

# TRmorph: A morphological analyzer for Turkish

Çağrı Çöltekin

**Draft:** August 23, 2013

This document describes the new/development version of TRmorph. As such, there may be some mismatches between what is documented here and how the analyzer behaves. This version is a complete overwrite of the previous version reported in Çöltekin 2010. If you are using the older version (you shouldn't), this document is probably useless for you.

## 1 Introduction

TRmorph is an open-source<sup>1</sup> finite-state morphological analyzer for Turkish. This document describes how to use the tools that comes with this package, as well as some implementation details that may be helpful for people who want to customize this open-source tool for their own needs. The complete source of the analyzer and a web-based demo can be accessed through <http://www.let.rug.nl/coltekin/trmorph>.

This document describes the current version of the analyzer. This version is a complete rewrite of the earlier version report in Çöltekin 2010. The earlier version of TRmorph was implemented using SFST (Schmid 2005), the current version is implemented with more popular finite state description languages *lexc* and *xfst* from Xerox (Beesley and Karttunen 2003), using Foma (Hulden 2009) as the main development tool.

The *lexc/xfst* implementation of TRmorph should compile with any *lexc/xfst* compiler without much additional effort. The only foma-specific notation used in the morphology description is about handling simple reduplication, which can also be handled with *twolc* rules, or compile-replace (Beesley and Karttunen 2003).<sup>2</sup>

<sup>1</sup>Current version of TRmorph is licensed under GNU Lesser General Public License. See README file in the TRmorph distribution for more information.

<sup>2</sup>TRmorph can be compiled with HFST (Lindén et al. 2009) without modification since HFST uses foma as the back end for parsing xfst files.

## 2 How to use it

### 2.1 Compilation from the source

To compile TRmorph from the source, you need a *lexc/xfst* compiler such as *foma*, a C preprocessor, GNU *make* and some standard UNIX utilities.

If all requirements are in place, to build analyzer/generator FST, you should type *make* in the main TRmorph distribution directory. The resulting binary automaton file will be *trmorph.fst*.

TRmorph comes with a set of other finite-state tools that are useful in various NLP tasks. Currently, tools for the following tasks are distributed together with TRmorph.

- stemming/lemmatization
- morphological segmentation
- hyphenation
- guessing unknown words

To compile these tools, you should specify the FST you want to build as an argument to *make*, e.g., *make stemmer* will build an binary automaton called *stem.fst*. These additional tools are described in Section 6.

### 2.2 Customizing TRmorph

TRmorph is an open source utility. As a result, you are free to modify the source according to your needs. Source

code includes some useful comments on what/how/where things are done. Furthermore, TRmorph can be customized for some common choices during the compilation. These options are typically related to more relaxed analysis. For example whether to allow non-capitalized proper names, or analyze (and generate) text written in all capitals, or set the decimal and thousand separator in numbers. These options are set in the file `options.h`. The file contains documentation along with the existing options. This feature is under development (as of July 2013), new options are currently being added, and existing options may not fully work as intended yet.

Another common need for customizing a morphological analyzer is to add or modify the lexical entries. The lexicon structure and format of the lexical entries are described in Section 4.

## 2.3 Trying it out

Assuming you have built the binary `trmorph.fst` using foma, you can simply start foma and use xfst commands implemented in foma to analyze and generate the words. Here is an example session:

```
$ foma
foma[0]: regex @"trmorph.fst";
2.1 MB. 62236 states, 135237 arcs, Cyclic.
foma[1]: up okudum
oku<V><past><1s>
foma[1]: down oku<V><past><2s>
okudun
```

The first line is typed at the shell prompt to start foma. The second line reads the FST specified in `trmorph.fst` into the foma environment. The fourth line asks for the analysis of the verb *okudum* ‘I read-PAST’, and the fifth line is the output of the analysis. The sixth line asks for the generation of the analysis string produced earlier, modifying the agreement marker to second person singular agreement.

Note that part of the output is removed for readability. We should also note that this example presents one of the rare cases where the analysis is unambiguous. Turkish morphological analysis is an ambiguous process, and TRmorph does not try to avoid it during the analysis (see

Section 5).<sup>3</sup>

Once you are convinced that the output may be useful for your purposes, you will probably want to use it for analyzing large amount of text. For batch analysis tasks foma’s `flookup` utility is a better fit.

To use the analyzer with the HFST, you need to compile the source automaton with HFST tools.

## 3 The tagset

The description of the morphology in TRmorph mostly follows Göksel and Kerslake 2005. However, there are some divergences, and tags used in TRmorph analyses does not necessarily match with any of the tags used in any grammar book. This section describes the tags used in the current version of TRmorph. The aim of this section is to help users understand the output of the system. Occasional discussion of the morphological process is included, but this section documents neither the morphology of the language nor the way it is implemented in TRmorph. Our focus in this section is to describe the tags one finds in the analysis strings produced by the analyzer (or tags one needs to use for generation). The index at the end of the document also allows easy access to the points where a particular tag is defined or mentioned in this document.

A clarification of the notation for the surface forms is in order before starting the documentation of the tagset and related suffixes. Suffixes in Turkish often contain under-specified vowels and consonants that are resolved according to morphophonological rules, like vowel harmony. These vowels and consonants are indicated with capital letters listed below.

**A** is realized as either ‘a’ or ‘e’.

**I** is realized as either ‘ı’, ‘i’, ‘u’ or ‘ü’.

**D** is realized as either ‘d’ or ‘t’.

**P** is realized as either ‘p’ or ‘b’.

**K** is realized as either ‘k’, ‘ğ’ or ‘y’.

<sup>3</sup>For most purposes, the output of the morphological analyzer needs to be disambiguated. There are quite a few morphological disambiguators for Turkish reported in the literature, but, as yet, there are no disambiguators that work with TRmorph output.

C is realized as either ‘c’ or ‘ç’.

A letter in parentheses indicate a buffer consonant or vowel, that may be dropped in certain contexts.

### 3.1 Part-of-speech tags

All analysis strings include at least one POS tag. A word may get multiple POS tags, and even the same POS tag repeatedly. This process tracks the changes to the syntactic function of a word with suffixation. It follows the idea of *inflectional groups* (IGs, Oflazer, Hakkani-Tür, et al. 1999; Oflazer, Say, et al. 2003) used in the Turkish NLP literature. In a nutshell, the suffixation process may result in changes to the syntactic function of a word, and each of these functions may have their own ‘inflectional’ features. The POS tag of the whole word is the final POS tag in the sequence. However, the other intermediate IGs are also important, and, at least in typical dependency parsing, these IGs may participate in dependency relations with the other words in the sentence. Example (1) demonstrates this process with the analysis of the word *evdekilerinki*.<sup>4</sup>

(1) ev<N><loc><ki><Adj><0><N><pl><gen><ki><Adj>

The example analysis in (1) includes the following IGs.

1. The initial noun ‘ev’ with the locative maker.
2. Addition of the suffix -ki makes an adjective.
3. The adjective becomes a (pro)nominal with a zero derivation, which is inflected for plural and genitive case.
4. Yet another -ki is suffixed, and the word becomes an adjective again.<sup>5</sup>

By default, TRmorph does not mark IG boundaries explicitly. However, one can easily trace the IG changes following the POS tags. All POS tag names start with a capital letter, while other tags always start with a lowercase letter or number. The tag before a new POS tag is always

<sup>4</sup>Rough translation: ‘the *ones* that belong to the *ones* in the house’, as in ‘the *books* that belong to the *people* in the house’.

<sup>5</sup>More likely reading of this example includes another zero derivation causing final POS to be again noun.

Table 1: The list of part of speech tags in TRmorph.

Tag	Description
<Alpha>	Symbols of the alphabet
<Adj>	Adjective
<Adv>	Adverb
<Cnj>	Conjunction
<Det>	Determiner
<Exist>	The words <i>var</i> and <i>yok</i>
<Ij>	Interjection
<N>	Noun
<Not>	The word <i>değil</i>
<Num>	Number
<Onom>	Onomatopoeia
<Postp>	Postposition
<Prn>	Pronoun
<Punc>	Punctuation
<Q>	Question particle <i>mi</i>
<V>	Verb

preceded with the tag of the suffix that cause the derivation. When there is no surface suffix, a zero-derivation tag <0> is inserted before the POS tag.

