

Part of speech tagging of Turkish

Ümit Mersinli* – Mustafa Aksan*

1. Introduction

Beginning with Hankamer (1989) and Köksal (1975), studies on computer aided processing of Turkish has emerged following multiple approaches. Turkish, with respect to its complex morphotactics, the in-root phonological alternations forced by harmony rules and other phonological constraints, and the number of homographs, keeps its challenging position in the Natural Language Processing (NLP) literature. After mentioning Oflazer et al. (1994) and Çiçekli (1997) as the applications of 1990s, Akın (2007) and Çöltekin (2010) are two current, rule-based, accessible implementations on the morphological analysis and annotation of Turkish.

This paper presents a root-driven, non-stochastic, graph-based approach to the Parts of Speech Tagging of Turkish. The implemented graphs representing the cascaded finite-state transducer used in modeling Turkish morphotactics is accessible through www.tudd.org and are adaptable to other formalisms. The NLP-dictionaries of the module are formed in synchrony with Turkish National Corpus Project¹ and thus, represent the lexicon of present-day Turkish. The architecture of dictionaries and pre-defined lexical features may also contribute to the research on forming the standard NLP-dictionary of Turkish.

1.1. Definitions

In this paper, the term POS Tagging is used in a broad sense that will cover the annotation of inflectional and derivational affixes as well as lexical categories of the base forms. In this respect, the term is reserved for morpheme tagging of Turkish as a preliminary step for morphosyntactic or grammatical tagging as in Treebanks or semantic tagging as in WordNets.

Considering the highly agglutinative nature of Turkish and the amount of homographs, the distinction between inflectional and derivational affixations in Turkish is a challenging issue as also stated in Sezer (2001). In this paper, this distinction is done rather with a computational point of view and thus will not be discussed in detail on a theoretical basis.

* Mersin University.

¹ Turkish National Corpus is funded by Scientific and Technological Research Council of Türkiye (TÜBİTAK). (Grant no: 108K242)

1.2. Data

Data of the study are derived from the ongoing Turkish National Corpus (TNC) Project held at Mersin University, Turkey. TNC, as a balanced and representative corpus, is not specifically restricted to any particular subject field, genre or register. Since it contains samples of both written and spoken language, the lexicon and graphs formed for the module represents present-day Turkish.

1.3. Software

Software used for annotation is NooJ as documented in Silberztein (2003). NooJ includes tools for corpus building and management, linguistic analysis, annotation and concordancing. In NooJ “the descriptions of natural languages are formalized as electronic dictionaries, as grammars represented by organized sets of graphs” (Silberztein 2003).

Following the NooJ formalism, Turkish module is also comprised of two basic components; dictionaries and graphs.

2. Dictionaries

To compile a NooJ dictionary (.nod) file; a Properties Definition file (.def), an Inflectional/Derivational Rule file (.nof) and pre-compiled Raw Dictionaries (.dic) are required (Silberztein 2003). The contents of the three mentioned file formats will be illustrated in the following sections of the study.

2.1. Tokenization

Data for tokenization are extracted from a subcorpus including over 100 texts representing different genres taken from TNC. The subcorpus included over 3,300,000 words forms and over 280,000 tokens when proper nouns, abbreviations and acronyms are excluded.

Below is a sample tokenization with NooJ. Proper Nouns, Acronyms and Abbreviations are filtered out manually from the list of word forms.

Figure 1. Tokenization in NooJ

Freq	Tokens
92920	bir
57959	ve
31966	da
31750	de
29874	bu
17037	için
16643	gibi
14124	Bu
13202	çok
12338	o
11750	daha
11137	ne
10443	sonra
9845	olarak
9786	kadar
9010	her
8974	ile
8153	ki
8112	Bir

2.2. Lemmatization

Stemming or affix stripping algorithms and their implementations that can be used for lemmatization in Turkish is out the scope of this study. In this respect, the filtered word forms taken in the previous step are lemmatized manually and an affix database including the affix combinations of Turkish is created. After the pre-tagging process for Lexical Categories, the base form of raw dictionaries is formed as in (1);

- | | | |
|-----|---------|--------------------|
| (1) | al, VB | <i>(take)</i> |
| | al, AJ | <i>(red)</i> |
| | yüz, NB | <i>(a hundred)</i> |
| | yüz, VB | <i>(swim)</i> |
| | yüz, NN | <i>(face)</i> |

Parts of Speech Tags for Lexical Categories are listed in Table 1.

