



• Turkic languages (SOV, agglutinative, vowel harmony)				
	Kyrgyz /qɯrɯkɯz/	Kazakh /qɑzɑq/	Tatar /tʰɑtɑr/	Kumyk /qumuq/
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	Bashqortostan	—
external influences				
Mongolic	moderate	moderate	light	light
Oghuz	—	—	light	moderate
Persian	heavy	heavy	heavy	heavy
Russian	heavy	heavy	heavy	heavy

- 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!

Gloss									
(1)	Кудай	Өзү	жаратканынын	баарына	карап,	өтө	жакшы	экенин	көрдү.
	Құдай	Өзінің	жаратқандарының	бәріне	қарап,	өте	жақсы	екенін	көрді.
	Аллаһ	Үзе	яраткан	нәрселәргә	карап,	аларның	бик	яхшы	икәнөн күрде.
	Аллаһы	Өзю	яратгъан	затлагъа	къарап,	олар	бек	яхшы	екенин гѣрген.
	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)									

Kyrgyz (kir)	Kazakh (kaz)	Tatar (tat)	Kumyk (kum)
Кудай өзү жаратканынын баарына карап, өтө жакшы экенин көрдү.	Кудай Өзүнүн жараткандарынын бәрине карап, өтө жакшы экенин көрдү.	Аллаһ Үзө яраткан нәрсәләргә карап, аларның бик яхшы икәнен күрдө.	Аллаһь Өзьә яратгъан затлагъа къарап, олар бек яхшы экенин гөргөн.
Кудай<n><nom> 03-prn<n><ref>px3sp<n> akar<v>-tv>ger<v>px3sp<-gen> baary<prn>-qnt<-px3sp>-dat> kara<v>-tv>sgna_perf> ,<cm> — өтө<adv> жакшы<adj> э<cop>-ger<v>px3sp>-acc> көр<v>-tv>ifi<-p3>-sg> .<sent>	Кудай<n><nom> 03-prn<n><ref>px3sp>-gen> akar<v>-tv>ger<v>px3sp<-pl>-px3sp>-gen> bäri<prn>-qnt<-px3sp>-dat> kara<v>-tv>sgna_perf> ,<cm> — өтө<adv> жакшы<adj> э<cop>-ger<v>px3sp>-acc> көр<v>-tv>ifi<-p3>-sg> .<sent>	Аллаһ<n><nom> Үз<prn>-ref<-px3sp><nom> ayar<v>-tv>ger<v>px3sp<-pl>-px3sp>-gen> närsä<n><pl>-dat> kara<v>-tv>sgna_perf> ,<cm> алар<prn>-pers<p3><pl>-gen> бик<adv> яхшы<adj> и<cop>-ger<v>px3sp>-acc> күр<v>-tv>past<-p3>-sg> .<sent>	Аллаһь<n><nom> Өз<prn>-ref<-px3sp><nom> arat<v>-tv>ger<v>px3sp<-pl>-px3sp>-gen> zat<n><pl>-dat> къара<v>-tv>sgna_perf> ,<cm> адар<prn>-pers<p3><pl>-nom> бек<adv> яхшы<adj> э<cop>-ger<v>px3sp>-acc> рёр<v>-tv>past<p3>-sg> .<sent>

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

- HFST transducers are trivially converted to **spell checkers**
 - Segmenter, e.g. көргөзгөндөрдөнсүнбү :
- көр>{G}{A}з>{G}{A}н>{L}{A}р>{D}{A}н>с{I}н>{B}{I}

case	form	1SG	2SG	3SP
nominative	—	-(I)M	-(I)H	-(c)I
accusative	-NI	-(I)mdI	-(I)ndI	-(c)IH
genitive	-NIH	-(I)mdIH	-(I)ndIH	-(c)IHnH
locative	-DA	-(I)mdA	-(I)ndA	-(c)IHdA
ablative	-DAH	-(I)mdAH,	-(I)ndAH,	-(c)IHdAH
		-(I)mdAH	-(I)ndAH	
dative	-GA	-(I)mA	-(I)HdA	-(c)IHdA