The example in (1) is also interesting because of the fact the suffix -ki may result indefinitely long words (see Section 3.4).

All part-of speech tags used in TRmorph are listed in Table 1. Most POS tags are self explanatory, and does not require much explanation. The following part of speech tags are somewhat unusual and deserves some explanation.

<Exist> is used for two words *var* ‘existent/present’ and *yok* ‘non-existent/absent’, where the latter is marked as <Exist:neg>, indicating that it is the negative form (see Section 3.2, for the details of this notation). These words behave mostly like nouns in their predicate function (with zero copula), but marking them simply as nouns does blur their function.

<Not> is used for *değil* ‘not’ only. Like *var* and *yok*, *değil* also behaves like nominal predicates. But again, marking it as noun or verb hides the fact that it has a special function.

⟨Q⟩ is used for the question particle -mı. The question particle is written separately from the predicate it modifies. However, the preferred analysis of question particle in TRmorph is together with the predicate. This ensures that it follows the correct form of the predicate it is attached to, and vowel harmony is applied correctly. However, since we do not assume that the input is tokenized with this assumption, this form make sure that the input is analyzed with the cost of precision. The question particle is discussed further in Section ??.

### 3.2 Subcategorization of lexemes

Besides the major major POS tags or word classes discussed above, TRmorph makes use of a set of subcategory tags to mark features that are part of a lexeme. Typically the subcategorization is applied to a root form in the lexicon, but some morphemes and POS tags after a derivation may also receive a subcategory tags. Subcategories defined here are features of a morpheme that do not have a surface realization. Representing these features using a different notation allows one to make this distinction, and the surface-analysis mapping becomes (almost) one-to-one. If a representation where all tags have a uniform notation is desired, the analyzer source can be modified accordingly, or easier, a simple regular expression based converter can be used.

The subcategories generally mark semantic differences, but they may also result in morphosyntactic differences. Lexical subcategorization in TRmorph output is marked using the syntax ⟨Cat:subcat<sub>1</sub>:subcat<sub>2</sub>:...⟩, where ‘Cat’ is a major category and ‘subcat<sub>1</sub>’, ‘subcat<sub>2</sub>’ and so on are sub categories. The order of subcategory tags are not important (although they are produced in a consistent order). A typical example of a subcategory is *proper nouns*, which are tagged as ⟨N:prop⟩.

The following lists subcategories used in TRmorph for all word classes that may be specified together with a subcategory.

**Nouns** Besides the tag ⟨N:prop⟩ marking proper names, abbreviated nouns are marked with the tag ⟨N:abbr⟩. For an abbreviated proper name, the tag is ⟨N:prop:abbr⟩.

**Conjunctions** are subcategorized as *coordinating*, *adverbial* or *subordinating* conjunctions, marked using tags ⟨Cnj:coo⟩, ⟨Cnj:adv⟩, ⟨Cnj:sub⟩ respectively.

The last one of these categories, ⟨Cnj:sub⟩, include only a limited set of conjunctions which come first in a subordinate clause. These words currently are *ki*, *eğer* and *şayet* (all borrowings from Persian). The other subordinating particles/words occur at the end of subordinate clauses, and they are marked as postpositions (⟨Postp⟩) described below. Furthermore, most of the subordination in Turkish is done through suffixation which is described in Section 3.14.

**Pronouns** Pronouns are further categorized as *personal*, *demonstrative* and *locative* pronouns, marked using ⟨Prn:pers⟩, ⟨Prn:dem⟩, ⟨Prn:locp⟩ respectively. Furthermore, the pronouns that form questions, like *kim* ‘who’, and *ne* ‘what’, are marked as ⟨Prn:qst⟩. Subcategory markers for both aspects can be present. For example *kim* ‘who’ would be marked as ⟨Prn:pers:qst⟩.

Besides the above subcategories, personal pronouns get person-number agreement markers. These markers can be useful in subject-predicate agreement as well as in other constructions (such as genitive-possessive construction involving pronouns). However, the agreement in Turkish is far from trivially determined (see Göksel and Kerslake 2005, pp.116–122). The markers ⟨Prn:pers:1s⟩, ⟨Prn:pers:2s⟩, ⟨Prn:pers:3s⟩, ⟨Prn:pers:1p⟩, ⟨Prn:pers:2p⟩ and ⟨Prn:pers:3p⟩ are tags used for the personal pronouns with person-number agreement. The agreement markers are further discussed in Section 3.11.

The reflexive pronoun *kendi(si)* is marked as ⟨Prn:pers:refl⟩.

Subcategorization of pronouns, particularly as personal pronouns, are sometimes not a clear decision. Subcategories of some pronouns are left unspecified even though they are often used as personal pronouns, and some pronoun marked as personal pronouns may refer to entities other than people.

**Determiners** are marked for definiteness. *Definite* determiners are marked ⟨Det:def⟩ and *indefinite* de-

terminers are marked  $\langle \text{Det:indef} \rangle$ . The question words that fill the same syntactic slot as determiners *ne kadar* ‘how much’ and *hangi* ‘which’ are tagged with  $\langle \text{Det:qst} \rangle$ .

Further subcategorization of determiners (for example quantifiers) can be implemented in the future.

**Postpositions** are always subcategorized in two dimensions. First subcategory is the syntactic category (POS) of the resulting postpositional phrase, either an *adjectival* or *adverbial* phrase, marked as  $\langle \text{Postp:adj} \rangle$  and  $\langle \text{Postp:adv} \rangle$  respectively. Note that unlike other POS tags, these subcategory markers start with a lowercase letter.

Postpositions choose their noun phrase complements. Besides the category of the resulting phrase, postpositions also include a tag specifying the requirement for the complement noun phrase. The tag marking required complement type is formed by a concise description of the requirement followed by the capital letter ‘C’. The postpositions that require the complement to be in *ablative*, *dative* and *instrumental* cases are marked  $\langle \text{Postp:ablC} \rangle$ ,  $\langle \text{Postp:datC} \rangle$ , and  $\langle \text{Postp:insC} \rangle$  respectively. Postpositions that require non-case marked complement are tagged  $\langle \text{Postp:nomC} \rangle$ . The postpositions that require the noun phrase to be suffixed with either -lI or -slz are marked with  $\langle \text{Postp:liC} \rangle$ .<sup>6</sup> Finally, postpositions that require numeric expressions as their complements are marked with  $\langle \text{Postp:numC} \rangle$ .

**Numbers** are tagged as  $\langle \text{Num:ara} \rangle$  for Arabic numerals, and  $\langle \text{Num:rom} \rangle$  for Roman numerals. Numbers that are spelled out are not marked with a subcategory marker (but still marked as  $\langle \text{Num} \rangle$ ). Besides numbers, the question word *kaç* ‘how many’ is also tagged as a number with a sub tag specifying that it is a question word, resulting in  $\langle \text{Num:qst} \rangle$ .

**Verbs** are currently not subcategorized in TRmorph.

Subcategorizing verbs as *transitive* and *intransitive*, or marking all types (cases) of noun phrase complements a verb can take is planned and some early steps are underway as of this writing (July 2013).

<sup>6</sup>Like -(y)lA that we mark as  $\langle \text{ins} \rangle$ , these suffixes are typically considered derivational suffixes, however their use resemble case markers.

**Adverbs** are not currently subcategorized, except a few adverbial question words for which the tag  $\langle \text{Adv:qst} \rangle$  is used.

**Exist** The tag  $\langle \text{Exist} \rangle$  exists only for two words *var* ‘existent/present’ and *yok* ‘non-existent/absent’. Since *yok* is the negative of *var*, it is tagged as negative:  $\langle \text{Exist:neg} \rangle$ .

