

Kypchak languages

FINITE-STATE MORPHOLOGICAL TRANSDUCERS FOR THREE KYPCHAK LANGUAGES

Jonathan North Washington

Ilnar Salimzyanov

Francis M Tvers

Also special thanks to



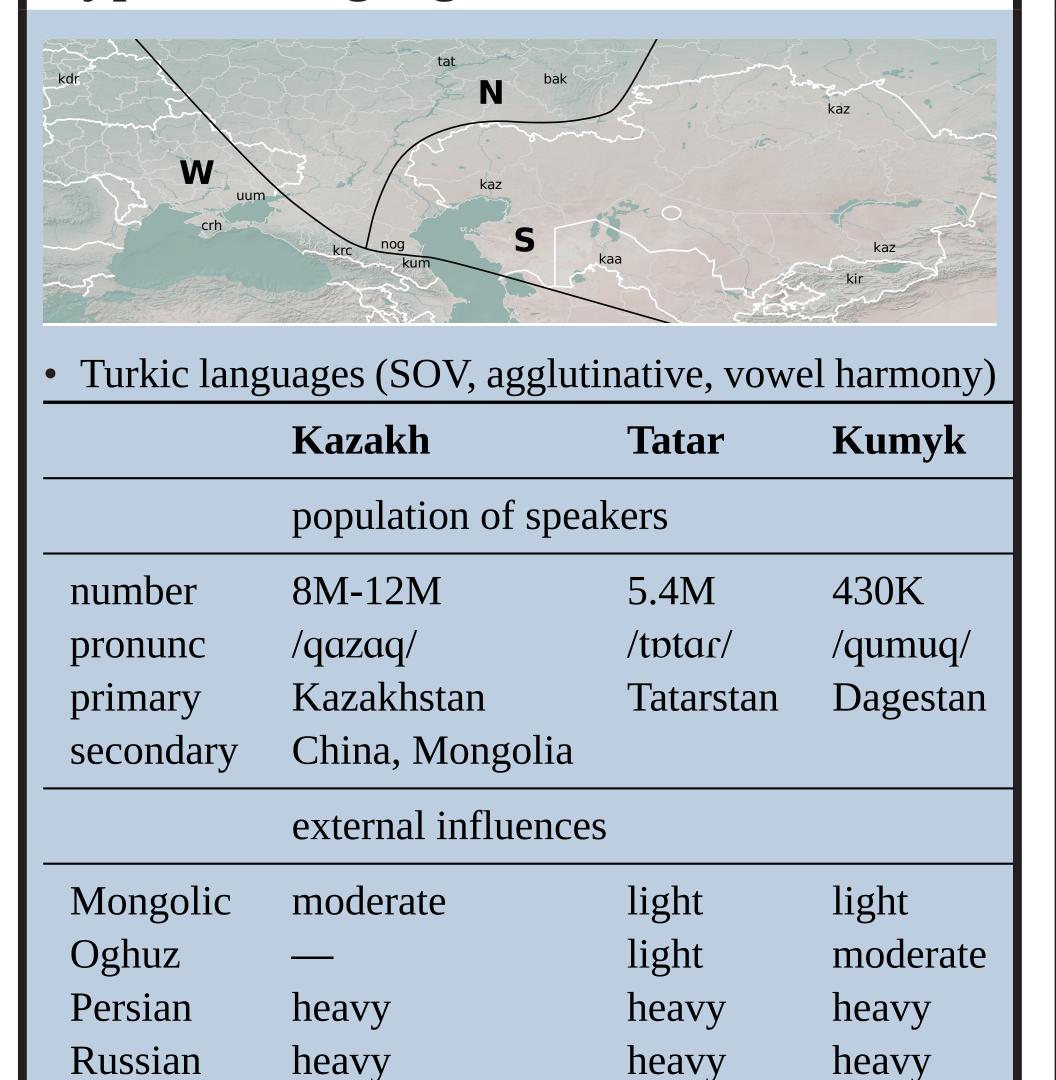
Indiana University

Some University jonwashi@indiana.edu

ilnar.salimzyan@gmail.com

r runcis ivi. Tycis
UiT Norgga Árktalaš Universitel
francis.tyers@uit.no

Aida Sundetova email@email



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><p1><sg>, алд<n><px1sg><nom>
- Transducers for Turkic languages.....
- Turkish (Çöltekin, 2010; Öflazer, 1994)
- Crimean Tatar (Altıntaş, 2001)
- Turkmen (Tantuğ et al., 2006)
- Kyrgyz (Tyers et al., 2012)
- GPL (=free and open)!

