

A free/open-source Kazakh-Tatar machine translation system

Ilmar Salimzyanov
Kazan Federal University
Kazan, Republic of Tatarstan
Russian Federation
ilmar.salimzyan@gmail.com

Jonathan North Washington
Departments of Linguistics
and Central Eurasian Studies
Indiana University
Bloomington, Indiana 47405 USA
jonwashi@indiana.edu

Francis Morton Tyers
Departament de Llenguatges
i Sistemes Informàtics
Universitat d'Alacant
E-03877 Alacant
ftyers@dlsi.ua.es

Abstract

This paper presents a bidirectional machine translation system between Kazakh and Tatar.

1 Introduction

This paper presents a prototype shallow-transfer rule-based machine translation system between Kazakh and Tatar.

The paper will be laid out as follows: Section 3 gives a brief description of the two languages; Section 2 gives a short review of some previous work in the area of Turkic-Turkic language translation; Section 4 describes the system and the tools used to construct it; Section 5 gives a preliminary evaluation of the system; and finally Section 6 describes our aims for future work and some concluding remarks.

2 Previous work

Within the Apertium project, work on MT systems between Turkic languages has been started (Turkish-Kyrgyz, Azeri-Turkish), but the Kazakh-Tatar system described by the present study is the closest to production-ready of them. Among these systems is a prototype Tatar-Bashkir machine translation system which was built by the authors of this paper (Tyers et al., 2012); due to the closeness of these languages, it proved to provide high accuracy in its translations, but being a prototype system by design, had relatively low coverage.

Besides these systems, several previous works on making machine translation systems between Turkic languages exist, although to our knowledge none are publicly available. Some MT systems have been reported that translate between Turkish and other Turkic languages, including Turkish-Crimean Tatar (Altintas, 2001b), Turkish-Azerbaijani (Hamzaoglu, 1993), Turkish-Tatar (Gilmullin, 2008), and Turkish-Turkmen

(Tantuğ et al., 2007), though none of these have been released to a public audience.

Within the Apertium project, work on several Turkic-language MT systems has been started, with the ultimate goal of creating independent but compatible finite-state transducers for each language.

3 Languages

Both Tatar and Kazakh belong to the Kypchak (or North-western) group of Turkic languages. The spoken and written languages share some level of mutual intelligibility to naïve native speakers, though this is somewhat limited, and is obscured by different orthographical conventions and some opaque correspondences.

Kazakh is primarily spoken in Kazakhstan, where it is the national language, sharing official status with Russian as an official language. Large groups of native speakers also exist in China, neighbouring Central-Eurasian republics, and Mongolia. The total number of speakers is at least 10 million people.

Tatar is spoken in and around Tatarstan by approximately 6 million people. It is co-official with Russian in Tatarstan — a republic within Russia. A majority of native speakers of both languages are bilingual in Russian.

3.1 Phonological differences

As closely related languages, Kazakh and Tatar share many phonological processes, including front-back vowel harmony systems, consonant voicing assimilation, and even a typologically rare consonantal nasal harmony system. However, the differing details of these processes and the existence of processes unique to each language render Kazakh and Tatar fairly different. For example, Kazakh has a ubiquitous system of desonorisation of the initial sonorants found in many common morphemes. Furthermore, Tatar has nasal assimilation of the initial /l/ of the plural-suffix.

3.2 Orthographic differences

The standard varieties of Kazakh and Tatar our system deals with are both written in Cyrillic, though their implementations of Cyrillic differ in many ways.

While Tatar and Kazakh both have a velar/uvular obstruent distinction (e.g., /k/ vs. /q/) that interacts with adjacent vowels, the Tatar orthography only has one series of letters (e.g., <к>), relying on adjacent vowels (and employing “hard” and “soft signs” when these fail) to differentiate the two, and Kazakh has two series of obstruents (e.g., <к> and <қ>).

Kazakh does not orthographically distinguish high unrounded vowels (/ɘ/ <и> and /ə/ <ы>) before glides (/w/ <у> and /j/ <й>) by writing the combination with one letter; i.e., /ɘj/ and /əj/ are both written <и>, while /əw/ and /əw/ are both written <у>. The quality of these vowels is necessary to know in order to predict the quality of following harmonising vowels. Additionally, Tatar and Kazakh both use “yoticed” vowels—i.e., when <о>, <у>, or <а> (along with <ə> in Tatar) follow /j/, a single character is used to represent both: <ё>, <ю>, and <я> respectively.¹

All of these orthographical conventions present acute challenges to designing accurate morphological transducers for the languages.

