





moderate

heavy





Jonathan North Washington

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







		T/	. Vazalsh	Totak	I/.
•	Turkic	languages ((SOV, agglutinati	ve, vowel ha	rmony)

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> ↔ алд>{I}м ↔ 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!

						Gl o	SS					
(1)	Кудай	Өзү	жаратка	нынын	баарына		карап,		ӨТӨ	жакшы	экенин	көрдү.
	Құдай	Өзінің	жаратқа	ндарының	бәріне		қарап,		өте	жақсы	екенін	көрді.
	Аллаһ	Үзе	яраткан		нәрсәләргә		карап,	аларны	ың бик	яхшы	икәнен	күрде.
	Аллагь	Овзю	яратгъан	H	затлагъа		къарап,	олар	бек	яхшы	экенин	гёрген.
	God	own-his	created		[everything/t]	hing-s]-to	looked.at,	they/th	eir very	good	being	saw.
	'God loo	oked at ev	erything	he had crea	ited and saw tl	hat it was	s very good.'	(Bibl	e, Genesi	is 1:31)		
							<i>v</i> 3	`		,		
Kyrg	yz (kir)			Kazakh (kaz)		····Out	Tatar (tat)	•••••			k (kum)	
_	•	тканынын ба шы экенин к			жаратқандарының қсы екенін көрді		Аллаһ Үзе яра аларның бик я					гъан затлагъа яхшы экенин гёрген.
θз<р жара баары кара	Кудай <n><nom></nom></n>			3sp> <gen></gen>	Aллah <n><nom> Y3<prn><ref>< ярат<v><tv><g нәрсә<n=""><pl>< кара<v><tv><g ,<cm=""> алар<pre> алар<pre> лом нәрсә нәр</pre></pre></g></tv></v></pl></g></tv></v></ref></prn></nom></n>	px3sp> <no pr_past> dat> na_perf></no 		0ьз<ри ярат< зат <n> къара< ,<cm></cm></n>	<pre>><n><nom> cn><ref><px></px><tv><gpr_><pl><dat> <v><tv><gna <="" pre=""></gna></tv></v></dat></pl></gpr_></tv></ref></nom></n></pre>	past>		
3 <c0< th=""><th>s<adj> o><ger_pas v><tv><ifi< th=""><th>t><px3sp><a ><p3><sg></sg></p3></a </px3sp></th><th>CC></th><th>өтe<adv> жақсы<adj> e<cop><ger_pa көр<v><tv><if .<sent></sent></if </tv></v></ger_pa </cop></adj></adv></th><th>ist><px3sp><acc></acc></px3sp></th><th></th><th>бик<adv> яхшы<adj> и<cop><ger_pa күр<v=""><tv><pa< th=""><th>st><px3sp< th=""><th>><acc></acc></th><th>бек<ас яхшы<а э<сор</th><th>dv> adj> ><ger_past> ><tv><past></past></tv></ger_past></th><th>- <px3sp><acc></acc></px3sp></th></px3sp<></th></pa<></tv></ger_pa></cop></adj></adv></th></ifi<></tv></ger_pas </adj></th></c0<>	s <adj> o><ger_pas v><tv><ifi< th=""><th>t><px3sp><a ><p3><sg></sg></p3></a </px3sp></th><th>CC></th><th>өтe<adv> жақсы<adj> e<cop><ger_pa көр<v><tv><if .<sent></sent></if </tv></v></ger_pa </cop></adj></adv></th><th>ist><px3sp><acc></acc></px3sp></th><th></th><th>бик<adv> яхшы<adj> и<cop><ger_pa күр<v=""><tv><pa< th=""><th>st><px3sp< th=""><th>><acc></acc></th><th>бек<ас яхшы<а э<сор</th><th>dv> adj> ><ger_past> ><tv><past></past></tv></ger_past></th><th>- <px3sp><acc></acc></px3sp></th></px3sp<></th></pa<></tv></ger_pa></cop></adj></adv></th></ifi<></tv></ger_pas </adj>	t> <px3sp><a ><p3><sg></sg></p3></a </px3sp>	CC>	өтe <adv> жақсы<adj> e<cop><ger_pa көр<v><tv><if .<sent></sent></if </tv></v></ger_pa </cop></adj></adv>	ist> <px3sp><acc></acc></px3sp>		бик <adv> яхшы<adj> и<cop><ger_pa күр<v=""><tv><pa< th=""><th>st><px3sp< th=""><th>><acc></acc></th><th>бек<ас яхшы<а э<сор</th><th>dv> adj> ><ger_past> ><tv><past></past></tv></ger_past></th><th>- <px3sp><acc></acc></px3sp></th></px3sp<></th></pa<></tv></ger_pa></cop></adj></adv>	st> <px3sp< th=""><th>><acc></acc></th><th>бек<ас яхшы<а э<сор</th><th>dv> adj> ><ger_past> ><tv><past></past></tv></ger_past></th><th>- <px3sp><acc></acc></px3sp></th></px3sp<>	> <acc></acc>	бек<ас яхшы<а э<сор	dv> adj> > <ger_past> ><tv><past></past></tv></ger_past>	- <px3sp><acc></acc></px3sp>
		• • • • • • •				Tag						
<n></n>	Noun		<iv></iv>	Intransitiv		'Nomin			Sentence		gna_perf	
<v></v>	Verb	un	<tv></tv>	Transitive	-	Genitiv Accusa		•	Past (Ger Past		ank nact	(Perfect) Vorbal adjective
<pre><pre><det></det></pre></pre>	_		<p3></p3>	Third pers Plural	on <acc></acc>	Dative	ilive <	ifi>		ness/Rec	gpr_past ent)	Verbal adjective (Past)
<adj></adj>	A 10	_	<ref></ref>	Reflexive		Quanti	fier <	px3sp> 3	3rd perso			
<adv></adv>	A 7	_		Personal	<cm></cm>	Comm			_	ar/Plural)		(Past)

