



FINITE-STATE MORPHOLOGIES & TEXT CORPORA AS RESOURCES FOR IMPROVING MORPHOLOGICAL DESCRIPTIONS

Francis M. Tyers

UiT Norgga Árkतालá Universitehta
francis.tyers@uit.no

Tommi Pirinen

Ollscoil Chathair Bhaile Átha Cliath
tommi.pirinen@computing.dcu.ie

Jonathan North Washington

Indiana University
jonwashi@indiana.edu



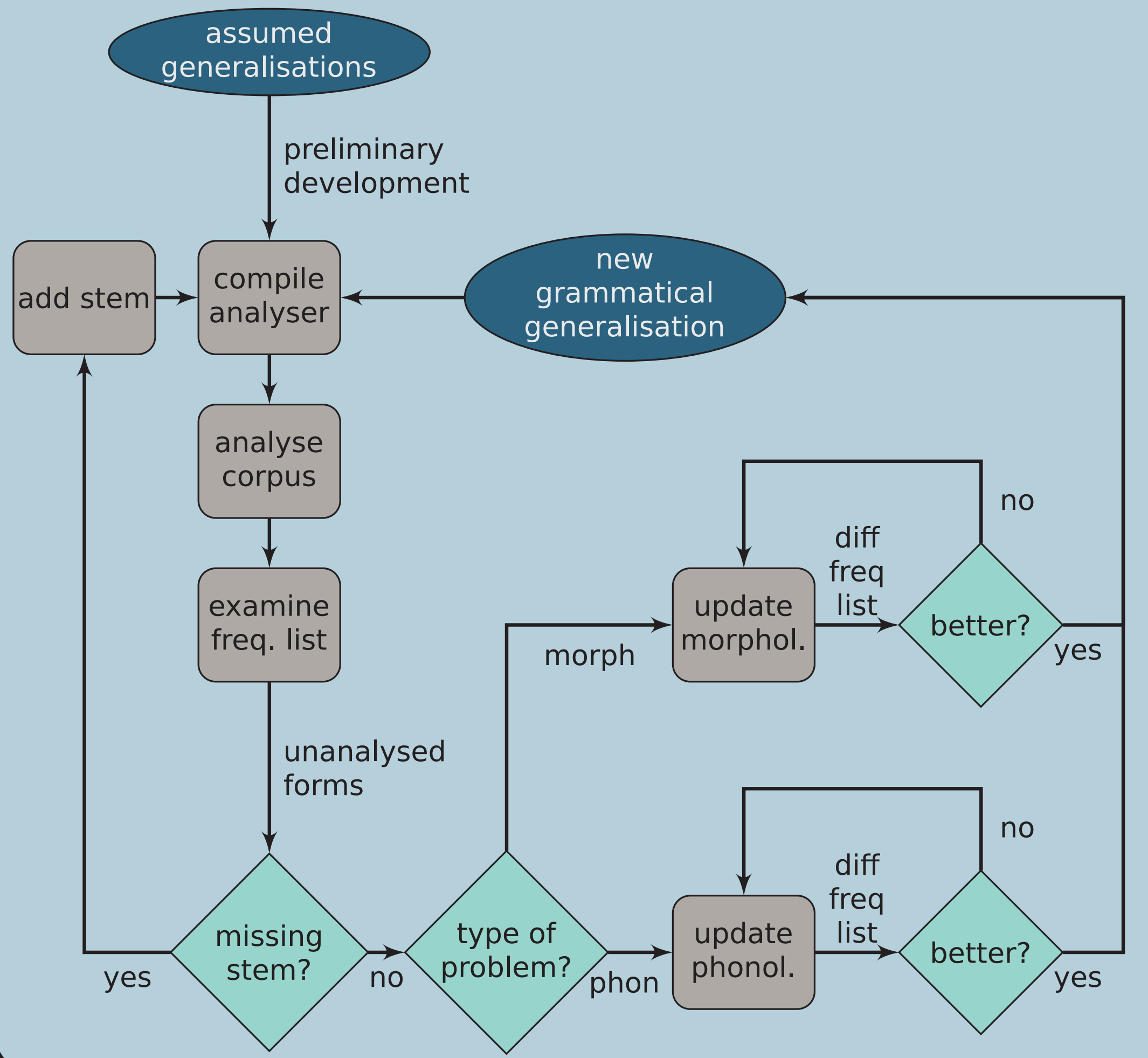
MORPHOLOGICAL DESCRIPTIONS

- **What are they?**
- Core part of grammatical descriptions
- Found in reference grammars and textbooks
- **What they describe**
- **Morphotactics** — the morphemes of a language and how they can be combined
- **Morphophonology** — the alternations between the various phonological and orthographical forms of each morpheme
- **Reasons for being incomplete**
- Restrictions of the medium (e.g., number of pages)
- Limitations of introspection and working with native speakers
- **Lack of automatic testing**

