

Converting Ottoman-script to Latin-script is also studied by Güngör et al. [5] and Korkut [6]. Both studies preferred detecting the root word by means of a large Ottoman lexicon and converting the suffixes using Ottoman-Latin suffix pairs instead of using a two-level morphological analyzer.

III. METHODS

A. Lexicon Collection

We made use of numerous Ottoman-Turkish dictionaries from online sources [7, 8] in order to collect Ottoman transcriptions of Turkish words and build a lexicon. Hayrat Nesriyat’s dictionary holds nearly 110,000 different words and their translations to Ottoman script, as seen in Figure 2. Instead of using Tevakkü’s conjugated verb list, we made use of root verbs and nouns from [7].

Kelimeid	Latince	Osmanlica
10453	ördümek	اوردورمك
10462	öpülmek	اوپولمك
10463	örtbas etmek	اورت باس ايتك
10468	ötürmek	اوتورمك
10470	ötümsüzleşmek	اوتوسمز ليشمك
10474	ötümlüleşmek	اوتوملوشمك
10480	ötmek	اوتتيرمك
10483	ötmek	اوتك
10492	ötüşmek	اوتوشمك
10493	öykülemek	اويكولمك
10499	övünmek	اوگونمك
10507	övlmek	اوكلنك
10509	övmek	اوگنك
10520	örülmek	اورولمك
10522	örüklemek	اوروكلمك
10530	örtünmek	اورتونمك
10535	örtülmek	اورتولمك

Figure 2. A sample from Hayrat Nesriyat’s Ottoman dictionary [7]

B. Morphological Analysis Method

Morphological analysis of Ottoman text may be performed with two methods:

1. Detecting the possible root words by using a dictionary and trying to extract remaining suffixes by comparing syllables to a suffix list.
2. Building a finite state machine with a lexicon composed of verbs, nouns and suffixes for two-level morphological recognition and generation.

We did not prefer the first method because it does not recognize the morphotactics of Turkish language in a structured way. Also, it does not allow morphological generation, that is picking a root word and defining some suffixes and synthesizing a surface level word.

We followed the second approach and therefore we used *foma* [10], a well-known and common tool for finite-state machine compiling as it is mostly suitable for our requirements. We inserted root verbs and nouns to a lexicon file and described inflection rules by defining a set of verb and

noun suffixes. A sample rule set of noun inflections can be seen in Figure 3.

For describing verb and noun inflection rules in *foma* environment, we are inspired by Coltekin’s study

```

LEXICON NPluralSingular
+N:0 NPossessive;
+N+Pl:لر NPossessive;

LEXICON NPossessive
#;
+P1S:م NCase; !(I)m
+P2S:ن NCase; !(I)n
+P3S:ى NCase; !(s)I
+P1P:من NCase; !(I)mIz
+P2P:ن NCase; !(I)nIz
+P3P:لر NCase; !lArI

LEXICON NCase
#;
+Acc:ى #; !(y)I
+Dat:ه #; !(y)A
+Abl:دن #; !DAn
+Loc:ده #; !DA
+Gen:ن #; !(n)In
+Inst:له #; !(y)lA

```

Figure 3. A sample from Ottoman noun inflection rules in *foma* environment

C. Transformations on suffixes

In Turkish, inflection of a word is not always straightforward. Same suffix may take different forms when being appended to different words as shown in Table I. Therefore we needed to define these kinds of suffix transformations in *foma* environment by means of functions.

For example, “+P3S” the possessive 3rd singular person suffix (-I) is preceded with the letter “s” if the root word’s last letter is a vowel. (kitab-**ı**, araba-**sı**). Same transformation is also applied in Ottoman script by prepending the letter “س” to the suffix “-ى”.

Another case is the necessity of inserting a vowel (-I) before possessive suffixes (-m, -n, -mIz, -nIz) when the root word ends with a consonant. This happens in Ottoman script by prepending the letter “-ى” before +P1S, +P1P, +P2P possessive suffixes. We describe these transformations in the “*PossessiveKaynastirmaInsertion*” function in *foma* as described below:

```

define Vowels [ و | ي | ا | ا | ا | ا | ا ];
define Consonants [ ح | ج | ج | ث | ت | پ | ب |
ا | خ | د | ذ | ر | ز | ز | س | ش | ص | ض | ط |
ن | م | ل | ل | ك | ك | ق | ق | ف | غ | ع | ط |
];

```

```
! (I)m, (I)n, (I)mIz, (I)nIz, (s)I insertion
define PossesiveKaynastirmaInsertion
[.] -> ي || NonVowels "^" _ [ ك ز م ا ن ] , ,
[.] -> س || Vowels "^" _ [ ي ] ;
```

After $-(s)I$ and $-lArI$, the suffixes of $-DA$ and $-DAn$ become $-nDA$ and $-nDAn$, respectively, and $-(y)A$ and $-(y)I$ become $-nA$ and $-nI$, respectively. These insertions are performed in Ottoman script by inserting the letter “-” as described in the “*NInsertion*” function below:

```
define NInsertion
[.] -> [ ا | د | ه | ن ] _ "^" [ ي | ل | ر ] "^" || ن
[.] -> ن || Vowels "^" _ [ ك ; ! (n)In
```

If the root word ends with a vowel, the letter “y” (“ي” in Ottoman) is inserted in suffixes such as: $-(y)I$, $-(y)A$, $-(y)IA$. This rule is described in “*YInsertion*” function below:

```
define YInsertion
[.] -> ي || Vowels "^" _ [ ا | ه | ل | ي ] ;
```

IV. EXPERIMENTS

At this stage of our project, we only implemented and experimented inflections of noun words in Ottoman. In order to clearly demonstrate consonant and vowel insertion scenarios of noun inflections we selected two example words:

1. چانطه - Çanta (ends with vowel)
2. كتاب - Kitab (ends with consonant)

All remaining results could be found in Appendix.