..... Framework: HFST..... Reimplementation of Xerox FST formalisms

- (lexc and twol)
- Also provides a wrapper around popular free/opensource FST toolkits: SFST, OpenFST, and Foma Development effort.....

Morphotactics

.... Morphological & orthographical words

- өнүктүрөбүзбү? 'will we develop [it]?' ӨНҮК<v><tv><caus><aor><pl>><pl>+бы<qst>
- келатсаң 'if you come'
- Keл<v><iv><prt impf>+жат<vaux><gna cnd><p2><sg>
- ...Irregular [noun + possessive + case] forms...
- Some combinations of possessive and case morphemes are distinct (i.e., not formed simply by concatenation):

dat	-GA	-(I)мАн -(I)мА	-(I)ңАн -(I)ңА	-(S) Ін А
abl	-DAн	-(І)мдАн,	-(I)ндAн,	-(S) І нАн
loc	-DA	-(І)мдА	-(І)ңдА	-(S) І ндА
gen	-NIH	-(І)мдІн	-(І)ңдІн	- (S)ІнІн
acc	-NI	-(І)мдІ	-(I)ңдI	-(S) I H
nom	_	-(I)M	-(I)ң	-(S)I
case	form	1SG	2SG	3SP

- Trade-off:
- morphophon. complicateder, morphotactics simpler
- underlying form used: {S}{I}{n}
- phonological rules delete {n}, {S} by context

one type of N-N compunds: N2 has <px3> and related morphology

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

COMPOUND ; ! "invitation"