MORPHOLOGICAL TRANSDUCERS

- **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 ↔ ал<v><tv><ifi><pl><sg>, алд<n><pxlsg><nom>
- **Turkic-language transducers we’ve made**
- Kyrgyz (Washington et al., 2012)
- Tatar & Bashkort (Tyers et al., 2012)
- Kazakh (Salimzyanov et al., 2013)
- Kumyk (Washington et al., 2014)
- ongoing work on Sakha, Tuvan, Khakas, and more!
- **Framework: HFST**
- Reimplements Xerox FST formalisms (lexc & twol)
- Also provides a wrapper around popular free/open-source FST toolkits: SFST, OpenFST, and Foma
- **Approach**
- two-level method (Koskenniemi, 1983)
- **morphotactics** implemented in lexc
- **morphophonology** implemented in twol (SPE-style rules)
- compiled separately; compose-intersected to single transducer
- алдым ↔ ал>{D}{I}>М ↔ ал<v><tv><ifi><pl><sg>
- алдым ↔ алд>{I}М ↔ алд<n><pxlsg><nom>

DEVELOPMENT CYCLE



EXAMPLE: TUVAN VOWEL HARMONY

- **Tuvan vowel harmony**
- High and low vowels harmonise in **backness**
кижи<n><pl><acc>: кижилерни, ном<n><pl><acc>: номнарны
- High vowels also harmonise in roundness
round: өр<n><acc>: өрнү, ой<n><acc>: ойну
unround: үе<n><acc>: үени, ай<n><acc>: айны
- **Initial implementation (high-vowels)**
- **Unanalysed forms in corpus**
- forms of loanwords ансамбль ‘ensemble’ and рубль ‘rouble’: the following vowel is always front unround

347 рубльди	28 рубльден	10 ансамблинге
99 рубльге	17 ансамбльдиң	8 ансамбльди
83 ансамблиниң	15 ансамбльдер	1 рубльдерни
60 ансамбли	14 ансамбли	1 рублин
54 рубльдиң	11 ансамблин	1 рубли

- **Revised generalisation**
- After cons. cluster ending in лъ, vowel always front unround
- New discovery? Appears not to be previously documented

..... **Revised implementation**

```
"{I} harmony"
%{I%}:Vy <=> :Vx [ :Cns* :RealCns 1/:0 | %- ]* _ ;
where Vx in ( ү и е э о а о у у я ё ю )
Vy in ( ү и и и у у у у у у у )
matched ;

"{I} always front when intervening Cns"
%{I%}:и <=> [ :BackVow :Cns* :Cns л: ь: :Cns* :RealCns 1/:0* _ ;
[ :BackVow :Cns* :Cns л: ь:0 1/: [ :0 - ь: ]* _ ;
where Vx in ( ү и е э о а о у у я ё ю )
Vy in ( ү и и и у у у у у у у )
matched ;

"{I} always front when intervening Cns"
%{I%}:и <=> [ :BackVow :Cns* :Cns л: ь:0 1/: [ :0 - ь: ]* _ ;
[ :BackVow :Cns* :Cns л: ь:0 1/: [ :0 - ь: ]* _ ;
```

EXAMPLE: CATEGORISATION

- **Unrestricted 0-derivation**
- Common overgeneralisation of adjective morphology:
 - all adjectives may act like nouns (i.e. take noun morphology)
 - all adjectives may be used as is as adverbs
- Commonly reported in Turkic language grammars
- Somfai Kara (2002, pp. 28-29): Adjectives morphologically do not differ from nouns
- All adjectives can be used as adverbs without any morphological changes

- **Attested system**
 - Some adjectives may not be substantivised
 - Some adjectives may not be used adverbially
 - Some adjectives do not have comparative forms
 - ...a range of “adjective classes” in most Turkic languages
 - **Implementation: proper categorisation**
 - only correct forms are analysed and generated (Tatar exs.)
- | Type | Gloss | <adj>(<comp>) | <adj>(<comp>)<subst> | <adj>(<comp>)<adv> |
|------|---------|----------------|----------------------|--------------------|
| A1 | ‘good’ | яхшы (яхшырак) | яхшы (яхшырак) | яхшы (яхшырак) |
| A2 | ‘old’ | iske (искерэк) | iske (искерэк) | — (—) |
| A3 | ‘dead’ | үлө (—) | үлө (—) | — (—) |
| A4 | ‘basic’ | төп (—) | — (—) | — (—) |
- Native speaker intuitions on morph. limitations often more restrictive than range of possible uses found in large text corpora.

FURTHER INFORMATION

- Part of Apertium Turkic project: http://wiki.apertium.org/wiki/Apertium_Turkic
- Transducers available live (turkic.apertium.org), **source code** under **GPL** from Apertium’s SVN repo
- Multilingual Turkic RBMT **mailing list** (>25 subscribers): apertium-turkic@lists.sourceforge.net
- And feel free to contact the authors any time!