A. Surface to Lexical

We performed surface to lexical experiments in order to evaluate the accuracy of FSM. Since inflectional forms of nouns might sometimes cause ambiguities, our FSM should be able to detect possible ambiguities and list all probable parsings according to the rules of Turkish. Figure 4 depicts such a scenario with our two sample words:

```
done:
2.1 kB, 40 states, 54 arcs, 431 paths.
redefined Lexicon: 2.1 kB, 40 states, 54 arcs, 431 paths.
redefined Grammar: 3.3 kB, 58 states, 93 arcs, 431 paths.
3.3 kB, 58 states, 93 arcs, 431 paths.
foma[3]: up
apply up> چانطه+N+P3P+Loc
چانطه+N+P3P+Loc
چانطه+N+Pl+P3S+Loc
apply up> كتاب+N+P3P+AbI
كتاب+N+P3P+AbI
كتاب+N+Pl+P3S+AbI
```

Figure 4. Listing noun inflectional ambiguities using FSM in foma environment

چانطه‌لرینده (Çantalarında) might express two different meanings:

1. Çanta+N+P3P+Loc (at their bag)
2. Çanta+N+Pl+P3S+Loc (at his/her bags)

B. Lexical to Surface Level

Another experiment is giving the lexical description of an inflection and checking the surface level result. As seen in Figure 5, the inflection of word “كتاب - Kitab” is given accurately in two scenarios:

1. كتاب+N+Pl+P2P+Gen: Plural noun, Possessive 2nd plural person and genitive inflection is translated to surface level as كتابلریڭیز (Kitablarınızın - ... of your books) accurately.
2. كتاب+N+P1P+Inst: Singular noun, Possessive 1st singular person, instrumental inflection is translated to surface level as کتابیمیزله (Kitabımızla = with our book) accurately.

```
apply down> كتاب+N+Pl+P2P+Gen
كتابلریڭیز
apply down> كتاب+N+P1P+Inst
کتابیمیزله
```

Figure 5. Morphological generation with word “كتاب - Kitab”

V. EVALUATION

We listed a brief part of our experiments in the previous section. However there were also some wrong inflections that we observed. For example, چانطه+N+Pl+P3P (Plural noun and Possessive 3rd plural person) is translated as چانطه‌لری (Çantaları) although it should be translated as (Çantaları). This happens because the ruleset in our FSM description lacks this exceptional case. These kinds of errors should be observed and corrected in the remaining part of our project.

VI. CONCLUSION

Although some exceptional cases and some other suffixes are obsolete in our current FSM description, we performed the major parts of noun inflections in Ottoman. We demonstrated some example results of two-level morphological recognition and generation (bottom-up and top-down translation). After this midway report we will focus on verb inflections which is a more difficult procedure because of the complicated inflectional nature of verbs in Turkish.

REFERENCES

- [1] Yalniz, I. Z., Altingovde, I. S., Gdkbay, U., & Ulusoy, . (2009). Ottoman archives explorer. *Journal on Computing and Cultural Heritage*, 2(3), 1–20. doi:10.1145/1658346.1658348.
- [2] Ethem F. Can, Pınar Duygulu, A line-based representation for matching words in historical manuscripts, *Pattern Recognition Letters*, Volume 32, Issue 8, 2011, Pages 1126–1138, ISSN 0167-8655, <https://doi.org/10.1016/j.patrec.2011.02.013>.
- [3] Ahmet Afşın Akın, Mehmet Dndar Akın. Zemberek, an open source NLP framework for Turkic languages. <https://github.com/ahmetaa/zemberek-nlp> Accessed on 03.12.2021.
- [4] Çağrı ltekin. A freely available morphological analyzer for Turkish. In *Proceedings of the 7th International Conference on Language Resources and Evaluation (LREC 2010)*, pages 820–827, 2010.
- [5] Ayşegl Gngr and İ. Emre Şahin. Dervaze: A spelling dictionary for digital translation of Ottoman documents. *International Journal of Languages' Education and Teaching*, 5(3):78–84, 2017. doi: 10.18298/ijlet.1824.
- [6] Joomy Korkut. Morphology and Lexicon-based Machine Translation of Ottoman Turkish to Modern Turkish. Draft, May 2019. <https://www.cs.princeton.edu/~ckorkut/papers/ottoman.pdf> Accessed on 03.12.2021.
- [7] Hayrat Neşriyat. “Osmanlıca İmla Kılavuzu” Android application. <https://play.google.com/store/apps/details?id=com.hayrat.imlakilavuzu> Accessed on 03.12.2021
- [8] Tevakku. “Osmanlıca Szlk” Android application. <https://play.google.com/store/apps/details?id=com.tevakku.osttr> Accessed on 03.12.2021
- [9] Oflazer, Kemal. (1993). Two-level Description of Turkish Morphology. *Literary and Linguistic Computing*. 9. 10.1093/lc/9.2.137.
- [10] Mans Hulden. 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 (EACL '09)*. Association for Computational Linguistics, USA, 29–32.