3.3 Morphological differences

There are a number of examples where the morphologies of Kazakh and Tatar are rather different, including morphemes in one language that don’t exist in the other, entirely different uses of the same morpheme combinations, and morphotactic differences (i.e., allowable ordering and placement of morphemes).

An example of a morpheme that doesn’t exist in one of the two languages is Kazakh {-E}т{I}н, which is used to form non-past verbal adjectives and verbal nouns. The semantically equivalent structure in Tatar is {-E}торган, which historically corresponds to the source of the Kazakh morpheme; however, the use of {-E}тұрған in modern Kazakh is different from that of {-E}т{I}н.

Another example of a far-reaching morphological difference between Tatar and Kazakh is the presence of a four-way distinction in Kazakh’s 2nd person system (both pronouns and agreement suffixes), where Tatar only has a two-way distinction. Kazakh has a distinct pronoun for all combinations of [±plural, ±formal], whereas Tatar collapses all pronouns except the [-plural, -formal] into one pronoun, as summarised in table 1.

¹Furthermore, in Tatar, /j/ followed by <ə> or <ы> in Tatar is represented by <е>, though <е> is also the non-word-initial variant of <ə>.

	-pl	+pl		-pl	+pl
-formal	сен	сендер	-formal	син	сез
+formal	сиз	сиздер	+formal	сез	сез

(a) Kazakh 2nd person pronouns (b) Tatar 2nd person pronouns

Table 1: The 2nd person pronoun systems of Kazakh and Tatar

This systematic difference would seem to be a minor issue, since, as is typical in pro-drop languages, pronouns are only used for emphasis and clarification. However, this difference between Tatar and Kazakh in the second person system runs much deeper than just the pronoun system. Since all finite verb forms morphologically agree in person and number with their subject and all possessed nouns agree in person and number with their possessor (even when there is no overt pronoun, in either situation), the Kazakh and Tatar systems of agreement suffixes reflect the same pattern; i.e., there are several sets of agreement morphemes which have a one-to-one correspondence with the pronouns in each language, resulting in several systems of suffixes in each language that have the same set of distinctions as in the 2nd person pronoun systems.

The past tense systems of Kazakh and Tatar have a many-to-many correspondence. As shown in table 2, at a basic level, in the past tense, Kazakh differentiates [±eyewitness]² (where [-eyewitness] is used for cases of both potentially unreliable information and newly discovered information) and [±recent], whereas Tatar has only three categories: eyewitness, non-eyewitness, and newly-acquired-information—all with no [±recent] distinction. As an example of the many-to-many correspondence that this results in, Tatar has a single non-eyewitness past tense morpheme (-GAн-) while Kazakh has a recent non-eyewitness past (-Ин-) and a distant non-eyewitness past (-GAн екен-). On the other hand, these two non-eyewitness past forms in Kazakh are used for both potentially unreliable information and newly acquired information, whereas in Tatar, non-eyewitness (-GAн-) and newly-acquired-information (-GAн- икән) past forms are distinguished.

Without regard to the semantic alignment of these forms, the morphotactics of the cognate Kazakh distant non-eyewitness past (-GAн екен-) and Tatar newly-acquired-information past (-GAн- икән) are different. Specifically, in both languages, the person agreement

²“Eyewitness” is a convenient term for this feature, though it may be better expressed as simply “reliability of knowledge” (which indeed often equates to whether the knowledge was acquired first-hand or not) in many cases.

	[+recent]	[-recent]
[+reliable]	-DI-	-ГАН-
[-reliable]	-Іп-	-ГАН екен-

(a) Kazakh past tense morphology

	[-newlyAcq'd]	[+newlyAcq'd]
[+reliable]	-DI-	
[-reliable]	-ГАН-	-ГАН- икән

(b) Tatar past tense morphology

Table 2: A comparison of the basic past-tense morphology of Kazakh and Tatar

takes the form of a person copula suffix, although in Kazakh this suffix follows the tense morphemes (e.g., барған екенсің “apparently you went”), whereas in Tatar this suffix intervenes between the two pieces of the “compound” tense morpheme (e.g., баргансың икән “I guess you went”).