EXAMPLE: TUVAN VELAR ELISION

- **Previous descriptions**
- Anderson and Harrison (1999, pp. 9, 22-23)
- **1.3.1 Velar Elision and Vowel Lengthening**
- /g/ > [ɣ] may appear in coda position of a monosyllabic word containing a short vowel: *ɔɣ* ‘yurt’ *ʕaɣ* ‘fat’; it never appears in coda-position in a mono-syllable with a long vowel. When the 3sg poss suffix is added (or any other vowel-initial suffix), [ɣ] is elided and the vowel of the word undergoes lengthening: *ɔɔ* ‘yurt-3, *ʕaa* ‘fat-3. The voiceless velar/uvular stop [k/q] patterns differently. First, [k/q] may appear finally in monosyllables containing either short or long vowels: *ɔk* ‘button’ *ɔɔk* ‘glottis’ *ʕaq* ‘time’. Further, [k/q] does not elide when 3sg poss affix is added: *ɔɣü* button-3 *ɔɔɣü* glottis-3 *ʕaɣi* time-3. However, in polysyllabic words, both /k/ and /g/ pattern together, both undergo elision: *inek* ‘cow’ *inee* cow-3 *ayaq* ‘bowl’ *ayaa* bowl-3 *uruy* ‘daughter’ *uruu* daughter-3 *biliɣ* ‘knowledge’ *bilii* knowledge-3. Both velar phonemes have analogs in the archiphonemic/morphophonological system, see 2.1.1.4 and 2.2.3.1 below.
- Disyllabic stems ending in a velar/uvular sound lose this in possessed forms, and lengthen the preceding vowel. Monosyllabic forms lose a voiced velar but retain a voiceless one (compare ‘yurt’ and ‘glottis’)

(31)	SG	1POSS	2POSS	3POSS	gloss
	<i>tavak</i>	<i>tavaam</i>	<i>tavaaŋ</i>	<i>tavaa</i>	'plate'
	<i>urux</i>	<i>uruum</i>	<i>uruuŋ</i>	<i>uruu</i>	'child'
	<i>balik</i>	<i>baliiim</i>	<i>baliiŋ</i>	<i>balii</i>	'fish'
	<i>biliɣ</i>	<i>biliim</i>	<i>biliŋ</i>	<i>bilii</i>	'understanding'
	<i>belek</i>	<i>beleem</i>	<i>beleŋ</i>	<i>bele</i>	'gift'
	<i>ɔɣ</i>	<i>ɔɔm</i>	<i>ɔɔŋ</i>	<i>ɔɔ</i>	'yurt'
	<i>ɔk</i>	<i>ɔɣüm</i>	<i>ɔɣüŋ</i>	<i>ɔɣüŋ</i>	'glottis'

- Исхаков and Пальмбах (1961, pp. 117-118)
- Конечный *к* в многосложных, а конечный *г* как в многосложных, так и в односложных словах перед аффиксами принадлежности обычно выпадают, создавая условия для образования долгих гласных (но иногда сохраняются):
карак ‘глаз’ — *караа* (<*карагы*), *инек* ‘корова’ — *инээ* (<*инеги*), *балык* ‘рыба’ — *балыы* (<*балыгы*), *бижи́к* ‘письмо’ — *бижи́и* (<*бижиги*), *кудук* ‘колодец’ — *кудуу* (<*кудугу*), *чустүк* ‘колыло’ — *чустүү* (<*чустүгү*), *даг* ‘гора’ — *даа* (<*дагы*), *суз* ‘вода’ — *суу* (<*сугу*), *одаг* ‘костер’ — *одаа* (<*одагы*).

- **More complete characterisation**
- /k/ <к> elides at the end of >1-σ stems intervocalically
инэк+I → инээ, өк+I → өгү
- /g/ <г> elides at the end of any stem intervocalically
өг+I → өө
- /ŋ/ <н> elides at the end of **some** stems intervocalically
соң+I → соо, чаң+I → чаңы, түң+I → түнү
- **Correct forms in the grammars**
- But not part of the morphophonological descriptions.