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% караз N-INFL-3PX-COMPOUND
; ! "invitation"

	letters	values	examples
kaz	и, у, ю	/əj, əw, jəw/ /əj, əw, jəw/	кюда 'chopping down' кюде 'getting dressed'
tat	е	э / C ₋ /j/+ы /j/+э	дэреслэр 'lessons' еллар 'years' егетлэр 'boys'
kum	ё, ю	/ø, y/ / C ₋ /jø, jy/ /jo, ju/	гюнлэр 'days' гёзлэр 'eyes' юреклер 'hearts' ёнкюлер 'darlings' юлдузлар 'stars' ёллар 'roads'

- Letters that represent front vowels in native words may sent “back” vowels in Russian words

```
LEXICON N1-RUS
:%{a%} N1 ;

LEXICON Nouns
артист:артист N1-RUS ; ! "artist"
галим:галим N1 ; ! "scientist"
```

5:5%{э}%{с}% NUM-DIGIT ; ! "бес"

"Deletion of й before yoticed vowels"

```

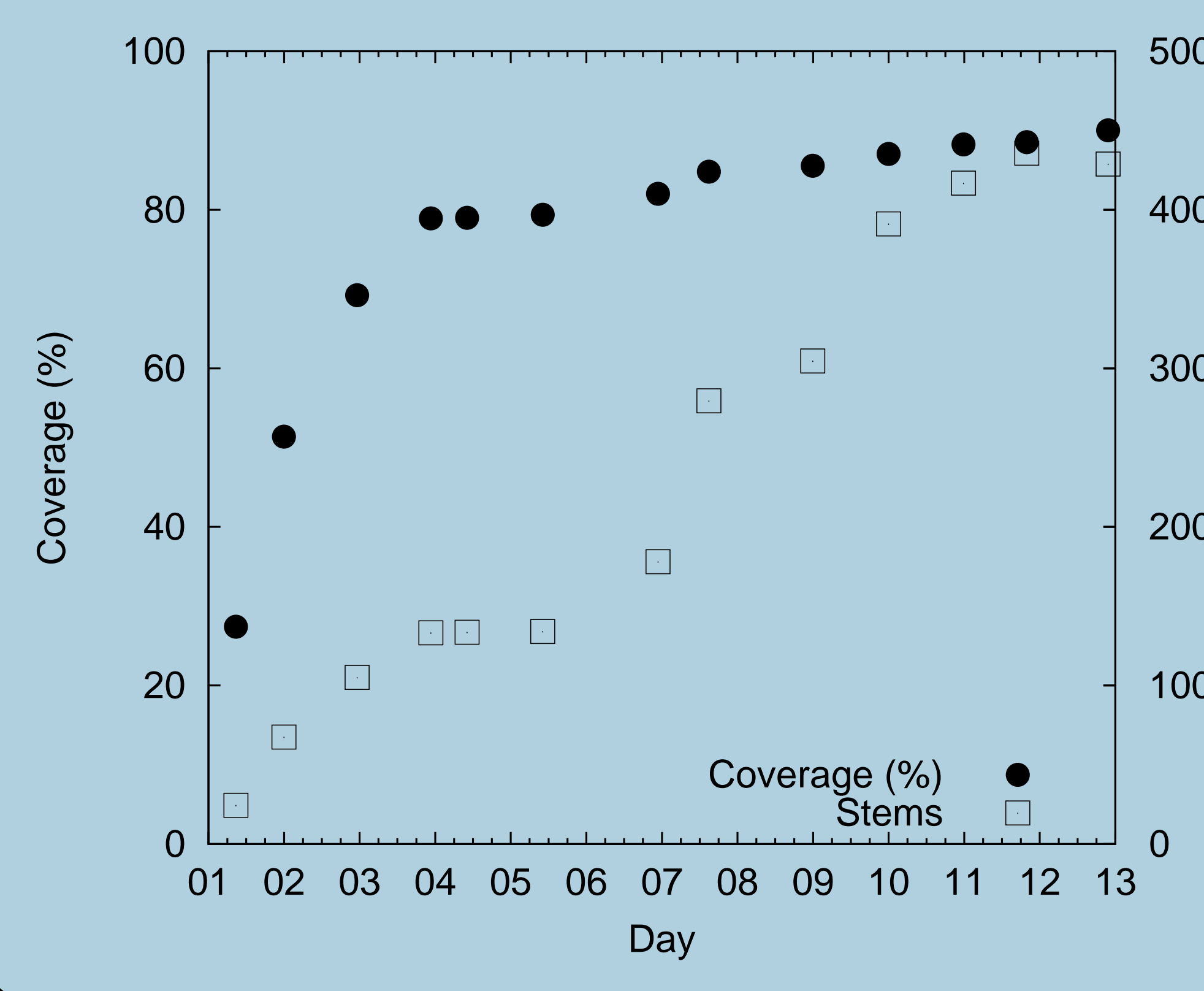
RdyOutVow = e ∪ E ∪ Ø ;
AbstractVow = {a:un} %[%a] %[%a] %[%a] ;

"A front unrounded harmony"
[[RdyVow]:e <=>
    [ [ :FrontVow [ [ :Vow :b ] ] :Cns :Cns* ]/0 :
    [ [ :RdyVow :a ] :Cns :Cns* ]/0 :
    [ [ [ [ :# ] :Vow ] :RdyOutVow ] :Cns :Cns* ]/0 :
    [ [ :RdyOutVow :i:0 ] :RdyOutVow :Cns :Cns* ]/1 :0 : -i:0 :
    [ [%[%]:0] %[%]:0 ] :Cns* ]/[ [ :0 : AbstractVow ] | - : ]*
except
    [ :RdyOutVow :Cns* %[%]:0 :Cns* ]/[ [ :0 : -[%%]:0 ]
    [ :Cns :ip [%%]:0 :Cns* ]/0 :
    [ [ :Vow - :RdyOutVow ] :RdyOutVow :Cns :Cns* ]/0 :
    [ :Vow ]/[ [ :0 : -i:0 ] | :# ] :