Another morphotactic difference between Kazakh and Tatar is found with the negative forms of the cognate -ГАН- past tenses. In Kazakh, the negative form of the non-recent reliable-information past tense is -ГАН емес-, whereas in Tatar, the negative form of the non-eyewitness past tense is -МАГАН-.

3.4 Syntactic differences

There are a number of syntactic differences in Tatar and Kazakh, which include differences in verb valencies in equivalent translations, FIXME, and FIXME.

An example of a difference in verb valencies is with the expression corresponding to “to like something” in Kazakh and Tatar. In Kazakh, the verb ұна is used, e.g. бауырсақ маған ұнайды “I like bawyrdaq”, where the subject “I” in English is expressed through a dative experiencer in Kazakh and the object “bawyrdaq” in English is the grammatical subject in Kazakh. Tatar, on the other hand, uses a verb whose arguments correspond to the arguments of “to like” in English, e.g. мин бавырсақ яратам “I like bawyrdaq”, where the first person pronoun is the grammatical subject and “bawyrdaq” is the grammatical direct object (with no accusative suffix since it is indefinite).

In Kazakh, a gerund (i.e., verbal noun) with case marking and sometimes person agreement in the form of possessive suffixes is used to make a verb phrase an argument to certain other main phrases. In Tatar, many of these phrases use an invariant infinitival form. Some examples are shown below:

- (1) a. Мен үйге баруым керек.
мен үй-ГА бар-у-Ім керек
I home-DAT go-GER-1SG need
‘I need to go home.’
- b. Мин өйгә барырга кирәк.
мин өй-ГА бар-ІргА кирәк
I home-DAT go-INF need
‘I need to go home.’

> Tatar infinitive (Kazakh ger+nom/acc/dat, Kazakh -ГАII vadv)

4 System

The system is based on the Apertium machine translation platform (Forcada et al., 2011).³ The platform was originally aimed at the Romance languages of the Iberian peninsula, but has also been adapted for other, more distantly related, language pairs. The whole platform, both programs and data, are licensed under the Free Software Foundation’s General Public Licence⁴ (GPL) and all the software and data for the 30 supported language pairs (and the other pairs being worked on) is available for download from the project website.

4.1 Architecture of the system

The Apertium translation engine consists of a Unix-style *pipeline* or *assembly line* with the following modules (see Fig. 1):

- A *deformatter* which encapsulates the format information in the input as *superblanks* that will then be seen as blanks between words by the other modules.
- A *morphological analyser* which segments the text in surface forms (SF) (*words*, or, where detected, multi-word lexical units or MWLUs) and for each, delivers one or more *lexical forms* (LF) consisting of *lemma*, *lexical category* and morphological information.
- A *morphological disambiguator* (constraint grammar) which chooses, using linguistic rules the most adequate sequence of morphological analyses for an ambiguous sentence.
- A *lexical transfer* module which reads each SL LF and delivers the corresponding target-language (TL) LF by looking it up in a bilingual dictionary encoded as an FST compiled from the corresponding XML file. The lexical transfer module may return more than one TL LF for a single SL LF.

³<http://www.apertium.org>

⁴<http://www.fsf.org/licenses/licenses/gpl.html>

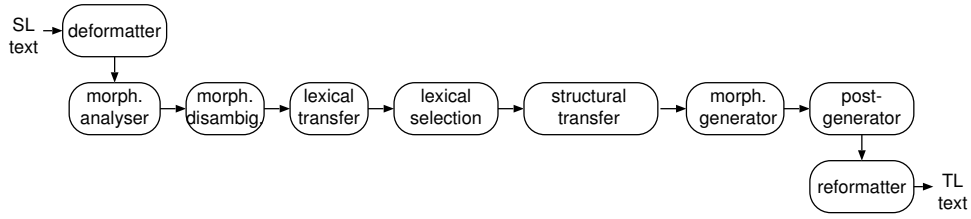


Figure 1: The pipeline architecture of the Apertium system.