чакыруу% кагазы:чакыруу% кагаз N-INFL-3PX-

Examp	le output						
Kazakh		Tataı	1	Kumyk			
	інің жаратқандарының арап, өте жақсы екенін көрді		h Үзе яраткан нәрсәләргә кара ның бик яхшы икәнен күрде.		Оьзю яратгъан затл , олар бек яхшы эке		
Құдай <n><nom></nom></n>		Yз<р ярат нәрс кара , <cm< td=""><td colspan="2">Аллаh<n><nom> Y3<prn><ref><px3sp><nom> ярат<v><tv><gpr_past> нәрсә<prn><itg><pl><dat> кара<v><tv><gna_perf> ,<cm></cm></gna_perf></tv></v></dat></pl></itg></prn></gpr_past></tv></v></nom></px3sp></ref></prn></nom></n></td><td colspan="3">Аллагь<n><nom> Oьз<prn><ref><px3sp><nom> ярат<v><tv><gpr_past> зат<n><pl><dat> къара<v><tv><gna_perf> ,<cm></cm></gna_perf></tv></v></dat></pl></n></gpr_past></tv></v></nom></px3sp></ref></prn></nom></n></td></cm<>	Аллаh <n><nom> Y3<prn><ref><px3sp><nom> ярат<v><tv><gpr_past> нәрсә<prn><itg><pl><dat> кара<v><tv><gna_perf> ,<cm></cm></gna_perf></tv></v></dat></pl></itg></prn></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></cm></gna_perf></tv></v></dat></pl></n></gpr_past></tv></v></nom></px3sp></ref></prn></nom></n>		
		бик< яхшы и<со	<adj> p><ger_past><px3sp><acc> v><tv><past><p3><sg></sg></p3></past></tv></acc></px3sp></ger_past></adj>	бек <ad яхшы<а э<сор></ad 	dj> <ger_past><px3sp><a <tv><past><p3><sg></sg></p3></past></tv></a </px3sp></ger_past>		
			Gloss				
(1) Алла	агь Оьзю яратгъан затл	паргъа к	ьарап, олар бек яхши	ы экенин	гёрген.		
God	own-his made thin	ig-s-to lo	ook-having, they very good	being	saw.		
^Yстөл/Y ^жана/жа ^отургуч ^астын/а ^карап/к ^жатат/ж ^,/, <cm> ^бирок/б ^Азамат/ ^аякта/а</cm>	стөл <n><nom>\$ н<v><iv><prc_impf>/жана тардын/отургуч<n><pl><q cт<n=""><px3pl><acc>/acт<n apa<v=""><iv><gna_perf>/ка aт<vaux><aor><pre><pre><pre><pre></pre></pre></pre></pre></aor></vaux></gna_perf></iv></n></acc></px3pl></q></pl></n></prc_impf></iv></v></nom></n>	a <adv>/ж gen>\$ n><px3sg apa<v><i жат<vaux >\$ >/аяк<n></n></vaux </i </v></px3sg </adv>	y> <acc>\$ v><prc_real>/кара<v><tv c=""><aor><p3><sg>/жат<vaux e=""><loc>/аякта<v><tv><imp></imp></tv></v></loc></vaux></sg></p3></aor></tv></v></prc_real></acc>	> <gna_per ><prc_irr< th=""><th>f>/kapa<v><tv>< e>\$ (intransitive verb</tv></v></th><th><prc_real></prc_real></th></prc_irr<></gna_per 	f>/kapa <v><tv>< e>\$ (intransitive verb</tv></v>	<prc_real></prc_real>	
^ 3Mec / 3 ^ . / . < sen	cop> <neg><p3><pl>/э<cop< td=""><td>p><neg><</neg></td><td></td><td></td><td></td><td></td></cop<></pl></p3></neg>	p> <neg><</neg>					
<n><n><np><v><cnjcoo><cnjcoo><cnjadv></cnjadv></cnjcoo></cnjcoo></v></np></n></n>	Noun Proper noun Verb Determiner Coord. conjunct. Adv. conjunct.	<pre><p2><p3><pant></pant></p3></p2></pre> <dem><m><sg></sg></m></dem>	Second person Third person Anthroponym Demonstrative Masculine Singular		3rd person poss. 3rd person poss. Negative Aorist Imperative > Verbal adverb (Person		
<adv> <vaux> <cop></cop></vaux></adv>	Adverb Auxiliary verb Copula	<pl><pl><nom></nom></pl><gen></gen></pl>	Plural 'Nominative' Genitive	<pre><pre><pre>c</pre></pre></pre>	>Participle (Impers >Participle (Irrealis >Participle (Realis	s)	

Locative

Accusative

<loc>

Morphophonology

Intransitive

Transitive

<iv>

<tv>

. Desonorisation

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

• 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 ;

......й+vowel letters.....

- [a o y] become [яёю] after й and й deletes
- й incorporated into the context of many rules
- + separate rules to change the characters
- + a rule to delete the original й

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

Part of speech Noun Verb Adjective Proper noun

Adverb

Pronoun

Encyclop

Evaluation

2795 2640 1143 1470 816 754 5361 5701 Numeral Conjunction Postposition

Kazakh

Determiner Total:

11224 10737

kaz

20131006 wpdump 20130225 wpdump

Comma

. Number of stems

Kumyk

2568

386

219

1443

4845

azattyq.org 2010

tat.tatar-inform.ru 20

yoldash.etnosmi.ru

 90.07 ± 1.91

 96.35 ± 2.17

Number of stems

Tatar

News RFE/RL kaz Татар-информ Елдаш

Religion quran + bible kkitap.net, kuran.kz ibt.org.ru, tanzil.net quran + nt genesis + nt ibt.org.ru

split into 10 equal parts; coverage calculated over each separately; standard deviation of mean calculated

...... Coverage measures

 Naïve coverage - percentage of surface forms in a given corpus receiving ≥ 1 analysis (surface forms may have missing analyses)

Mean ambiguity - average number of analyses for each surface form found in analyed corpus

.........Coverage results (as of r36739)...... Corpus **Tokens** Coverage (%) Language Wikipedia 25.6M 85.61 ± 1.37 92.12 ± 2.72 3.8M News Kazakh 851K Religion 92.49 ± 1.66

150K

Average

Wikipodia

Further information