Some verbs, nouns, adjectives, adverbs and conjunctions are formed by more than one written words. Some of these are adjacent words, like the adverb *apar topar* ‘hurriedly’, but some may be split like the conjunction *ya*, as in *ya evdedir ya iş yerinde* ‘s/he is either at home or the office’. Furthermore, some of individual ‘words’ in such constructions cannot be used by themselves, like *topar* above. If the non-split multi-word expressions are input to the analyzer together, they are analyzed like other words of the same class. However, if they are input word-by-word, a sub tag  $\langle \text{:partial} \rangle$  is added to the main POS tag. For example *apar* and *topar* are tagged as  $\langle \text{Adv:partial} \rangle$  and *ya* is tagged as  $\langle \text{Cnj:partial} \rangle$  (more precisely  $\langle \text{Cnj:coo:partial} \rangle$ ). Currently, the tags  $\langle \text{N:partial} \rangle$ ,  $\langle \text{Adj:partial} \rangle$  and  $\langle \text{V:partial} \rangle$  are used for parts of nouns, adjectives and verbs respectively.

### 3.3 Nominal morphology and noun inflections

Nouns, pronouns, adjectives and adverbs in Turkish form the larger class of nominals. Most adverbs and adjectives can function as nouns (or pronouns). For example *mavi* ‘blue’ may mean ‘the blue one’, or *hızlı* ‘fast’ may mean ‘the fast one’. In TRmorph this is handled by allowing any adjective or adverb to become a noun with a *zero derivation*. A zero derivation is always marked with the tag  $\langle 0 \rangle$  followed by the new POS tag, in this case  $\langle \text{N} \rangle$ .

Nouns can be suffixed with the plural suffix, one of the possessive suffixes and one of the case suffixes. All of these inflections are optional. When not marked with any of these suffixes, the default is singular, no possessive marking, an no case marking (or nominal), respectively. When these suffixes co-occur, they have to occur in the order listed, shown in Figure 1. The full list of noun inflections are presented in Table 2.

If there is a plural marker, analysis string after the  $\langle \text{N} \rangle$  will include the tag  $\langle \text{pl} \rangle$ . TRmorph does not mark for



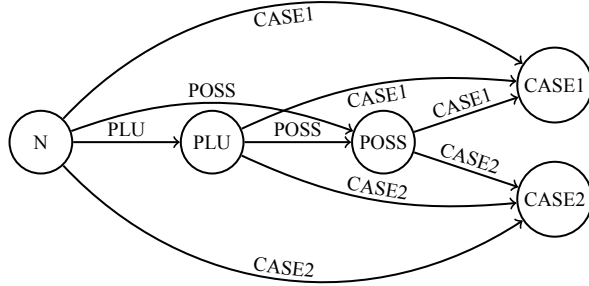


Figure 1: Automata depicting noun inflections. The edge CASE1 represents the locative and ablative suffixes, CASE2 represents all other case-like suffixes. The reason for the differentiation is due to the fact that the state CASE1 can be followed by the suffix -ki.

singular. If a noun is not marked for plural, it is assumed to be singular.

Possessive markers follow either the nominal stem, or the plural marker. Besides marking for possession, these suffixes have a number of other functions, such as deriving pronouns when attached to adjectivals (e.g., determiners or adjectives).

The first five suffixes in Table 2 are commonly recognized cases in Turkish. The *instrumental/commutative* marker also behaves like case suffixes. There are two more suffixes, namely -lı and -sız that can occupy the same slot, which are marked with tags [⟨li⟩](#) and [⟨siz⟩](#) respectively.

The -(s)ı suffix, listed as [⟨p3s⟩](#) in Table 2, is highly ambiguous. One of its many functions that may be confused with the possessive suffix is forming noun compounds. One of types of noun compounds in Turkish is formed by a bare noun and a (head) noun with suffix -(s)ı, for example *at arabası* ‘horse carriage’.<sup>7</sup> TRmorph marks this usage of the suffix -(s)ı as [⟨ncomp⟩](#). In this use, this suffix always causes ambiguities. Besides the fact that a noun suffixed with -(s)ı can either be marked for possession or

<sup>7</sup>Even though one can assume that this use is somewhat related to possession, it is not strictly possessive marking (the horse does not own the carriage). Furthermore, although the surface suffix is skipped after a compound marker, a compound of this form may also be in possessive form (e.g., ‘someone’s horse carriage’).

Table 2: Noun inflections.

	Function	surface	tag
Possessive	Plural	-lAr	<a href="#">⟨p1⟩</a>
	First person singular	-(l)m	<a href="#">⟨p1s⟩</a>
	Second person singular	-(l)n	<a href="#">⟨p2s⟩</a>
	Third person singular	-(s)ı	<a href="#">⟨p3s⟩</a>
	First person plural	-(l)mız	<a href="#">⟨p1p⟩</a>
	Second person plural	-(l)niz	<a href="#">⟨p2p⟩</a>
	Third person plural	-lArı	<a href="#">⟨p3p⟩</a>
Case	Accusative	-(y)ı	<a href="#">⟨acc⟩</a>
	Dative	-(y)A	<a href="#">⟨dat⟩</a>
	Ablative	-DA n	<a href="#">⟨abl⟩</a>
	Locative	-DA	<a href="#">⟨loc⟩</a>
	Genitive	-(n)ın	<a href="#">⟨gen⟩</a>
	Instrumental/commutative	-(y)ıA	<a href="#">⟨ins⟩</a>

as the head of a noun compound, since one of the two -(s)ı suffixes following each other is deleted from the surface form, it can also be both (a noun compound marked for possession).

The case (or case-like) suffixes change the role of the noun (or the noun phrase headed by the noun) in the sentence. For example a locative marked noun phrase may function as an adverb (*saat dokuzda görüşürüz*) or an adjective (*yedi yaşında çocuk*). However, following the common practice in the literature we do not attempt to mark possible POS changes after case-like markers.

### 3.4 The suffix -ki

The suffix -ki, tagged as [⟨ki⟩](#), attaches to locative or ablative marked nouns. The resulting word functions as an adjective or a pronoun. In both cases, TRmorph marks the transition to an adjective. For example, *evdeki* is analyzed as ‘ev(N)⟨loc⟩⟨ki⟩⟨Adj⟩’. Since all adjectives are allowed to become a noun through a zero derivation, the pronoun reading is intended to be represented by this change. For example, the intended analysis for *evdeki kitap* ‘the book in the house’ is ‘ev(N)⟨loc⟩⟨ki⟩⟨Adj⟩’, while analysis for *evdeki uyuyor* ‘the one/person in the house is sleeping’ appends ‘⟨O⟩⟨N⟩’ at the end of the analysis string.

The (pro)noun formed by -ki can further be suffixed with other nominal suffixes. Although the number of iterations using -ki rarely exceed 2 in practice, there is no principled limit. As a result, length of a Turkish word is unbounded in-principle.

### 3.5 Tags related to nominal predicates

Any nominal in Turkish may become a predicate with one of the copular suffixes -(y)DI, -(y)mış, -(y)sA or -(y). These suffixes correspond to *past*, *evidential*, *conditional*, and *present* predicates involving the copula ‘be’. The copular markers has to precede one of the verbal person agreement markers. For example *öğrenciydik* ‘we were students’, *öğrenciymişler* ‘they were [evidentially] students’, *öğrenciysen* ‘if you are/were a student’, *öğrenciyim* ‘i’m a student’. Since the third person singular agreement suffix is null on the surface and the buffer -(y)- does not surface in this case, any nominal without additional copular or person suffixes serve as a nominal predicate with present copula and third person singular agreement. Additionally, since a predicate with third person singular agreement also agrees with a third person plural subject, we additionally mark such a noun as having present copula and third person plural agreement (for example, *babam öğretmen*, *anem ve ablam doktor* ‘my father is a teacher, my mother and older sister are doctors’).