- A *lexical selection* module which chooses, based on context rules the most adequate translation of ambiguous source language LFs.
- A *structural transfer* module which performs local syntactic operations, is compiled from XML files containing rules that associate an *action* to each defined LF *pattern*. Patterns are applied left-to-right, and the longest matching pattern is always selected.
- A *morphological generator* which delivers a TL SF for each TL LF, by suitably inflecting it.
- A *reformatter* which de-encapsulates any format information.

4.2 Morphological transducers

The morphological transducers are based on the Helsinki Finite State Toolkit (Linden et al., 2011), a free/open-source reimplementation of the Xerox finite-state toolchain, popular in the field of morphological analysis. It implements both the **lexc** formalism for defining lexicons, and the **twol** and **xfst** formalisms for modeling morphophonological rules. It also supports other finite state transducer formalisms such as **sfst**. This toolkit has been chosen as it — or the equivalent XFST — has been widely used for other Turkic languages (Çöltekin, 2010; Altintas, 2001a; Tantuğ et al., 2006; Washington et al., 2012; Tyers et al., 2012), and is available under a free/open-source licence.

The morphologies of both languages are implemented in **lexc**, and the morphophonologies of both languages are implemented in **twol**.

Use of **lexc** allows for straightforward definition of different word classes and subclasses. For example, Tatar (but not Kazakh) has two classes of verbs: one which takes a harmonised high vowel in the infinitive (the default), and one which take a harmonised low vowel in the infinitive. This was implemented in **lexc** with two similar continuation lexica for verbs: one pointing at a lexicon with an A-initial infinitive ending, and another pointing at a lexicon with an I-initial infinitive ending.

Use of **twol** allows for phonological processes present in the languages, like vowel harmony and desonorisa-

tion, to be implemented in a straightforward manner. For example, in Tatar, the A and I archiphonemes found in the infinitive are harmonised to one of two vowels each, depending on the value of the preceding vowel; the basic form of this process can be implemented in one twol rule.

The same morphological description is used for both analysis and generation. To avoid overgeneration, any alternative forms are marked with one of two marks, **LR** (only analyser) or **RL** (only generator). Instead of the usual compile/invert to compile the transducers, we compile twice, once the generator, without the **LR** paths, and then again the analyser without the **RL** paths.

4.3 Bilingual lexicon

The bilingual lexicon currently contains 9,269 stem to stem correspondences and was build mostly by hand. Part of toponyms and other proper names was translated semi-automatically looking up in Wikipedia links (also some of the Russian loanwords were the result of intersection of Russian and Kazakh wordlists - "автомобиль" etc).

Entries consist largely of one-to-one stem-to-stem correspondences with part of speech, but also include some entries with ambiguous translations (see e.g., Fig. 2).

4.4 Disambiguation rules

The system has a morphological disambiguation module in the form of a Constraint Grammar (CG) (Karlsson et al., 1995). The version of the formalism used is vislcg3.⁵

The grammar currently has only four rules, but given the closeness of the languages, the majority of ambiguity may be passed through from one language to the other.

4.5 Lexical selection rules

Likewise, lexical selection is not a large problem between Tatar and Bashkir, but a number of rules can be written for ambiguous words; for example, the Tatar word *борын* ‘nose (person), nose (ship)’ can be translated into Bashkir as either *манай* ‘nose (person)’ or

⁵http://beta.visl.sdu.dk/constraint_grammar.html

```

<e><p><l>күрал<s n="">/<l><r>корал<s n="">/<r></p></e>
<e><p><l>күрал<s n="">/<l><r>чара<s n="">/<r></p></e>
<e><p><l>борын<s n="">/<l><r>морон<s n="">/<r></p></e>
<e><p><l>ераклык<s n="">/<l><r>алыслык<s n="">/<r></p></e>
<e><p><l>ераклык<s n="">/<l><r>йыраклык<s n="">/<r></p></e>

```

Figure 2: Example entries from the bilingual transfer lexicon. Kazakh is on the left, and Tatar on the right

морон ‘nose (ship)’. A lexical selection rule chooses the translation *танау* if the immediate context includes a proper name.