```

	Kyrgyz	Kazakh	Tatar	Kumyl
begun	Apr. 2011	Dec. 2010	Dec. 2011	Oct. 2010
80% cov.	Aug.? 2011	Aug. 2012	Aug. 2012	Oct. 2012
time	4 months	19 months	7 months	1 week

- (various periods of intermission, various rewrites)
- Kazakh transducer based on Kyrgyz transducer
- Kyrgyz transducer currently being rewritten based on insights gained while writing other Turkic transducers
- Kumyk transducer based on Kazakh, Tatar transducers: ~1 week to reach 80% coverage, +1 week to reach 90%



Type	Gloss	<adj> (<comp>)	<adj> (<comp>) <subst>	<adj> (<comp>) <
------	-------	----------------	------------------------	------------------

gloss	'today'	'this year'	'yesterday'	'just now'
-------	---------	-------------	-------------	------------

<n><abl> form	бүгүндөн	быйылдан	—	—
---------------	----------	----------	---	---

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

	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

Wikipedia	5.3M	84.51 ± 2.27	3.56
-----------	------	------------------	------

Kyrgyz	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

	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

	Average	286K	1.53
(r50260)		89.12 ± 1.60	2.26
Kumuk	News	91.10 ± 0.86	1.53
	Religion	92.47 ± 1.03	1.53

Runway	Length	227K	52.47 \pm 1.65	1.53
(r50300)	Average		91.78 \pm 0.94	1.53

selected & proofed unique random surface forms from news corpora

Kazakh	1000	98.61	57.98
Tajik	1000	95.93	95.95

Tatal	1000	95.03	83.03
Kumyk	500	96.57	69.11

- Other Turkic lgs: Uzbek, Uyghur, Chuvash, Yakut, Tuva, etc.