Table 1. Parts of speech tagset for Turkish.

TAG	POS	EXAMPLE
<VB>	Verb	<i>git, gel, dur, bak, kal, sus, gör, dök</i>
<NN>	Noun	<i>gece, hava, renk, fark, dost, oyun</i>
<PN>	Pronoun	<i>bu, kendi, hepsi, herkes, kim, öteki</i>
<NB>	Number	<i>iki, üç, beş, sekiz</i>
<AJ>	Adjective	<i>mavi, yeni, düz, dürüst, zeki</i>
<AV>	Adverb	<i>acaba, asla, bazen</i>
<PP>	Postposition	<i>gibi, göre, için, kadar, karşı, rağmen</i>
<ITJ>	Interjection	<i>aferin, sağol, haydi, hoşçakal, lütfen</i>
<CJ>	Conjunction	<i>ama, çünkü, meğer, üstelik</i>
<ON>	Onomatopoeia	<i>takır, vızıl, gürül</i>
<NP>	Proper Noun	<i>Atatürk, Mersin, Ümit</i>
<AB>	Abbreviation, Acronym	<i>TBMM, TDK</i>
<MI>	Affirmative particle	<i>mi, mı, mu, mü</i>

2.3. Phonemic Alternations

Considering the in-root phonemic alternations as in (2) and (3), phonological rules are defined textually in an Inflectional/Derivational Rule file (phonology_TR.nof) and thus prefixed with “FLX=” as in “FLX=soften_t”.

- (2) akıl → aklında
mind → mind:GEN+LOC
- (3) tıp → tıbbın
medicine → medicine:POSS

The in-root phonemic alternations are listed in Table 3 using the operators in Table 2.

Table 2. Rule operators for NooJ Inflectional/Derivational Grammars.

	delete last character / backspace	<L>	go left
<B2>	delete last two characters	<R>	go right
<D>	duplicate last character	+	OR

Table 3.

tag	rule	example
double	<D>	<i>af > affi, zam > zamma</i>
drop	<L><R>	<i>akıl > aklını, fikir > fikrimin</i>
dropsoften	<B2>b	<i>kayıp > kaybına, kutup > kutbuna</i>
compound1		<i>anaokulu > anaokulları</i>
compound2	<B2>	<i>elyazısı > elyazuları, başağrısı > başağrıları</i>
compound3	<B2>ç	<i>ipucu > ipuçları</i>
compound4	<B2>k	<i>ayçiçeği > ayçiçekleri</i>
soften_ch	c	<i>ağaç > ağacı, süreç > süreci</i>
soften_k	ğ	<i>emek > emeği, diyalog > diyalogu</i>
soften_p	b	<i>kitap > kitabı, mektup > mektubu</i>
soften_t	d	<i>cilt > cilde, dört > dördünü</i>
soften_t_er	d + de	<i>et > eder, git > gider</i>
soften_t_ar	d + da	<i>tat > tadar</i>
softendouble	b<D>	<i>tıp > tıbbın, muhip > muhibbi</i>
change_an	<B2>an	<i>ben > bana, sen > sana</i>
add_er	e	<i>üz > üzer</i>
add_ar	a	<i>yap > yapar</i>

2.4. Lexical Features

The Raw Dictionary is compiled with the rules declared in the Inflectional/Derivational Rule file that includes predefined features in the Properties Definition file.

In the current release of the Turkish module, instead of defining lexical properties that can have multiple features as in (4), we preferred to use a binary format for lexical features as in (5).