Another example is the word *катлаулы* ‘layered’. It is always translated to Bashkir as *катмарлы*, except in the collocation *катлаулы мәсьәлә* ‘difficult matter/problem’, which is translated as *катлаулы мәсьәлә*.

5 Evaluation

Lexical coverage of the system is calculated over a freely available corpus of Bashkir, the Bashkir Wikipedia,⁶ and over two freely available corpora of Tatar, the Tatar Wikipedia⁷ and the New Testament in Tatar. The version of the translation tested was r37137 from the Apertium SVN.⁸

As shown in Table 4, the coverage is still far too low to be of use as a general broad-domain MT system, but we hope that it shows that a good proportion of the morphology of both languages is in place.

To get an idea of the kind of performance that could be expected from the system, we translated a simple story from Tatar to Bashkir and vice versa. The story may be found online,⁹ and was used for pedagogical purposes in a recently workshop on MT for the languages of Russia.

Table 5 presents the Word Error Rate, an edit metric based on the Levenshtein distance (Levenshtein, 1966). This measure was calculated once all the stems in the text had been added to the system, thus presents an upper bound on the current performance of the transfer lexicon, and the disambiguation and transfer rules. The difference in the number of unknown words between translating Tatar→Bashkir and vice versa is because certain forms were not found due to lack of corresponding morphophonological rules.

We calculate the WER instead of other MT evaluation metrics such as BLEU as the WER is geared towards a particular task, that of measuring postedition effort. The translations of the story into Tatar and Bashkir were done

in parallel to make them as close as possible, so using BLEU would give an over-optimistic view of the quality.

5.1 Error analysis

The majority of errors are currently due to mistakes and gaps in the morphophonology component; some minor problems still remain involving:

- Combinations of case and possessive suffixes,
- Orthographical representations of phonology,
- Vowel harmony processing on clitics (e.g., *да/да* ‘and’) after unknown words.

6 Concluding remarks

To our knowledge we have presented the first ever MT system between Tatar and Bashkir, and the first ever MT system involving Bashkir. The system is available as free/open-source software under the GNU GPL and the whole system may be downloaded from SVN.¹⁰

We plan to continue development on the pair; the main work will be expanding the dictionaries with new lists of stems, and providing bilingual correspondences. The long-term plan is to integrate the data created with other open-source data for Turkic languages in order to make transfer systems between all the Turkic language pairs. Related work is currently ongoing with Chuvash–Turkish and Turkish–Kyrgyz.

Acknowledgements

We would like to thank the anonymous reviewers for their helpful comments in improving the paper. This work has been partially funded by Spanish Ministerio de Ciencia e Innovación through project TIN2009-14009-C02-01.

References

- Altintas, Kemal. 2001a. A morphological analyser for Crimean Tatar. *Proceedings of Turkish Artificial Intelligence and Neural Network Conference*.
- Altintas, Kemal. 2001b. Turkish To Crimean Tatar Machine Translation System. Master’s thesis, Bilkent University.

⁶<http://ba.wikipedia.org/articles.xml.bz2> bawiki-20111210-pages-

⁷<http://tt.wikipedia.org/articles.xml.bz2> ttwiki-20111215-pages-

⁸<https://apertium.svn.sourceforge.net/svnroot/apertium>

⁹<https://apertium.svn.sourceforge.net/svnroot/apertium/branches/xupaixkar/rasskaz>

