundetova Ağarahim Sultanmuradov







_	Turkic languages (S			
	Kvrgvz	Kazakh	Tatar	\mathbf{K}_{1}

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	_
external influe	nces			
3 £ 10			3. 3	1. 1

	Obnaz				moderate
	Persian	heavy	heavy	heavy	heavy
	Russian	heavy	heavy	heavy	heavy
'					

- 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!

Аллаh Үзе яраткан	нәрсәләргә	карап, аларның бик ях	кшы икәнен күрде.
Аллагь Оьзю яратгъаг	н затлагъа	къарап, олар бек ях	кшы экенин гёрген.
God own-his created	[everything/thing-s]-t	o looked.at, they/their very go	ood being saw.
'God looked at everything	he had created and saw that it was	s very good.' (Bible, Genesis 1	.:31)
		put	
Kyrgyz (kir)	Kazakh (kaz)	Tatar (tat)	Kumyk (kum)
Кудай Өзү жаратканынын баарына карап, өтө жакшы экенин көрдү.	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.
Кудай <n><nom> θз<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><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></ger_past></tv></v></nom></px3sp></ref></prn></nom></n>	<pre>Kyдaй<n><nom> 03<prn><ref><px3sp><gen></gen></px3sp></ref></prn></nom></n></pre>	Аллаh <n><nom> Y3<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><sent> .<sent></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></n></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><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><sent> .<sent></sent></sent></p3></past></tv></v></acc></px3sp></ger_past></cop></adj></adv></nom></pl></p3></pers></prn></cm></gna_perf></tv></v></dat></pl></n></gpr_past></tv></v></nom></px3sp></ref></prn></nom></n>
<pre><n> Noun</n></pre>	Intransitive <nom> 'Nomi Transitive <gen> Genitive Third person <acc> Accusa Plural <dat> Dative Reflexive <qnt> Quanti Personal <cm> Comm</cm></qnt></dat></acc></gen></nom>	native' <sent> Sentence ve <past> Past (General ative <ifi> Past</ifi></past></sent>	<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>

...... Desonorisation (kaz & kir)........ • {N} desonorises to д after a consonant

- алма- $\{N\}\{I\}$ \rightarrow алманы 'apple-ACC' сыр- $\{N\}\{I\}$ → сырды 'secret-ACC'
- $\{L\}$ desonorises to д 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 %>: __ ;

..... Epenthesis 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>
- келатсаң 'if you come'
- кел<v><iv><prc impf>+жат<vaux><gna cnd><p2><sg>

өтө жакшы экенин көрдү.

өте жақсы екенін көрді.

.... 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 accusative	 -NI	-(I)м -(I)мдI	-(I)ң -(I)ңдI	-(c)I -(c)Iн
genitive	-NIH	-(І́)мдІн	-(I)́ндІн	-(c)IнIн
locative ablative	-DA -DAн	-(I)мдА -(I)мдАн,	-(I)ңдА -(I)ндАн,	-(с)ІндА -(с)ІнАн
dative	-GA	-(Î)мАн -(I)мА	-(I)ңАн -(I)ңА	-(с)ІнА
$\frac{\text{uative}}{\text{A.I.N.D.G have va}}$				

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

................Noun-noun compounds......

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"

. Ambiguous characters Have front- and back-vowel readings in native words

қиюда 'chopping down' **kaz** и, у, ю дәресләр 'lessons' tat еллар 'years' егетләр 'boys' /ø, y/ / C _

- solution: hairy twol rules cover majority of examples
- unaccounted-for words get a harmony-forcing character
- adjust rules for harmony-forcing characters
- Letters that represent front vowels in native words may reprecent "hack" wowels in Russian words

Sent	Dack vowers in Russi	all wolus
	native word example	Russian word example
kaz tat kum	елдің 'country's' галимнәр 'scientists' сёзлер 'words'	Назарбаевтың 'Nazarbayev's' артистлар 'artists' самолётлар 'airplanes'

• solution: separate continuation lexicon (messy rules)

:%{¾%} N1 ;

LEXICON N1-RUS

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-DIGIT ; !	"төрт"
5:5%{9%}%{c%} NUM-DIGIT ; !	"бес"
3%0:3%0%{a%}%{3%} NUM-DIGIT	; ! "отыз"

- + vowel letters..... • [a o y] become [я ё ю] after й and й deletes
- й incorporated into the context of many rules
- additional rules to change the characters and delete original й

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

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

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] | %>:] _ ;

		k	Cyrgy	yz		Kazak	kh	Ta	ıtar			Kı	ımyk	
beg 80%	gun % cov		Apr. 2 Aug.?			Dec. 2 Aug. 2			ec. 2				et. 2013 et. 2013	
tim	.e	4	mon	iths	1	19 mo	nths	7 1	mon	iths		1 v	veek	_
• Ku	Kazakh transducer based on Kyrgyz transducer Kumyk transducer based on Kazakh, Tatar transducers ~1 week to reach 80% coverage, +1 week to reach 90%													
	100												5000	
	80	_		• •	•	•	•						4000	
Coverage (%)	60	_					·	·				-	3000	Stems
Cover	40	_				·						_	2000	St
	20		·									-	1000	
	0	01 02	2 03	04	05	06 07	Cov 08	erage S 09	e (% tem: 10		12	13	0	

..... 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

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

Number of stems							
Part of speech _	Number of stems						
	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 Ouran + Bible Wikipedia tat.tatar-inform.ru Quran + New Testament voldash.etnosmi.ru Genesis + New Testamen

..... 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** probability that a provided analysis is valid
- **Recall** probability that a certain valid analysis is among those provided by the transducer

.

Evaluation results	5
--------------------	---

	Corpus	Tokens	Coverage (%)	Amb.
	Wikipedia	5.3M	84.51 ± 2.27	3.56
Kyrgyz	News	4.1M	91.43 ± 0.51	4.19
TY1gy2	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
Kazakii -	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
Idlaf	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.