..... 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'
- "L Desonorisation"
- %{L%}:д <=> :VoicedLowSonCns %>: __ ;

кыз- $\{L\}\{A\}$ р → кыздар 'girl—PL'

- "N Desonorisation"
- %{N%}:д <=> :VoicedCns %>: __ ; Epenthesis
- Turn {y} into a harmonised high vowel when a vowel doesn't follow the following consonant:
- $myp{y}H \rightarrow mypyh 'nose'$
- $myp{y}H+{I}M \rightarrow mypдym 'my nose'$ %{y%}:Vy <=> [:LastVowel :Cns* :Cns]/[:0] __
- [:Cns [.#. | :Cns]]/[:0 | %>:]; where Vy in (иүииүыыууыуу) LastVowel in (иүеэөяаёоыюу) matched ;
- HFST transducers are trivially converted to **spell checkers**
- Segmenter, e.g. көрүшкөндөрүмдөнсүңбү: $\kappa \Theta p > \{I\} \sqcup > \{G\} \{A\} H > \{L\} \{A\} p > \{I\} M > \{D\} \{A\} H > C\{I\} H > \{B\} \{I\} \}$

- 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>
- 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	-(I)мдIн	-(I)́ндІн	-(с)ІнІн
locative ablative	-DA -DAн	-(I)мдА -(I)мдАн,	-(I)ңдА -(I)ндАн,	-(c)ІндА -(c)ІнАн
dative	-GA	-(I)мАн -(I)мА	-(I)ңАн -(I)ңА	-(с)ІнА
.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 simplerNoun-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
- аба% ырайы:аба% ырай N-INFL-3PX-COMPOUND
- "weather" чакыруу% кагазы:чакыруу% кагаз N-INFL-3PX-COMPOUND ; ! "invitation"