TRmorph handles this process by allowing any noun to first become a verb by a zero derivation, and then marking it with the appropriate copula and the person agreement marker. The tags for copula are *<cpl:pres>*, *<cpl:past>*, *<cpl:evid>* and *<cpl:cond>* for *present*, *past*, *evidential* and *conditional* copula respectively. Last three tags are also possible after a verb with a tense/aspect/modality suffix, and is discussed further in Section 3.12. Example analyses for the examples discussed above would be as follows:

<i>öğrenciydik</i>	<i>&lt;N&gt;&lt;O&gt;&lt;V&gt;&lt;cpl:past&gt;&lt;1p&gt;</i>
<i>öğrenciymişler</i>	<i>&lt;N&gt;&lt;O&gt;&lt;V&gt;&lt;cpl:evid&gt;&lt;3p&gt;</i>
<i>öğrenciysen</i>	<i>&lt;N&gt;&lt;O&gt;&lt;V&gt;&lt;cpl:cond&gt;&lt;2s&gt;</i>
<i>öğrenciyim</i>	<i>&lt;N&gt;&lt;O&gt;&lt;V&gt;&lt;cpl:pres&gt;&lt;1s&gt;</i>
<i>öğretmen</i>	<i>&lt;N&gt;&lt;O&gt;&lt;V&gt;&lt;cpl:pres&gt;&lt;3s&gt;</i>
<i>doktor</i>	<i>&lt;N&gt;&lt;O&gt;&lt;V&gt;&lt;cpl:pres&gt;&lt;3p&gt;</i>

Besides copular suffixes, the suffix -(y)ken (making adverbials from verbs, discussed in Section 3.14) may occupy the same slot as the copular suffixes, although its

use is more restricted.

The nominal predicate with a copula and person agreement may be followed by the marker Göksel and Kerslake 2005 call ‘generalizing modality marker’, the suffix -Dlr. It is particularly common with *<3s>* as it disambiguates between the noun and the predicate reading. The tag for this marker in TRmorph is *<dir>*.

### 3.6 Number inflections

The suffix -(ş)Ar, tagged *<dist>*, attached to numbers form *distributive* numerals. Besides the numbers (written as numerals or spelled out), question word *kaç* ‘how many’ may also get this suffix, and tagged with *<dist>*.

The ordinal numerals are formed using the suffix -(l)ncl, and tagged as *<ord>*. Ordinals are also specified by a ‘dot’ after Arabic or Roman numerals. TRmorph currently does not handle this notation.

Percent sign before a numeral is treated like a prefix, and tagged as *<perc>*.

### 3.7 Verbal voice suffixes

Turkish verbs can be suffixed with one or more of the voice suffixes *reflexive*, *reciprocal*, *causative* and *passive*. The tags used for these functions are *<rfl>*, *<rcp>*, *<caus>* and *<pass>*, respectively. The first two are rather unproductive while causative and passive forms are productive. Furthermore, causative suffix can be used repetitively.<sup>8</sup> With some verbs, use of double causative suffix yields the same semantics as a single causative suffix. TRmorph does not treat these cases separately. If surface string has double causative suffixes, the analysis will include two *<caus>* tags, regardless of its semantics.

Despite the fact that most grammar books list voice suffixes under inflectional morphology, TRmorph treats them as derivations, i.e., a *<V>* tag follows the voice related tags.

### 3.8 Compound verbs

A verbal stem (possibly including voice suffixes) may be followed by a set of suffixes that form compound

<sup>8</sup> Again, although this is limited in practice, there is no principled limit on the number of causative suffixes that one can string one after another.

Table 3: Suffixes that make compound verbs.

Suffix	Tag	Expresses
-(y)Abil	<a href="#">abil</a>	ability
-(y)Iver	<a href="#">iver</a>	immediacy
-(y)Agel	<a href="#">agel</a>	habitual/long term
-(y)Adur	<a href="#">adur</a>	repetition/continuity
-(y)Ayaz	<a href="#">ayaz</a>	almost
-(y)Akal	<a href="#">akal</a>	stop/freeze in action
-(y)Agör	<a href="#">agor</a>	somewhat like <a href="#">iver</a>

verbs. These suffixes are related to some stand-alone verbs. These suffixes are listed in Table 3.

The first three suffixes in this Table 3 are relatively productive, the others are rare or their use are mostly lexicalized. Although not frequent in use, more than one these suffixes may attach to the same stem, for example *çıkıverebilir* ‘he/sh/it may possibly come out/show up’ analyzed as ‘çık([V](#))[iver](#)([abil](#))[aor](#)([3s](#))’.

The form of [abil](#) in a negative verb is -(y)A, and unlike the rest of the suffixes listed in Table 3 it follows the negative marker.

Currently, these suffixes are not marked as derivations (they do not cause an IG change).

### 3.9 Negative marker

Negation of a verbal predicate is indicated with the suffix -mA, and marked simply as [neg](#). For nominal predicates do not get this suffix, instead the particle *değil* is used.

### 3.10 Tense/aspect/modality markers

A verb with a set of suffixes described above either becomes a finite verb by taking one of the tense, aspect and modality (TAM) markers followed by a person-number agreement suffix, or it can be subject to subordination and becomes nominalized.

The list of TAM suffixes, the corresponding tags and brief descriptions are given in Table 4.

Table 4: Tense/aspect/modality markers. The usage of suffix -(y)A to express conditional aspect is informal, and rather restricted. Aorist suffix is highly irregular. The choice of -Ar and -Ir depends on the stem. The -z form occurs only after negative marker, and it is not realized on the surface if it precedes first person agreement suffixes.

Tag	Suffix	Description
<a href="#">evid</a>	-mİş	evidential past (perfective)
<a href="#">fut</a>	-(y)AcAk	future
<a href="#">obl</a>	-mAlI	obligative
<a href="#">impf</a>	-mAktA	imperfective
<a href="#">cont</a>	-(I)yOr	imperfective
<a href="#">past</a>	-DI	past (perfective)
<a href="#">cond</a>	-sA,-(y)A	conditional
<a href="#">opt</a>	-(y)A	optative
<a href="#">imp</a>	-	imperative
<a href="#">aor</a>	-Ar,-Ir,-z,-	aorist

### 3.11 Person and number agreement

After TAM markers a finite verb requires one of the person and number agreement markers. The surface form of the person-number agreement markers change depending on the suffixes they follow. Table 5 lists the person agreement markers and their surface form according the TAM of the verb they attach to.

### 3.12 Copular markers and -DIr

The copular suffixes discussed in Section 3.5 can also be attached to a verb after a TAM marker, typically forming complex tenses. These suffixes are -(y)DI, -(y)mİş and -(y)sA, tagged as [cpl:past](#), [cpl:evid](#) and [cpl:cond](#), respectively.

The conditional copula -(y)sA can co-occur with other copular markers. When there is a copular suffix, person-number agreement suffixes normally attach after the first copula. However the third person plural suffix may be after the TAM marker or second copular suffix as well.

Similar to the nominal predicates with a copula, copular suffixes may be followed the ‘generalizing modality marker’ -DIr tagged as [dir](#).



Table 5: Verbal person agreement markers. The first character of the person agreement tags is a number indicating the person (1<sup>st</sup>, 2<sup>nd</sup> or 3<sup>rd</sup>), and second one indicates the number (singular or plural). The suffixes listed in the column marked ‘TAM1’ follow the TAM markers *<evid>*, *<fut>*, *<obl>*, *<impf>* and *<cont>* as well as the evidential copula *<cpl:evid>* and nominal predicates. The same set of suffixes also follow positive verbs with *<aor>* without a negative marker. The suffixes on the column marked ‘TAM2’ are used after *<past>* and *<cond>* as well as the corresponding copular markers *<cpl:past>* and *<cpl:cond>*.

Tag	TAM1	TAM2	optative	imperative
<i>&lt;1s&gt;</i>	-(y)lm	-m	-(y)lm	*
<i>&lt;2s&gt;</i>	-sln	-n	-sln	-
<i>&lt;3s&gt;</i>	-	-	-	-sln
<i>&lt;1p&gt;</i>	-(y)lz	-K	-llm	*
<i>&lt;2p&gt;</i>	-slnlz	-nlz	-slnlz	-(y)ln,-(y)lnlz
<i>&lt;3p&gt;</i>	-lAr	-lAr	-lAr	-slnlAr,-

### 3.13 The question particle

Question particle -ml, tagged as *<q>*, is normally written separately. However, it has an intimate relationship between the verb or the nominal predicate it attaches to. First, a few exceptions aside, it is attached to a tensed verb without a person agreement. In this case, the person agreement and the suffixes that may follow must be attached to the question particle. In this particular case, the verb will often be analyzed wrongly as having the agreement marker *<3s>* or *<3p>*, since a predicate with null person agreement suffix may agree with third person singular or plural subjects. Second, the question particle follows the vowel harmony rules, and the underspecified vowel on -ml is realized based on the last vowel of the verb. As a result the question particle can only be analyzed (and generated) with precision only together with the word it is attached to.