(4) al,VB+PHON=end_l

(5) al,VB+end_l

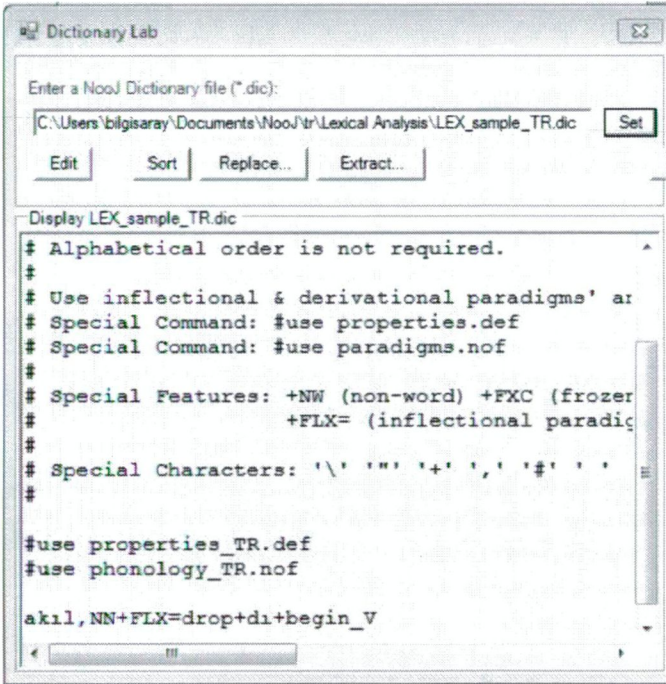
In both of the entries – for use in passivization constraints – it is stated that the Verb “al” ends with consonant “L”, whereas in (5) we can also add more phonological features such as “begin_V” (begins with a vowel) for use in duplications like “ev mev”. (6) is a sample entry of the final version of raw dictionaries.

(6) akı1,NN+FLX=drop+begin_V

2.5. Compilation

The Raw Dictionaries including a declaration of the related Properties Definition (.def) and Inflectional/Derivational Rule file (.nof) are compiled through the menus Lab → Dictionary as in Figure 1.

Figure 1. Dictionary compilation pane in NooJ.



3. Graphs

After NooJ dictionary files (.nod) are compiled, morphotactics of Turkish is modeled with a Morphological Grammar file (.nom) graphically. NooJ graphs let the user design cascaded finite-state transducers through a graphical interface.

3.1. Overall architecture

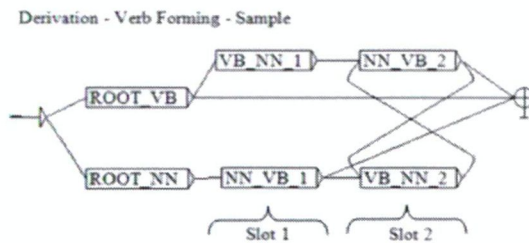
Turkish morphological graph is designed to include both derivational and inflectional affixes since, in most cases, the distinction is problematic due to homophonous affixes such as [-mA] serving in both derivational and inflectional processes as in (7) to (12).

- (7) saçma (*ridiculous*)
→ adjective forming derivational affix
- (8) soruşturma (*investigation*)
→ noun forming derivational affix
- (9) gitme (*don't go OR going*)
→ negative OR gerundive
- (10) dövme (*must forge OR with a tattoo*)
→ part of affix "mAlI" OR noun forming derivational affix
- (11) yapmadan (*without doing OR from doing*)
→ part of adverbial affix "mAdAn" or gerundive
- (12) gidememe (*to my being not able to go*)
→ negative + gerundive + possessive + dative

3.2. Derivation

The derivational subgraph presented in Figure 2 is organized in two slots covering the recursive affixations as in "yaptırttırdı" (s/he **caused** someone to **make** some other one to **get** someone to do it).

Figure 2. Derivational subgraph

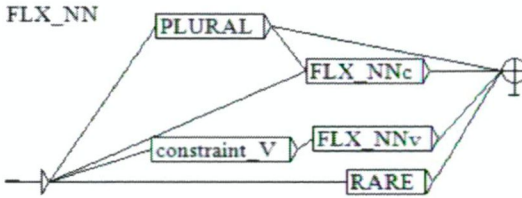


NooJ graphs can be organized in a cascaded manner, so the labels "NN_VB_1" and "NN_VB_2" include the nodes including affixes that derive verbs from nouns.