		Ambiguous	s characters	П	_		
Have	front- and	d back-vowel	readings in native words	ш			
	letters	values	examples	Ш		oegu 30%	
kaz	и, у, ю	/wej, we, jew/ /wej, we, je/	киюд <mark>а</mark> 'chopping down' киюд <mark>е</mark> 'getting dressed'	Ш	_	ime	
tat	e	э / С _ /j/+ы /j/+э	дәресләр 'lessons' еллар 'years' егетләр 'boys'		•]	Kaz Kun ~1 v	ny
cum	ё, ю	/ø, y/ / C _ /jø, jy/	гюнлер 'days' гёзлер 'eyes' юреклер 'hearts'				100
	5, 1 6	/ jø, jy/ / jo, ju/	ёнкюлер 'darlings' юлдузлар 'stars' ёллар 'roads'				80
unac	counted-fo	or words get a	ver majority of examples harmony-forcing character		(%)		60
J		5	ing characters words	Ш	Covera)))	4(
Lette	rs that rep	resent front vo	wels in native words may repre-	ш			
sent '	"back" vov	wels in Russia	n words	ш			20
	native wo	ord example	Russian word example	ш			
kaz tat kum	елдің 'сол галимнәр сёзлер 'w	'scientists'	Назарбаевтың 'Nazarbayev's' артистлар 'artists' самолётлар 'airplanes'				(
soluti	ion: separa	ate continuatio	n 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-DIGIT ; ! "	төрт"
5:5%{9%}%{c%} NUM-DIGIT ; ! "(бес"
3%0:3%0%{a%}%{3%} NUM-DIGIT ;	! "отыз"

- [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 i

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

..... + vowel letters.....

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

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

[:Vow]/[[:0 - й:0] | %>:] _ ;

		K	yrgy	/ Z		Ka	zak	h	7	 Cata	ır			k	ζu	myk	_
_	gun % cov		pr. 2 ug.?		1		c. 20 g. 20	010 012		Dec. Aug			_			t. 2013 t. 2013	
tim	ie	4	mon	ths		19	mor	nths	7	mo	ont	hs		1	W	/eek	
• Ku	zakh ımyk wee	trans	duce	er ba	asec	d on	Ka	zak	h, 7	Γata	r t	ra	nsd				
	100	·····							• • •			·	,			5000	
	80	_		• •			•	•	•]	•				4000	
Coverage (%)	60	-						·	·						-	3000	Stems
Cover	40	-					·								-	2000	Ste
	20	• -] [1000	
	0								vera	Ste	ms					0	
	0	1 02	03	04	05	06	07 Day		08	1()	11	12	2	13		

- 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
- if properly categorised, only correct forms are analysed and generated

	Type	Gloss	<adj>(<comp>)</comp></adj>	<adj>(<comp>)<subst></subst></comp></adj>	<adj>(<comp>)<advl></advl></comp></adj>
l	A1 A2 A3 A4	ʻgood' ʻold' ʻdead' ʻbasic'	яхшы (яхшырак) иске (искерэк) үле (—) төп (—)	яхшы (яхшырак) иске (искерэк) үле (—) — (—)	яхшы (яхшырак) — (—) — (—) — (—)
				Adverhs	

- Certain time adverbs have special attributive and genetive
 - б□г□н, быйыл, кече, жана,

Part of speech		Number o	of stems	
Turt or specen	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 c	orpora		
Wikipedia	a News	R	eligion	
Kyrgyz Wikipedia	azattyk.org	Bi	ble	

Kazakh Wikipedia azattyq.org Quran + Bible Tatar 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

	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

Selected & probled unique failuoin surface forms from he							
Language	Forms	Precision (%)	Recall (%)				
Kyrgyz	500						
Kazakh	1000	98.61	57.98				
Tatar	1000	95.03	85.65				
Kumyk	500	96.57	69.11				

selected & proofed unique random surface forms from news corno

 91.78 ± 0.94 1.53

- Disambiguation, more stems, clean up transducers
- Machine translation between these languages

(r50300) Average

- Bring other Kypchak transducers to comparable performance: Qaraqalpaq, Bashqort, Nogay, Crimean Tatar
- Other Turkic lgs: Uzbek, Uyghur, Chuvash, Yakut, Tuvan, etc.