If tokenized together with the predicate, TRmorph will swallow the space in between the predicate and the -ml and analyze it altogether. In this case the lowercase tag *<q>* is used. Furthermore, it is a common spelling mistakes to

write the question particle together with the related word. TRmorph can be instructed to accept this common mistake during compile time, in which case the tag will again be *<q>*.

### 3.14 Subordination

A set of suffixes attached to an ‘untensed’ verb, a verb without any TAM markers, result in the phrase headed by the verb to become a subordinate clause. TRmorph follows the description in Göksel and Kerslake 2005, and makes the distinction between three different forms of subordination. First, a set of suffixes produce *verbal nouns* from a non-finite verb. The resulting words function as the head of the noun phrases, and with some limitation they can receive all nominal inflections. The second group forms *participles*, which form relative clauses. Participles can also take nominal inflections with few restrictions. The last group, *converbs*, form adverbials and they are more restricted in terms of the morphemes attached to them. The suffixes that form different types of subordinating suffixes overlap significantly. As a result, producing ambiguous analyses.

TRmorph uses the tag structure *<type:subtype>* for marking subordinating suffixes. The first part, *type*, is one of *vn*, *part* and *cv* for verbal nouns, participles and converbs, respectively. The second, *subtype*, part indicate a further distinction of the function of the suffix, a relevant linguistic abbreviation, but sometimes a version of the surface form of the suffix. The tags used for all three types of subordinating suffixes are listed in Table 6.

Since verbal nouns, participles and converbs derive nominal, adjectival and adverbial phrases, respectively, POS tags, *<N>*, *<Adj>* and *<Adv>*, follow these tags.

Some of the suffixes have multiple functions and may derive more than one type of subordinate clauses. Furthermore, TRmorph will produce some spurious ambiguity because of the fact that any adjective, hence a word suffixed with an participle, is allowed to become a noun with a zero derivation.

The list in Table 6 follows Göksel and Kerslake 2005. The main exception is the suffixes listed by Göksel and Kerslake 2005 as converbial suffixes that require a postposition. Since the postposition in these cases will signal the adverbial function of postpositional phrase, TRmorph

Table 6: Subordinating suffixes and tags used for subordinating suffixes.

Tag	Suffix
<a href="#">⟨vn:inf⟩</a>	-mA
<a href="#">⟨vn:inf⟩</a>	-mAK
<a href="#">⟨vn:is⟩</a>	-(y)İş
<a href="#">⟨vn:past⟩</a>	-Dİk
<a href="#">⟨vn:fut⟩</a>	-(y)AcAk
<a href="#">⟨vn:res⟩</a>	-(y)An
<a href="#">⟨part:past⟩</a>	-Dİk
<a href="#">⟨part:fut⟩</a>	-(y)AcAk
<a href="#">⟨part:pres⟩</a>	-(y)An
<a href="#">⟨cv:ip⟩</a>	-(y)İp
<a href="#">⟨cv:meksizin⟩</a>	-mAksİzİn
<a href="#">⟨cv:ince⟩</a>	-(y)İncA
<a href="#">⟨cv:erek⟩</a>	-(y)ArAk
<a href="#">⟨cv:eli⟩</a>	-(y)Alİ
<a href="#">⟨cv:dikce⟩</a>	-DİkCA
<a href="#">⟨cv:esiye⟩</a>	-(y)AsİyA
<a href="#">⟨cv:den⟩</a>	-dAn
<a href="#">⟨cv:den⟩</a>	-zdAn
<a href="#">⟨cv:cesine⟩</a>	-CAsİnA
<a href="#">⟨cv:ya⟩</a>	-(y)A
<a href="#">⟨cv:ken⟩</a>	-(y)ken

does not mark the complement of the postposition as a converb.

Most of these suffixes attach to an untensed verb. Except, the suffix -(y)ken which behaves much like the copular suffixes discussed above. Furthermore, the -(y)A in its subordinating function is typically used together with reduplication, e.g., *koşa koşa* ‘run-(y)A run-(y)A = hurriedly’, but also occurs in words like *diye*, where it does not need reduplication.<sup>9</sup>

Besides the subordinating suffixes (participles) discussed above, some of the TAM markers, namely [⟨aor⟩](#), [⟨evid⟩](#), [⟨fut⟩](#), [⟨impf⟩](#) and less commonly [⟨cont⟩](#), derive adjectivals resembling relative clauses. TRmorph handles this by analyzing any verb with one of these TAM markers

<sup>9</sup>One may also analyze *diye* as a postposition, as it’s use as subordinator is semantically unlike the others uses of -(y)A.

without further suffixes (e.g., agreement markers) as an adjective. For example, the word *görölmüş* in *görölmüş mektup* ‘see-PASV-EVID letter = the letter that was seen’ is analyzed as ‘gör(V)⟨pass⟩(V)⟨evid⟩(Adj)’.

### 3.15 Productive derivational morphemes

Almost all the tags and relevant morphological process above are described as part of inflectional morphology in most grammar books. The suffixes described here are the ones that are traditionally considered derivational suffixes. Some of these suffixes, for example -İİ and -sİz discussed earlier, may attach to word forms that are already inflected by other suffixes. Others normally attach only to the stem and produce another stem.

Of these suffixes, the noun–verb derivation suffix -İA causes a large number of ambiguous analyses since it is part of many other suffixes. These, for example, include the plural suffix -İAr whose remainder -r also matches a verbal suffix (aorist). Hence, including -İA in the analysis causes an increase in the analyses of any plural noun. Currently, TRmorph analyzes -İA only after onomatopoeia. The rest of the verbs derived from nouns using this suffix are lexically specified.

TRmorph does not limit the number of derivational suffixes that can be stringed one after another other, even though multiple derivations of this sort is a lot more restricted.

Besides the sources of possible erroneous over-analyses listed above, the derivational morphology specification in TRmorph over-generates in some cases. In particular, any form of the diminutive suffix is allowed to attach to any noun, although most nouns are used only one of the diminutive suffixes. The ambiguity and overgeneration are discussed in Section 5.

## 4 The lexicon

TRmorph contains a root lexicon which is created extracting root forms from a large web corpora, and checking the possible forms against online dictionaries, and the lexicon of the earlier version which was based on Zemberek A. A. Akin and M. D. Akin 2007. The result is also checked and corrected manually as part of the development process.

Table 7: Derivational morphemes analyzed by TRmorph. The column ‘Derivation’ lists the POS changes using a two letter symbols. The first letter is the original POS, and the second one is the POS after the suffixation. Here, N, J, A, M, V and O stand for noun, adjective, adverb, number, verb and onomatopoeia, respectively.

Tag	Suffix	Derivation
<li>	-ll	NA NJ
<siz>	-slz	NA NJ
<lik>	-llk	NN JN AN
<dim>	-Clk	NN
	-cAk	
	-(l)cAk	
	-cAğlz	
<ci>	-Cl	NN NJ
<ca>	-CA	NA AA JJ MJ
<yici>	-(y)lcl	VJ
<cil>	-Cll	NJ
<gil>	-gil	NN
<lan>	-lAn	JV
<las>	-lAş	NV JV
<vis>	-ylş	VN
<esi>	-(y)Asl	VJ
<sal>	-sAl	VJ
<la>	-lA	NV OV

The lexicon files are located under the directory `lexicon` and included (through C preprocessor) as a single root lexicon. The files under `lexicon/` are simply a list of root forms and their *continuation classes*. Continuation classes can be any `LEXICON` declaration in the file `morph.lexc`, but typical continuation classes are the main word (POS) categories, such as N, Adj and V. The lexical exceptions are specified after the main category information. For example, `V_AorAr` for verbal roots that take the exceptional `-Ar` form of the aorist suffix. Likewise, `N_comp` is used for lexicalized nominal `-sl` compounds since when these words are pluralized the plural marker is inserted between the word and the suffix `-sl`.