3.3. Inflection

The inflectional subgraph includes the nominal and verbal inflectional paradigms of Turkish and thus is organized as in two subgraphs. Below is the nominal inflection graph.

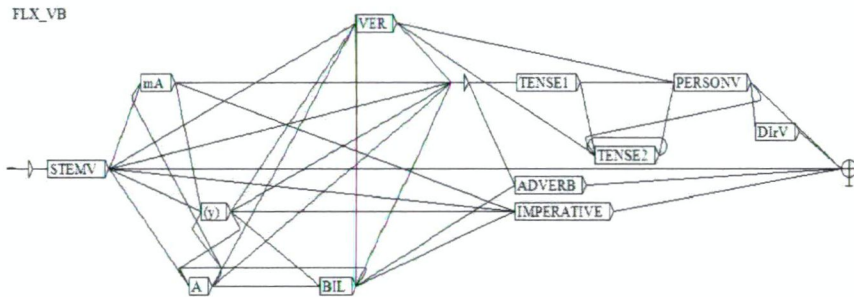
Figure 3. Nominal inflection subgraph



In Figure 3, graphs “FLX_NNc” and “FLX_NNv” are for base forms ending with a consonant and base forms ending with a vowel. This distinction prevents the parser produce artificial ambiguities caused by mostly buffer phonemes.

Figure 4 presents the graph for verbal paradigm.

Figure 4. Finite-state transducer graph for verbal inflection in Turkish



4. Implementation

In the following sections, the performance of the implemented module will be demonstrated.

4.1. Annotation

Below are sample annotations in (10).

- (13) *ölümsüzleştirtirilemeyebileceklerimizdenmişsinizcesine*
 öl,VB+(I)m_NN+sIz_AJ+lAş_VB+Dir_VB+t_VB+Dir_VB+(I)l_VB
 +A+mA+(y)+A+bil+AcAk_NN+lAr+I+mIz[Poss]+DAn[ABL]
 +mIş[Per]+sInIz[2Ppl]+cAsInA_AV
- (14) *okulunki buradaki kitaba benzemiyor*
 okul,NN+I+n+ki[PN]
 bura,NN+DA[LOC]+ki[AJ]
 kitap,NN+A[DAT]
 benze,VB+mA+yor

- (15) *yenilik*
 yen,VB+(I)l_VB+Ik_AJ
 ye,VB+(I)n_VB+(I)l_VB+Ik_AJ
 yeni,AJ+Ik_NN

4.2. Concordancing

Concordancing is done through the "Locate" menu in NooJ. NooJ regular expressions are stated between "<" and ">" symbols to indicate that searching will take place in the text annotation structure.

Below are sample concordance lines for the given search patterns.