¹⁰<https://apertium.svn.sourceforge.net/svnroot/apertium/nursery/apertium-tt-ba>

(Kazakh) Input	Ол енді ол дыбысты анығырақ ести бастады
Mor. analysis	[^] Ол/ол<det><dem>/Ол<prn><dem><nom>/ Ол<prn><pers><p3><sg><nom>\$ [^] енді/ен<n><acc>/ен<v><iv><ifi><p3><pl>/ен<v><iv><ifi><p3><sg>/ енді<adv>\$ Ол/Ол<det><dem>/Ол<prn><dem><nom>/Ол<prn><pers><p3><sg><nom>\$ дыбысты/дыбыс<n><acc>\$ анығырақ/анық<adj><comp>/ анық<adj><comp><advl>/анық<adj><comp><subst><nom>\$ ести/ести<v><tv><prc_impf>\$ бастады/баста<v><tv><ifi><p3><pl>/баста<v><tv><ifi><p3><sg>/баста<vaux><ifi><p3><pl>/ баста<vaux><ifi><p3>
Mor. disambiguation	[^] нава<n><nom>\$ [^] бүген<adv>\$ [^] бик<adv>\$ [^] эйбэт<adj>\$ [^] ,<cm>\$ [^] жылы<adj>\$ [^] гына<postadv>\$ [^] .<sent>\$
Lex. transfer (+ selection)	[^] нава<n><nom>/нава<n><nom>\$ [^] бүген<adv>/бөгөн<adv>\$ [^] бик<adv>/бик<adv>\$ [^] эйбэт<adj>/эйбэт<adj>\$ [^] ,<cm>/,<cm>\$ [^] жылы<adj>/йылы<adj>\$ [^] гына<postadv>/гына<postadv>\$ [^] .<sent>/.<sent>\$
Struct. transfer	[^] нава<n><nom>\$ [^] бөгөн<adv>\$ [^] бик<adv>\$ [^] эйбэт<adj>\$ [^] ,<cm>\$ [^] йылы<adj>\$ [^] гына<postadv>\$ [^] .<sent>\$
Mor. generation	нава бөгөн бик эйбэт, йылы гына.

Table 3: Translation process for the phrase *нава бүген бик эйбэт, жылы гына* ‘The weather today is very nice, it is very warm’.

Corpus	Tokens	Coverage
Tatar New Test.	163,603	72.04%
Tatar Wikipedia	37,123	70.19%
Bashkir Wikipedia	12,267	65.99%

Table 4: Naïve vocabulary coverage over the three corpora.

Corpus	Direction	Tokens	Unknown	WER
story	tt→ba	311	9	8.97%
	ba→tt	312	1	7.72%

Table 5: Word error rate and over the small test corpus.

Forcada, Mikel L., Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O’Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M. Tyers. 2011. *Aperium: a free/open-source platform for rule-based machine translation*. *Machine Translation*, 25(2):127–144.

Gilmullin, R. A. 2008. The Tatar-Turkish Machine Translation Based On The Two-Level Morphological Analyzer. In *Interactive Systems and Technologies : The Problems of Human-Computer Interaction*, pages 179–186, Ulyanovsk.

Hamzaoğlu, Ilker. 1993. Machine translation from Turkish to other Turkic languages and an implementation for the Azeri language. Master’s thesis, Bogazici Uni-

versity.

Karlsson, F., A. Voutilainen, J. Heikkilä, and A. Anttila. 1995. *Constraint Grammar: A language independent system for parsing unrestricted text*. Mouton de Gruyter.

Levenshtein, V. I. 1966. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics—Doklady* 10, 707–710. *Translated from Doklady Akademii Nauk SSSR*, pages 845–848.

Linden, Krister, Miikka Silfverberg, Erik Axelsson, Sam Hardwick, and Tommi Pirinen, 2011. *HFST—Framework for Compiling and Applying Morphologies*, volume Vol. 100 of *Communications in Computer and Information Science*, pages 67–85.

Tantuğ, A. Cüneyd, Eşref Adalı, and Kemal Oflazer. 2006. Computer Analysis of Turkmen Language Morphology. pages 186–193. *Advances in natural language processing, proceedings (LNAI)*.

Tantuğ, A. Cüneyd, Eşref Adalı, and Kemal Oflazer. 2007. A MT system from Turkmen to Turkish employing finite state and Statistical Methods. In *Proceedings of MT Summit XI, Copenhagen, Denmark*.

Tyers, Francis, Jonathan North Washington, Ilmar Salimzyan, and Rustam Batalov. 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 Eight Interna-*

tional Conference on Language Resources and Evaluation (LREC'12), Istanbul, Turkey, may.

Washington, Jonathan, Mirlan Ipasov, and Francis Tyers. 2012. A finite-state morphological transducer for kyrgyz. In Chair), Nicoletta Calzolari (Conference, Khalid Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, Istanbul, Turkey, may. European Language Resources Association (ELRA).

Çöltekin, Çağrı. 2010. A freely available morphological analyzer for Turkish. *Proceedings of the 7th International Conference on Language Resources and Evaluation (LREC2010)*, pages 820–827.