The lexical forms are similar to the written forms of the relevant stem. However, a set of special ‘multi-character’ symbols are used for providing information necessary for

morphophonological processing. A large group of these symbols are concerned with ‘final stop devoicing’ (or voicing depending on your view point). The consonants `ç`, `t`, `k`, `p` and `g` at the end of some of the roots are replaced with their voiced counterparts if they precede a suffix that starts with a vowel. These root forms are lexically marked by replacing the consonants above with multi-character symbols `^c`, `^t`, `^k`, `^p` and `^g`, respectively.

Besides the voicing changes of consonants, some borrowings end with a ‘palatalized’ consonant that affects vowel harmony process. For example *saat* ‘watch/clock’ is inflected as *saat-i* ‘watch-ACC’ instead of *saat-ı* as vowel harmony suggests. These words are indicated by the vowel before such a consonant by a three-letter multi-character symbol. These symbols always start with `^p` and a capitalized version of the relevant vowel. For example, the word *saat* is listed as *sa^pAt* in the lexicon.

One last class of similar special symbols are so-called *silent vowels and consonants*. These are particularly useful for abbreviations and numerals, but also some names of foreign origin. The suffixes that follow such words are also subject to morpho-phonological process like vowel harmony. However, this cannot be derived from their written form. For example correct inflected form of *ABD-DAT* ‘USA-DAT’ is *ABD^ye*, not *ABD^ya*. The way to solve this problem is to insert a silent (front-unrounded) vowel after the abbreviated form. The multi-character symbols `^sBUV` `^sBRV` `^sFUV` `^sFRV` `^sVC` and `^sUC` are used for silent vowels and consonants (see the comments in file `lexicon/abbreviation` for more information).

A somewhat inconsistent notation is used for three morphological processes. First, the multi character symbol `@DEL@` is inserted before a vowel that is deleted if a suffix starting with a vowel follows. Second, the last consonant in some borrowings are duplicated if they follow a suffix that start with a vowel. These root forms are marked by inserting the multi-character symbol `@DUP@` before the duplicated consonant. And the last symbol `@DELS@` is used in lexical entries of a few borrowed words which delete `s` in the suffix `-sl`.<sup>10</sup>

<sup>10</sup>These multi-character symbols are both inconsistent with the others, and they may be confused with ‘flag diacritics’ at first sight (TRmorph does not use any flag diacritics). This notation in the lexicon may change in the future version of TRmorph.

## 5 Ambiguity and overgeneration

This section discusses the ambiguous analyses in TRmorph, and also touches upon a related but different problem, overgeneration.

The morphological analysis of Turkish text is an inherently ambiguous process. However, the design choices made in a morphological analyzer affects the number of ambiguous analyses produced per word. TRmorph, by design, does not try to reduce the number of ambiguous analyses. In general, TRmorph produces more ambiguous analyses than the others (mainly based on Oflazer 1994) reported in the literature.

The following is a list of cases where one finds ambiguous morphological analyses in TRmorph. Some of these cases are not specific to TRmorph, and for example, noted by Oflazer and Tür 1997 as well. This list may be useful for the users who may wish to disambiguate the output of the analysis using rule-based methods, or it may also be useful in the process of designing statistical disambiguators.

1. Ambiguous root forms, for example *yüz* can be analyzed as:

- (a) yüz<N> ‘face’
- (b) yüz<Num> ‘hundred’
- (c) yüz<V><imp><2s> ‘swim’

2. A root form is the same as a shorter root and one or more suffixes, for example *buna* can be analyzed as

- (a) bu<Prn:dem><dat> ‘this-DAT’
- (b) buna<V><imp><2s> ‘become senile-IMP’
- (c) bun<N><dat> ‘trouble-DAT’

Note that the root ‘bun’ is a very rare/regional word, and the imperative verb reading is also very unlikely. However the best option for the analyzer is to produce all these analyses, and let the later stages disambiguate between them.

3. The surface form of a suffix is a combination of two other suffixes. For example, the word *evleri* can be

- (a) ev-leri ‘ev<N><p3p> = their house’
- (b) ev-ler-i ‘ev<N><p1><acc> = houses-ACC’

Furthermore, the same word can also be analyzed as

- (a) ‘ev<N><p1><p3s>’
- (b) ‘ev<N><p1><p3p>’
- (c) ‘ev<N><ncomp><p3p>’
- (d) ‘ev<N><ncomp><p1>’
- (e) ‘ev<N><ncomp><p1><p3p>’
- (f) ‘ev<N><ncomp><p1><p3s>’
- (g) ‘ev<N><ncomp><p1><p3p>’

The reason for these analyses has to do with the sources of ambiguity explained in items 6 and 7.

4. An analysis with multiple morphemes is also a (derived) lexicalized form. For example the word *konuşma* can be analyzed as

- (a) konuşma<N> ‘speech’
- (b) konuş<V><vn:inf><N> infinitive ‘to speak’, e.g., as in *konuşmamızı isemiyorlar* ‘The do not want us to speak’
- (c) konuş<V><neg><imp><2s> ‘speak-NEG-IMP = don’t talk’

5. different affixes surfacing the same way, *evin* can be

- (a) ev-(n)ln ‘ev<N><gen> = of the house’
- (b) ev-(l)n ‘ev<N><p2s> = your house’

6. The same surface suffix has multiple functions. For example, the word *doktorlar* can be,

- (a) doktor<N><p1> ‘doctors’
- (b) doktor<N><0><V><cpl:pres><3p> ‘they are doctors’

7. Some suffixes are not realized on the surface in the neighborhood of some other suffixes. These are generally, but not always, the suffixes having the same or similar surface forms. For example, *evleri* (the example in item 3) may be analyzed as

- (a) ev<N><p3p> as in *Annem ve babamın evleri İstanbul’da* ‘My parents’ house is in İstanbul’
- (b) ev<N><p1><p3p> as in *Annem ve babamın bütün evleri deniz manzaralı* ‘All houses of my parents have a sea view’.

since in case of  $\langle p1 \rangle$  (-lAr) and  $\langle p3p \rangle$  (-lArI) are combined, the plural suffix -lAr does not realized on the surface.<sup>11</sup>

This particular source causes an extremely large number of ambiguous analyses because the multi functional suffix -(s)I is omitted in case it precedes (or follows) another -(s)I, but also a -lArI, -lI, -lIk, -slz, -Cl or -Clk. Since some of these suffixes may follow each other, and -(s)I itself has multiple functions, a word like *bağım-sız-lık-çı-lığ-ı-nı* causes a combinatorial expansion of ambiguous analyses because of the fact that at every suffix boundary marked with a dash in the example there may be a -(s)I suffix being deleted. This is further amplified by the fact that -(s)I may express  $\langle ncomp \rangle$  or  $\langle p3s \rangle$  and any of the resulting words may also have a null suffix expressing third person singular or plural agreement on a nominal predicate.<sup>12</sup> Most of these analyses will be semantically not plausible. However, there is no clear way of ruling them out at the analysis stage. The following illustrates the problem with a more tangible example, using the word *arabasız* which can be analyzed as one of the following (and more).

- (c) araba $\langle N \rangle$  $\langle siz \rangle$  $\langle Adj \rangle$  ‘without a car’
- (d) araba $\langle N \rangle$  $\langle p3s \rangle$  $\langle siz \rangle$  $\langle Adv \rangle$  ‘without a his/her car’
- (e) araba $\langle N \rangle$  $\langle ncomp \rangle$  $\langle siz \rangle$  $\langle Adv \rangle$ , e.g., in *at arabasız* ‘without a horse carriage’
- (f) araba $\langle N \rangle$  $\langle ncomp \rangle$  $\langle p3s \rangle$  $\langle siz \rangle$  $\langle Adv \rangle$ , e.g., in *at arabasız* ‘without his/her horse carriage’

Besides the ambiguity described above, overgeneration is another problem that one faces when the FST is used for generating surface forms. Unlike analysis, generation is almost always deterministic in Turkish. Nevertheless, there are a few cases where TRmorph produces multiple surface strings for a single analysis string. The following provides a (likely incomplete) list of cases where TRmorph is expected to overgenerate, i.e., either produce

multiple (correct) surface strings for the same input, or produce incorrect surface strings in generation mode.