<oku,VB>

<p>bir teorik yapının sunduğu yöntemlerle bağlantılan kurdurabiliyor muyuz "Raskolnikov"u durduran, bir romanın kısaltılmış versiyonunu okumayı kendine hakaret olarak algılayan " zaman çarşıda olur yakından işitirdik. bir destan. Eve gidip iyice iyice okumak ve anlamak isterdim. Sony Reader PRS-500, elektronik kitap seyirlik unsura dönüştürme biçiminde uygulandığını zaman çarşıda olur yakından işitirdik. bir destan. Eve gidip iyice iyice okumak ve anlamak isterdim. bir teorik yapının sunduğu yöntemlerle bağlantılan kurdurabiliyor muyuz "Raskolnikov"u durduran, bir romanın kısaltılmış versiyonunu okumayı kendine hakaret olarak algılayan " elbette. Hele bir parça Osmanlıca mi bileceksin? Bakanlıkta şu kadar para eder de bir rahmet uzak dumalı. Televizyon seyretmeyim, gazete ya canım! O kadar biyoloji seyirlik unsura dönüştürme biçiminde uygulandığını belirtmişim. "Bir de bu kitabını kez de aynı abartısıyla karşılaştım. algılama zorluğu çekenler için. Üniversite verilerine ulaşamıyor. Ürün. saniyede 24 MB</p>	<p>okuyor okurken okumayı okur Okuma okumak Okuduğumda okumanızı okuduklarımızdan Okuma okumak Okuduğumda Okuduğumda okuyor okurken okumayı okur okuyabilen okumuş okur okumayım okumuşum okuduklarımızdan oku Okumaya okumanın okuyabiliyor</p>	<p>. Kullanım değerinin yerine çoktan ? Kurdumak zorunda mıyız? "Karanlığın Yüreği kendine hakaret olarak algılayan "okur "un klasikten ne umduğunu, ne yazması olmayan biriydi. Nedenini bilemem ve anlamak isterdim. Okuduğumda bir bir türlü pazar yerindeki o . RSS haberlerini takip etmenizi ve öğreniyoruz. Foucault'un "Hapishanenin Doğuşu yazması olmayan biriydi. Nedenini bilemem ve anlamak isterdim. Okuduğumda bir bir türlü pazar yerindeki o . Kullanım değerinin yerine çoktan ? Kurdumak zorunda mıyız? "Karanlığın Yüreği kendine hakaret olarak algılayan "okur "un klasikten ne umduğunu, ne varsa tanıdık çevrede, kesin söylüyorum kişi geceler boyu çalışmış didinmiş bana diye. Ama ben satamam . İyi de maaş zammına itiraz . kendi laboratuvarım davar. Olmadı mutfakta öğreniyoruz. Foucault'un "Hapishanenin Doğuşu !" dercesine Grass'in kitabını verdiler ara verdim. Gerek çinilmiş gözlemler maliyeti nedir? 22 yaşında okul bitecek ve 10 MB yazabiliyor. 512MB, iGB</p>
---	--	---

Figure 5. Sample concordance line 1

<VB+r[Aor]> <VB+mA+z[Aor]>

<p>"kültürsüz kumazlığın" eline Dışarı çıkınca yola adınızı "kültürsüz kumazlığın" eline 1. Devrimin ağırlıklı bir anlayış adamları insanların bir yanlış 2. Bir ekonomik kriz ortamına 3. İnceleme heyeti'nden onay 4. Merkezleri baskılanır ve bebek 5. Adar kaçınılması ve bebeğin 6. İbi, tek eşliler çiftleşme sona 7. Gezegene ilgili bir çevredir. 8. Bir çevreyle karşı karşıyayız. 9. Spılan fizyonun keşfedildiğini 10. Hastanın gözlerinde bu isteği 11. Duğuna göre, promosyon da 12. J da kansız Klytemnestra eve 13. Kı siteye bulaşmış durumda.</p>	<p>geçer geçmez atar atmaz geçer geçmez ister istemez bulur bulmaz girer girmez alır almaz doğar doğmaz doğar doğmaz erer emez Doğar doğmaz Doğar doğmaz duyar duymaz sezer sezmez ister istemez döner dönmez Bulaşır bulaşmaz</p>	<p>bir şekilde "toplumsa cenk bağlıyor. Çocuk bir şekilde "toplumsa beden yapının dengi gülmediklerine, ancak toplumdaki "günah k donör araştırmalarına ağlayamaz. Bu da ol annesinin sütü ile be yeni bir eşin peşine c böyle bir çevreyle ka , belki bütün organiz , bu elementin zincir r onlardan önce davrz hekimlere yöneliyor l banyoda şişleyerek i sistemi çökerten ve l</p>
--	--	---

Figure 6. Sample concordance line 2.