(A. & H., p. 35)	(<i>soŋ</i> >) <i>soo-</i>	'behind'
(И. & П., p. 447)	§ 521. <i>Соң</i> ‘конец’; <i>соо</i> < <i>соң</i> + аффикс принадлежности 3-го лица: <i>ажылдың соо</i> ‘конец работы’ (ср. <i>баштың соо</i>	

ONGOING AND FUTURE WORK

- Current annotation of a dependency treebank for Kazakh: a similar process that helps identify errors in morphological analyses and automatic disambiguation (constraint grammar)
 - Improved rule-based morphological disambiguation
 - Improved morphological analyser
 - A syntactic dependency treebank
- Future work:
 - Annotate more gold standard text in more Turkic languages
 - Research cross-linguistic techniques for analysing Turkic

..... Unanalysed forms before implementation

- **инээ, өгү, өө, соо, чаңы**

1590 соонда	6 инээниң	2 өөңерге	1 өөңнүң
610 те	6 соонче	2 өөңер	1 өөң
153 өөңүң	6 соонга	2 өөмнү	1 өөвүсче
48 соондан	3 инээм	2 өөвүстен	1 өөвүстүң
46 өөм	3 өөңден	2 өөвүске	1 өөвүстү
20 соон	3 өөңде	2 өөвүс	1 өөвүсте
18 өө	3 өөмнүң	1 инээңни	
17 өөңге	2 инээн	1 инээмниң	
9 соондагы	2 инээ	1 инээм	
8 өөн	2 өөңче	1 инээвистиң	

..... Initial implementation

```
"Intervocalic voiced velar deletion"
r:0 <=> :Vow/:0* _ [ %>: :Vow 1/:0* ;

"Intervocalic voiceless velar deletion"
k:0 <=> :Vow/:0* _ [ %>: :Vow 1/:0* ;
except
[ .#. | %- ] [ ( :Cns* ) ( :Vow* ) :Vow 1/:0 _ [ %>: :Vow 1/:0* ;
```

..... Unanalysed forms after initial implementation

- **инээ, өгү, өө, соо, чаңы**

1590 соонда	48 соондан	6 соонче
610 те	26 соон	6 соонга

..... Second attempt at implementation

```
"Intervocalic voiced velar deletion"
Cx:0 <=> :Vow/:0* _ [ %>: :Vow 1/:0* ;
where Cx in ( р қ ) ;

"Intervocalic voiceless velar deletion"
k:0 <=> :Vow/:0* _ [ %>: :Vow 1/:0* ;
except
[ .#. | %- ] [ ( :Cns* ) ( :Vow* ) :Vow 1/:0 _ [ %>: :Vow 1/:0* ;
```

..... Unanalysed forms after second implementation

- **инээ, өгү, өө, соо, чаңы**

610 те	56 түнү	5 түңүнден
203 бажыңыңа	47 чаңы	4 чаңыңдан
164 бажыңының	44 чаңын	4 чаңында
102 бажыңы	13 түңүн	3 түңүнүң
86 бажыңын	12 чаңының	2 түңүңге

..... Final implementation

```
"Intervocalic voiced velar deletion"
Cx:0 <=> :Vow/:0* _ [ %>: :Vow 1/:0* ;
except
:Vow _ [ %>{&}: :Vow 1/:0* ;
where Cx in ( р қ ) ;

"Intervocalic voiceless velar deletion"
k:0 <=> :Vow/:0* _ [ %>: :Vow 1/:0* ;
except
[ .#. | %- ] [ ( :Cns* ) ( :Vow* ) :Vow 1/:0 _ [ %>: :Vow 1/:0* ;
```

..... Unanalysed forms after final implementation

- **инээ, өгү, өө, соо, чаңы**

610 те

- All forms of velar deletion are now analysed correctly!
- Can now move on to other unanalysed forms...

REFERENCES

- Anderson, Gregory David and K. David Harrison (1999). *Tyvan*. Vol. 257. Languages of the World/Materials. München: Lincom Europa.
- Koskenniemi, K. (1983). *Two level morphology: a general computational model for wordform recognition and production*. Helsinki: Helsingin yliopisto.
- Salimzyanov, Ilmar, Jonathan North Washington, and Francis Morton Tyers (2013). “A free/open-source Kazakh-Tatar machine translation system”. In *Proceedings of Machine Translation Summit XIV*. Nice, France.
- Somfai Kara, Dávid (2002). *Kazak*. Vol. 417. Languages of the World/Materials. München: Lincom Europa.
- Tyers, Francis, Jonathan North Washington, Ilmar Salimzyanov, and Rustam Bat-alov (2012). “A prototype machine translation system for Tatar and Bashkir based on free/open-source components”. In *Proceedings of the First Workshop on Language Resources and Technologies for Turkic Languages at the Eighth International Conference on Language Resources and Evaluation (LREC’12)*. İstanbul, Turkey.
- Washington, Jonathan, Mirlan Ipasov, and Francis Tyers (2012). “A finite-state morphological transducer for Kyrgyz”. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*. İstanbul.
- Washington, Jonathan North, Ilmar Salimzyanov, and Francis M. Tyers (2014). “Finite-state morphological transducers for three Kypchak languages”. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*. Reykjavik, Iceland.
- Исхаков, Ф. Г. and А. А. Пальмбах (1961). *Грамматика Тувинского языка: фонетика и морфология*. Москва: Издательство восточной литературы.