1. One of the clear cases where overgeneration occurs is the diminutive,  $\langle dim \rangle$ . The diminutive suffix in Turkish is one of -Clk, -cAk, -(l)cAk, -cAğlız. TRmorph allows attaching any of these suffixes to any noun. This is unlikely to cause problems during the analysis. However, it will certainly produce incorrect surface forms.
2. The  $\langle p3s \rangle$  suffix -(s)I may also be used for marking third person plural possessive ( $\langle p3p \rangle$ ). For example *ev-i* in *Ali ve Ayşe’nin evi* ‘The house of Ali and Ayşe’ should be tagged as  $\langle p3p \rangle$ . On the other hand, the suffix -lArI is also used to express  $\langle p3p \rangle$ . As a result any analysis string with the symbol  $\langle p3p \rangle$  will generate both surface options.
3. A similar case of overgeneration is with the null agreement suffix which should generally be tagged as  $\langle 3s \rangle$ . However, such a predicate may also agree with a  $\langle 3p \rangle$  subject. Consequently, a null-agreement suffix on a predicate is tagged as both  $\langle 3s \rangle$  and  $\langle 3p \rangle$ . Since  $\langle 3p \rangle$  can also be expressed with the suffix -lAr, a analysis string with  $\langle 3p \rangle$  also generates multiple surface forms.
4. Another known case of overgeneration is related to the relaxed analysis of alternative spellings or common misspellings. In the simplest case, every word will be generated once capitalized and once all lower-case. If ‘all capitals’ option is enabled, another surface form which is in all capital letters will be produced.
5. Similarly, if the analyzer is instructed to accept the proper noun suffixes without an apostrophe, in the generation mode the surface form with and without apostrophe will be included. As a result, some of the options may need to be tuned if the FST is to be used for generation.
6. Some symbols, like apostrophe have multiple representations in Unicode definition. As a result any word that require an apostrophe will result in surface form for each alternative symbol.

<sup>11</sup>One can also explain this as  $\langle p3p \rangle$  being realized as -l in this particular context.

<sup>12</sup>Most straightforward reading of the word is dative form of the noun phrase can roughly be translated as ‘his/her state of being a supporter of independence’. With this root, The total number of analyses is 25560.



7. After a small set of borrowings like *cami* ‘mosque’, the ‘s’ in the suffix -(s)l is deleted according to official rules. However, this seems to be out of fashion in current use, and use of ‘s’ (even in text) is more common than its deletion. Since TRmorph accepts both surface strings, this will cause generating multiple strings.

There are also a few other cases where some (sizeable number of) speakers diverge from the canonical forms. An example is the redundant use of genitive suffix after a pronoun, before the suffix -(y)lA, e.g., the surface form of ‘sen(Prn:pers:1s)(ins)’ should be *sen-in-le* where the suffix -in is redundant. Some speakers tend not to use -in in such constructions. TRmorph accepts both use, hence the generation will be ambiguous.

8. Some borrowed words include a few vowels with circumflex, namely *â*, *û* and *î*. Except for a few words where use of circumflex help disambiguation between different words, these vowels have been replaced by their non-circumflexed version in modern use. TRmorph allows this replacement even if the lexical form of the word should include a circumflex.<sup>13</sup> This also results in overgeneration, since any analysis string with a circumflexed vowel will have a surface form with and without circumflex.

## 6 Other tools

### 6.1 Stemming and lemmatization

In morphologically complex languages like Turkish, proper stemming requires analyzing the given word and stripping off the analysis symbols such that only the stem remains.

Although one can do this easily by filtering analyzer output, TRmorph includes a simple wrapper automaton for convenience. The automaton is defined in the file `stemmer.fst`. You need to type `make stemmer` to produce the binary `stem.fst`. This binary file can be used

<sup>13</sup>One can also allow circumflexed vowels to be used for their non-circumflexed counterparts in the lexicon. This is useful if one needs to analyze somewhat older text. Enabling this option will also cause overgeneration.

the way analyzer is used. Given a surface word, this automaton will produce the lexical form as the analysis string.

Optionally, one can keep the first tag, which is the syntactic category of the stem. Note that stemmer takes the lexical form as the ‘stem’, even if the lexical form has derivational suffixes immediately following the root form. Another compile time option related to stemmer causes the verbs to be suffixed with correct form of infinitive marker -mAk. This form of the verbs are what the dictionaries use as head words. Both options can be set in the file `options.h`.

Note that ambiguity is less of a problem for the stemmer. However, in examples like *buna* discussed on page 12, there will be multiple stem forms produced (*bu*, *bun* and *buna* in this case).

### 6.2 Unknown word guesser

TRmorph includes a rudimentary guesser for guessing unknown words. To produce the automaton for this function, you should type `make guesser`, which would produce the file `guess.fst`. The usage of the automaton is again similar to the others. The surface strings of the FST is the (unknown) words, while analysis level is either the full analysis strings with possibly unknown root words that may lead to the surface form, or only the root word and its part of speech tag.

The guesser uses the same machinery as the analyzer, except the lexicon is replaced with a FSA that accepts a somewhat restricted set of strings as potential words. Since unknown words will likely include affixes, one may have a better chance of determining the root form of the word, and in most cases the class of the root word.

Depending on its application, the guesser be restricted further according to features of the words that can be coded into a finite state lexicon. For example, one may check whether the words fit into the syllable structure of the language, but this may miss the words of foreign origin, which are likely candidates for being unknown words. Currently only general restriction the guesser include the minimum and maximum root-word length that can be set in the file `options.h`.

The guesser may also be adjusted to return full analysis string(s) or only the root form followed by the POS tag. Again, these options can be set in `options.h`.

Other customizations can be achieved by adjusting the file `guesser.lexc`.

The guesser is a standalone FST, to use it in combination with the analyzer, two automata can be combined with *priority union* such that guesser is only invoked if the analyzer fails. This can be achieved either as a simple wrapper xfst file, or if you are using foma's `flookup` utility specifying both FST files on the command line like `flookup -a trmorph.fst guesser.fst`.

### 6.3 Morphological segmentation

Morphological segmentation is the task of finding morpheme boundaries on the surface strings. TRmorph distribution includes an automaton description for segmenting the words into their morphemes. To build the segmenter you need to type `make segmenter` and the resulting binary will be called `segment.fst`.

TRmorph marks the root and morpheme boundaries on the surface string to aid morpho-phonological rules. These boundaries are deleted from the surface string in the normal analyzer FST. The segmentation FST relies on this and the following trick for segmenting a given word to its surface morphemes: The given input string is first analyzed with the regular analyzer FST. Then the analysis strings are passed to a slightly modified FST in generation mode, which does not delete the boundary markers from the surface string.

It should be noted that the surface morpheme boundaries are not always determined uniquely. It is especially difficult to decide whether some buffer vowels or consonants belong to the morpheme preceding or following them. TRmorph consistently attaches these buffer letters to the morpheme that follow the boundary.

Because of the way it is implemented currently, the segmenter output needs to be post processed to obtain the desired result. The segmenter will produce multiple identical segmented strings, and there will also be some incorrect segmentations due to overgeneration discussed in Section 5. The output should be post-processed to remove multiple identical segmentations. The incorrect segmentations due to overgeneration can be eliminated by comparing the segmented string with the original one. An example post processing script is provided as `scripts/segment-filter.py`.

### 6.4 Hyphenation and syllabification

Hyphens in Turkish are inserted at the syllable boundaries. Because of the regular syllable structure and transparency of the orthography, this process does not require any dictionary lookup, or morphological analysis. Since the hyphenation problem is easy to solve with a FST, a standalone FST defined in xfst language included in the TRmorph distribution.

To build the hyphenation FST you need to type `make hyphenate` and the resulting binary will be called `hyphenate.fst`.

The surface string of the FST is Turkish words (or strings resembling words) and analysis string is the words where a hyphen '-' is inserted between the syllables, or at the points where one can insert a hyphen.