<NN+A[DAT]> doğru

<p>2. E bitkiler de genellikle aynı 3. ışık ışık, demet demet bu 4. danır: "Bereketli kıldığımız 5. tünden kola ve oradan da 6. dan başlar. Ayak sırtından 7. Enerjiyi dalak ve pankreas 8. bölgesi boyunca kann ve 9. rın gelip beni kurtarmasını, 10. lduğumuzu varsayalım. Bir 11. sol alt köşeden sağ üst 12. rfini koyarsak, hecelemeyi 13. ji bacaklarına geçirirler ve 14. üzerindeki çizgilerden orta 15. sin. Sonra dört köşeyi orta 16. yeri, sürekli kılan ve daha 17. ak istenildiğinde, bizleri bu 18. Ülke kalkınması aşağıdan 19. almak için koridora çıktım, 20. andım durdum koşuğlarda. 21. ede kalacaktır. İşimiz bitti, 22. :edir. Rezervlerin tükenme 23. ardından A noktasından B 24. izyinde, A noktasından C 25. sektir. b. Beklentiler aynı</p>	<p>tarafa doğru yöne doğru yere doğru omuza doğru bileğe doğru mendiyenine doğru göğse doğru bana doğru noktaya doğru köşeye doğru gene doğru geriye doğru noktaya doğru noktaya doğru fazlasına doğru fazlasına doğru yöne doğru yukarıya doğru odama doğru Öğlene doğru akşama doğru noktasına doğru noktasına doğru yöne doğru</p>	<p>samaşık yapırlar. Evet m aktığını görüyoruz; ama m , Süleyman'ın emriyle yürü uzanır. Köprücük kemiğin uzanır. Ayağın ve bacağı aktır. Mide ve bağırsak t ilerler. Buradan yine orta t yolları göstermesini ümit e ilerliyoruz. Ö noktaya ula gittiğim düşünelim. Sol alt yapabiliriz. Yukardaki uyd ilerleyerek lastiğin gergin l 3 cm kalana kadar kesin. birbirlerinin üzerine gelece yönlendiren başlangıç ba yönlendiren başlangıç ba sınırlayan bir çekim kuv olur. Yirmi birinci yüzyılın c koşmaya başladım. İşte c Hoca beni yine kapiya, a ambulansı andıran bir ara azalma eğilimine girdiğini i seyretmesi gereken döviz bir kırılma yaşanacaktır. birbirleriyle tutarlı iseler ke</p>
---	--	---

Figure 7. Sample concordance line 3.

5. Conclusion

As presented and demonstrated in this paper, the design of a corpus-driven non-stochastic annotation module for Turkish showed that this preliminary step in Turkish NLP still has unsolved problems and therefore needs further applications. Further areas of NLP such as

information extraction, morphosyntactic annotation or semantic annotation require a full-coverage standard NLP dictionary of Turkish and an accompanying transducer as well.

References

- Akın, M. D. & Akın, A. A. 2007. Türk dilleri için açık kaynaklı doğal dil işleme kütüphanesi: ZEMBEREK. *Elektrik Mühendisliği* 431, 38.
- Çiçekli, İ. & Temizsoy, M. 1997. Automatic creation of a morphological processor in logic programming environment. In: Drogemuller, R. (ed.) *Proceedings of the 5th International Conference on the Practical Application of Prolog (PAP'97)*. 22nd-24th April 1997. London, UK. 95–106.
- Çöltekin, Ç. 2010. A Freely Available Morphological Analyzer for Turkish. *Proceedings of the 7th International Conference on Language Resources and Evaluation (LREC2010)*. Valletta, Malta, May 2010.
- Hankamer, J. 1989. Morphological parsing and the lexicon. In: Marslen-Wilson, W. (ed.) *Lexical representation and process* Cambridge, MA: MIT Press. 392–408.
- Köksal, A. 1975. *A first approach to a computerized model for the automatic morphological analysis of Turkish*. Ph.D. dissertation. Hacettepe University, Ankara.
- Oflazer, K. & Göçmen, E. & Bozşahin, C. 1994. *An Outline of Turkish Morphology*. Technical Report. Middle East Technical University, Ankara.
- Sezer, E. 2001. Finite inflection in Turkish. In: Taylan, E. E. (ed.) *The Verb in Turkish*. Amsterdam: Benjamins. 1–47.
- Silberztein, M. 2003. *Nooj Manual*. January 10, 2010, from <http://www.nooj4nlp.net>