## References

- Akın, Ahmet Afşin and Mehmet Dündar Akın (2007). "Zemberek, an open source NLP framework for Turkic Languages". Available at <http://zemberek.googlecode.com/>. URL: <http://zemberek.googlecode.com/>.
- Beesley, Kenneth R and Lauri Karttunen (2003). "Finite-state morphology: Xerox tools and techniques". In: *CSLI, Stanford*.
- Çöltekin, Çağrı (2010). "A Freely Available Morphological Analyzer for Turkish". In: *Proceedings of the 7th International Conference on Language Resources and Evaluation (LREC2010)*. Valetta, Malta.
- Göksel, Aslı and Celia Kerslake (2005). *Turkish: A Comprehensive Grammar*. London: Routledge.
- Hulden, Mans (2009). "Foma: a finite-state compiler and library". In: *Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics: Demonstrations Session*. Association for Computational Linguistics, pp. 29–32.
- Lindén, Krister, Miikka Silfverberg, and Tommi Pirinen (2009). "HFST Tools for Morphology—An Efficient Open-Source Package for Construction of Morphological Analyzers". In: *State of the Art in Computational Morphology*. Ed. by Cerstin Mahlow and Michael Pitrowski. Communications in Computer and Informa-

- tionScience. Springer, pp. 28–47. ISBN: 978-3-642-04130-3.
- Oflazer, Kemal (1994). “Two-level description of Turkish morphology”. In: *Literary and Linguistic Computing* 9 (2).
- Oflazer, Kemal, Dilek Hakkani-Tür, and Gökhan Tür (1999). “Design for a Turkish Treebank”. In: *Proceedings of the Workshop on Linguistically Interpreted Corpora, at EACL’99*. Bergen, Norway.
- Oflazer, Kemal, Bilge Say, Dilek Zeynep Hakkani-Tür, and Gökhan Tür (2003). “Building a Turkish treebank”. In: *Text, Speech and Language Technology*. Springer, pp. 261–277. ISBN: 9781402013348.
- Oflazer, Kemal and Gökhan Tür (1997). “Morphological Disambiguation by Voting Constraints”. In: *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics*, pp. 222–229.
- Schmid, Helmut (2005). “A programming language for finite state transducers”. In: *Proceedings of the 5th International Workshop on Finite State Methods in Natural Language Processing (FSMNLP 2005)*. Helsinki, pp. 308–309.

## Index

- ⟨1p⟩, 7, 9  
⟨1s⟩, 7, 9  
⟨2p⟩, 9  
⟨2s⟩, 7, 9, 12  
⟨3p⟩, 7, 9, 12, 13  
⟨3s⟩, 7–9, 13  
  
⟨0⟩, 3, 5–7, 12  
  
⟨abil⟩, 8  
⟨abl⟩, 6  
⟨acc⟩, 6, 12  
⟨Adj⟩, 3, 3, 6, 9, 10, 13  
⟨Adj:partial⟩, 5  
⟨adur⟩, 8  
⟨Adv⟩, 3, 9, 13  
⟨Adv:partial⟩, 5  
⟨Adv:qst⟩, 5  
⟨agel⟩, 8  
⟨agor⟩, 8  
⟨akal⟩, 8  
⟨Alpha⟩, 3  
⟨aor⟩, 8–10  
⟨ayaz⟩, 8  
  
⟨ca⟩, 11  
⟨caus⟩, 7  
⟨ci⟩, 11  
⟨cil⟩, 11  
⟨Cnj⟩, 3  
⟨Cnj:adv⟩, 4  
⟨Cnj:coo⟩, 4  
⟨Cnj:coo:partial⟩, 5  
⟨Cnj:partial⟩, 5  
⟨Cnj:sub⟩, 4  
⟨cond⟩, 8, 9  
⟨cont⟩, 8–10  
⟨cpl:cond⟩, 7–9  
⟨cpl:evid⟩, 7–9  
⟨cpl:past⟩, 7–9  
⟨cpl:pres⟩, 7, 12  
⟨cv:cesine⟩, 10  
⟨cv:den⟩, 10  
⟨cv:dikce⟩, 10  
⟨cv:eli⟩, 10  
⟨cv:erek⟩, 10  
⟨cv:esiye⟩, 10  
⟨cv:ince⟩, 10  
⟨cv:ip⟩, 10  
⟨cv:ken⟩, 10  
⟨cv:meksizin⟩, 10  
⟨cv:ya⟩, 10  
  
⟨dat⟩, 6, 12  
⟨Det⟩, 3  
⟨Det:def⟩, 4  
⟨Det:indef⟩, 5  
⟨Det:qst⟩, 5  
⟨dim⟩, 11, 13  
⟨dir⟩, 7, 8  
⟨dist⟩, 7  
  
⟨esi⟩, 11  
⟨evid⟩, 8–10  
⟨Exist⟩, 3, 3, 5  
⟨Exist:neg⟩, 3, 5  
  
⟨fut⟩, 8–10  
  
⟨gen⟩, 3, 6, 12  
⟨gil⟩, 11  
  
⟨Ij⟩, 3  
⟨imp⟩, 8, 12  
⟨impf⟩, 8–10  
⟨ins⟩, 5, 6, 14  
⟨iver⟩, 8  
  
⟨ki⟩, 3, 6  
  
⟨la⟩, 11  
⟨lan⟩, 11  
⟨las⟩, 11  
⟨li⟩, 6, 11  
⟨lik⟩, 11

⟨loc⟩, 3, 6  
 ⟨N⟩, 3, 3, 5–7, 9, 12, 13  
 ⟨N:abbr⟩, 4  
 ⟨N:partial⟩, 5  
 ⟨N:prop⟩, 4  
 ⟨N:prop:abbr⟩, 4  
 ⟨ncomp⟩, 6, 12, 13  
 ⟨neg⟩, 8, 12  
 ⟨Not⟩, 3, 3  
 ⟨Num⟩, 3, 5, 12  
 ⟨Num:ara⟩, 5  
 ⟨Num:qst⟩, 5  
 ⟨Num:rom⟩, 5  
 ⟨obl⟩, 8, 9  
 ⟨Onom⟩, 3  
 ⟨opt⟩, 8  
 ⟨ord⟩, 7  
 ⟨p1p⟩, 6  
 ⟨p1s⟩, 6  
 ⟨p2p⟩, 6  
 ⟨p2s⟩, 6, 12  
 ⟨p3p⟩, 6, 12, 13  
 ⟨p3s⟩, 6, 12, 13  
 ⟨part:fut⟩, 10  
 ⟨part:past⟩, 10  
 ⟨part:pres⟩, 10  
 ⟨pass⟩, 7, 10  
 ⟨past⟩, 8, 9  
 ⟨perc⟩, 7  
 ⟨p1⟩, 3, 5, 6, 12, 13  
 ⟨Postp⟩, 3, 4  
 ⟨Postp:ablC⟩, 5  
 ⟨Postp:adj⟩, 5  
 ⟨Postp:adv⟩, 5  
 ⟨Postp:datC⟩, 5  
 ⟨Postp:insC⟩, 5  
 ⟨Postp:liC⟩, 5  
 ⟨Postp:nomC⟩, 5  
 ⟨Postp:numC⟩, 5  
 ⟨Prn⟩, 3  
 ⟨Prn:dem⟩, 4, 12  
 ⟨Prn:locp⟩, 4  
 ⟨Prn:pers⟩, 4  
 ⟨Prn:pers:1p⟩, 4  
 ⟨Prn:pers:1s⟩, 4, 14  
 ⟨Prn:pers:2p⟩, 4  
 ⟨Prn:pers:2s⟩, 4  
 ⟨Prn:pers:3p⟩, 4  
 ⟨Prn:pers:3s⟩, 4  
 ⟨Prn:pers:qst⟩, 4  
 ⟨Prn:pers:refl⟩, 4  
 ⟨Prn:qst⟩, 4  
 ⟨Punc⟩, 3  
 ⟨Q⟩, 3, 4, 9  
 ⟨q⟩, 9  
 ⟨rcp⟩, 7  
 ⟨rfl⟩, 7  
 ⟨sal⟩, 11  
 ⟨siz⟩, 6, 11, 13  
 ⟨V⟩, 3, 7, 8, 10, 12  
 ⟨V:partial⟩, 5  
 ⟨vn:fut⟩, 10  
 ⟨vn:inf⟩, 10, 12  
 ⟨vn:past⟩, 10  
 ⟨vn:res⟩, 10  
 ⟨vn:yis⟩, 10  
 ⟨yici⟩, 11  
 ⟨yis⟩, 11